

ACL2010

Uppsala, Sweden
July 11–16, 2010

The 48th Annual Meeting of the Association for Computational Linguistics

Conference Proceedings

Companion Volume



UPPSALA
UNIVERSITET



ACL 2010

**48th Annual Meeting of the
Association for Computational Linguistics**

Companion Volume to the Proceedings of the Conference

**including Short Papers
Proceedings of the Student Research Workshop
Demonstration Notes and
Tutorial Abstracts**

11-16 July 2010
Uppsala University
Uppsala, Sweden

Production and Manufacturing by
Taberg Media Group AB
Box 94, 562 02 Taberg
Sweden

©2010 The Association for Computational Linguistics

Order copies of this and other ACL proceedings from:

Association for Computational Linguistics (ACL)
209 N. Eighth Street
Stroudsburg, PA 18360
USA
Tel: +1-570-476-8006
Fax: +1-570-476-0860
acl@aclweb.org

ISBN 978-1-932432-68-8 / 1-932432-68-X

ACL 2010

**48th Annual Meeting of the
Association for Computational Linguistics**

Proceedings of the Conference Short Papers

11-16 July 2010
Uppsala University
Uppsala, Sweden

Production and Manufacturing by
Taberg Media Group AB
Box 94, 562 02 Taberg
Sweden

©2010 The Association for Computational Linguistics

Order copies of this and other ACL proceedings from:

Association for Computational Linguistics (ACL)
209 N. Eighth Street
Stroudsburg, PA 18360
USA
Tel: +1-570-476-8006
Fax: +1-570-476-0860
acl@aclweb.org

Table of Contents

<i>Paraphrase Lattice for Statistical Machine Translation</i> Takashi Onishi, Masao Utiyama and Eiichiro Sumita	1
<i>A Joint Rule Selection Model for Hierarchical Phrase-Based Translation</i> Lei Cui, Dongdong Zhang, Mu Li, Ming Zhou and Tiejun Zhao	6
<i>Learning Lexicalized Reordering Models from Reordering Graphs</i> Jinsong Su, Yang Liu, Yajuan Lv, Haitao Mi and Qun Liu	12
<i>Filtering Syntactic Constraints for Statistical Machine Translation</i> Hailong Cao and Eiichiro Sumita	17
<i>Diversify and Combine: Improving Word Alignment for Machine Translation on Low-Resource Languages</i> Bing Xiang, Yonggang Deng and Bowen Zhou	22
<i>Efficient Path Counting Transducers for Minimum Bayes-Risk Decoding of Statistical Machine Translation Lattices</i> Graeme Blackwood, Adrià de Gispert and William Byrne	27
<i>The Same-Head Heuristic for Coreference</i> Micha Elsner and Eugene Charniak	33
<i>Authorship Attribution Using Probabilistic Context-Free Grammars</i> Sindhu Raghavan, Adriana Kovashka and Raymond Mooney	38
<i>The Impact of Interpretation Problems on Tutorial Dialogue</i> Myroslava O. Dzikovska, Johanna D. Moore, Natalie Steinhauser and Gwendolyn Campbell ...	43
<i>The Prevalence of Descriptive Referring Expressions in News and Narrative</i> Raquel Hervas and Mark Finlayson	49
<i>Preferences versus Adaptation during Referring Expression Generation</i> Martijn Goudbeek and Emiel Kraemer	55
<i>Cognitively Plausible Models of Human Language Processing</i> Frank Keller	60
<i>The Manually Annotated Sub-Corpus: A Community Resource for and by the People</i> Nancy Ide, Collin Baker, Christiane Fellbaum and Rebecca Passonneau	68
<i>Correcting Errors in a Treebank Based on Synchronous Tree Substitution Grammar</i> Yoshihide Kato and Shigeki Matsubara	74
<i>Evaluating Machine Translations Using mNCD</i> Marcus Dobrinksat, Tero Tapiovaara, Jaakko Väyrynen and Kimmo Kettunen	80
<i>Tackling Sparse Data Issue in Machine Translation Evaluation</i> Ondřej Bojar, Kamil Kos and David Mareček	86
<i>Exemplar-Based Models for Word Meaning in Context</i> Katrin Erk and Sebastian Pado	92

<i>A Structured Model for Joint Learning of Argument Roles and Predicate Senses</i> Yotaro Watanabe, Masayuki Asahara and Yuji Matsumoto	98
<i>Semantics-Driven Shallow Parsing for Chinese Semantic Role Labeling</i> Weiwei Sun	103
<i>Collocation Extraction beyond the Independence Assumption</i> Gerlof Bouma	109
<i>Automatic Collocation Suggestion in Academic Writing</i> Jian-Cheng Wu, Yu-Chia Chang, Teruko Mitamura and Jason S. Chang	115
<i>Event-Based Hyperspace Analogue to Language for Query Expansion</i> Tingxu Yan, Tamsin Maxwell, Dawei Song, Yuexian Hou and Peng Zhang	120
<i>Automatically Generating Term Frequency Induced Taxonomies</i> Karin Murthy, Tanveer A Faruque, L Venkata Subramaniam, Hima Prasad K and Mukesh Mohania 126	
<i>Complexity Assumptions in Ontology Verbalisation</i> Richard Power	132
<i>Word Alignment with Synonym Regularization</i> Hiroyuki Shindo, Akinori Fujino and Masaaki Nagata	137
<i>Better Filtration and Augmentation for Hierarchical Phrase-Based Translation Rules</i> Zhiyang Wang, Yajuan Lv, Qun Liu and Young-Sook Hwang	142
<i>Fixed Length Word Suffix for Factored Statistical Machine Translation</i> Narges Sharif Razavian and Stephan Vogel	147
<i>Unsupervised Discourse Segmentation of Documents with Inherently Parallel Structure</i> Minwoo Jeong and Ivan Titov	151
<i>Coreference Resolution with Reconcile</i> Veselin Stoyanov, Claire Cardie, Nathan Gilbert, Ellen Riloff, David Buttler and David Hysom	156
<i>Predicate Argument Structure Analysis Using Transformation Based Learning</i> Hirotoshi Taira, Sanae Fujita and Masaaki Nagata	162
<i>Improving Chinese Semantic Role Labeling with Rich Syntactic Features</i> Weiwei Sun	168
<i>Balancing User Effort and Translation Error in Interactive Machine Translation via Confidence Measures</i> Jesús González Rubio, Daniel Ortiz Martínez and Francisco Casacuberta	173
<i>Improving Arabic-to-English Statistical Machine Translation by Reordering Post-Verbal Subjects for Alignment</i> Marine Carpuat, Yuval Marton and Nizar Habash	178
<i>Learning Common Grammar from Multilingual Corpus</i> Tomoharu Iwata, Daichi Mochihashi and Hiroshi Sawada	184
<i>Tree-Based Deterministic Dependency Parsing — An Application to Nivre’s Method —</i> Kotaro Kitagawa and Kumiko Tanaka-Ishii	189

<i>Sparsity in Dependency Grammar Induction</i>	
Jennifer Gillenwater, Kuzman Ganchev, João Graça, Fernando Pereira and Ben Taskar	194
<i>Top-Down K-Best A* Parsing</i>	
Adam Pauls, Dan Klein and Chris Quirk	200
<i>Simple Semi-Supervised Training of Part-Of-Speech Taggers</i>	
Anders Søgaard	205
<i>Efficient Optimization of an MDL-Inspired Objective Function for Unsupervised Part-Of-Speech Tagging</i>	
Ashish Vaswani, Adam Pauls and David Chiang	209
<i>SVD and Clustering for Unsupervised POS Tagging</i>	
Michael Lamar, Yariv Maron, Mark Johnson and Elie Bienenstock	215
<i>Intelligent Selection of Language Model Training Data</i>	
Robert C. Moore and William Lewis	220
<i>Blocked Inference in Bayesian Tree Substitution Grammars</i>	
Trevor Cohn and Phil Blunsom	225
<i>Online Generation of Locality Sensitive Hash Signatures</i>	
Benjamin Van Durme and Ashwin Lall	231
<i>Optimizing Question Answering Accuracy by Maximizing Log-Likelihood</i>	
Matthias H. Heie, Edward W. D. Whittaker and Sadaoki Furui	236
<i>Generating Entailment Rules from FrameNet</i>	
Roni Ben Aharon, Idan Szpektor and Ido Dagan	241
<i>Don't 'Have a Clue'? Unsupervised Co-Learning of Downward-Entailing Operators.</i>	
Cristian Danescu-Niculescu-Mizil and Lillian Lee	247
<i>Vocabulary Choice as an Indicator of Perspective</i>	
Beata Beigman Klebanov, Eyal Beigman and Daniel Diermeier	253
<i>Cross Lingual Adaptation: An Experiment on Sentiment Classifications</i>	
Bin Wei and Christopher Pal	258
<i>Using Anaphora Resolution to Improve Opinion Target Identification in Movie Reviews</i>	
Niklas Jakob and Iryna Gurevych	263
<i>Hierarchical Sequential Learning for Extracting Opinions and Their Attributes</i>	
Yejin Choi and Claire Cardie	269
<i>Jointly Optimizing a Two-Step Conditional Random Field Model for Machine Transliteration and Its Fast Decoding Algorithm</i>	
Dong Yang, Paul Dixon and Sadaoki Furui	275
<i>Arabic Named Entity Recognition: Using Features Extracted from Noisy Data</i>	
Yassine Benajiba, Imed Zitouni, Mona Diab and Paolo Rosso	281
<i>Extracting Sequences from the Web</i>	
Anthony Fader, Stephen Soderland and Oren Etzioni	286

<i>An Entity-Level Approach to Information Extraction</i> Aria Haghighi and Dan Klein	291
<i>A Semi-Supervised Key Phrase Extraction Approach: Learning from Title Phrases through a Document Semantic Network</i> Decong Li, Sujian Li, Wenjie Li, Wei Wang and Weiguang Qu	296
<i>Domain Adaptation of Maximum Entropy Language Models</i> Tanel Alumäe and Mikko Kurimo	301
<i>Decision Detection Using Hierarchical Graphical Models</i> Trung H. Bui and Stanley Peters	307
<i>Using Speech to Reply to SMS Messages While Driving: An In-Car Simulator User Study</i> Yun-Cheng Ju and Tim Paek	313
<i>Classification of Feedback Expressions in Multimodal Data</i> Costanza Navarretta and Patrizia Paggio	318
<i>Optimizing Informativeness and Readability for Sentiment Summarization</i> Hitoshi Nishikawa, Takaaki Hasegawa, Yoshihiro Matsuo and Genichiro Kikui	325
<i>Last but Definitely Not Least: On the Role of the Last Sentence in Automatic Polarity-Classification</i> Israela Becker and Vered Aharonson	331
<i>Automatically Generating Annotator Rationales to Improve Sentiment Classification</i> Ainur Yessenalina, Yejin Choi and Claire Cardie	336
<i>Simultaneous Tokenization and Part-Of-Speech Tagging for Arabic without a Morphological Analyzer</i> Seth Kulick	342
<i>Hierarchical A* Parsing with Bridge Outside Scores</i> Adam Pauls and Dan Klein	348
<i>Using Parse Features for Preposition Selection and Error Detection</i> Joel Tetreault, Jennifer Foster and Martin Chodorow	353
<i>Distributional Similarity vs. PU Learning for Entity Set Expansion</i> Xiao-Li Li, Lei Zhang, Bing Liu and See-Kiong Ng	359
<i>Active Learning-Based Elicitation for Semi-Supervised Word Alignment</i> Vamshi Ambati, Stephan Vogel and Jaime Carbonell	365
<i>An Active Learning Approach to Finding Related Terms</i> David Vickrey, Oscar Kipersztok and Daphne Koller	371
<i>Learning Better Data Representation Using Inference-Driven Metric Learning</i> Paramveer S. Dhillon, Partha Pratim Talukdar and Koby Crammer	377
<i>Wrapping up a Summary: From Representation to Generation</i> Josef Steinberger, Marco Turchi, Mijail Kabadjov, Ralf Steinberger and Nello Cristianini	382

Paraphrase Lattice for Statistical Machine Translation

Takashi Onishi and Masao Utiyama and Eiichiro Sumita

Language Translation Group, MASTAR Project

National Institute of Information and Communications Technology

3-5 Hikaridai, Keihanna Science City, Kyoto, 619-0289, JAPAN

{takashi.onishi,mutiyama,eiichiro.sumita}@nict.go.jp

Abstract

Lattice decoding in statistical machine translation (SMT) is useful in speech translation and in the translation of German because it can handle input ambiguities such as speech recognition ambiguities and German word segmentation ambiguities. We show that lattice decoding is also useful for handling input variations. Given an input sentence, we build a lattice which represents paraphrases of the input sentence. We call this a paraphrase lattice. Then, we give the paraphrase lattice as an input to the lattice decoder. The decoder selects the best path for decoding. Using these paraphrase lattices as inputs, we obtained significant gains in BLEU scores for IWSLT and Europarl datasets.

1 Introduction

Lattice decoding in SMT is useful in speech translation and in the translation of German (Bertoldi et al., 2007; Dyer, 2009). In speech translation, by using lattices that represent not only 1-best result but also other possibilities of speech recognition, we can take into account the ambiguities of speech recognition. Thus, the translation quality for lattice inputs is better than the quality for 1-best inputs.

In this paper, we show that lattice decoding is also useful for handling input variations. “Input variations” refers to the differences of input texts with the same meaning. For example, “*Is there a beauty salon?*” and “*Is there a beauty parlor?*” have the same meaning with variations in “*beauty salon*” and “*beauty parlor*”. Since these variations are frequently found in natural language texts, a mismatch of the expressions in source sentences and the expressions in training corpus leads to a decrease in translation quality. Therefore,

we propose a novel method that can handle input variations using paraphrases and lattice decoding. In the proposed method, we regard a given source sentence as one of many variations (1-best). Given an input sentence, we build a paraphrase lattice which represents paraphrases of the input sentence. Then, we give the paraphrase lattice as an input to the Moses decoder (Koehn et al., 2007). Moses selects the best path for decoding. By using paraphrases of source sentences, we can translate expressions which are not found in a training corpus on the condition that paraphrases of them are found in the training corpus. Moreover, by using lattice decoding, we can employ the source-side language model as a decoding feature. Since this feature is affected by the source-side context, the decoder can choose a proper paraphrase and translate correctly.

This paper is organized as follows: Related works on lattice decoding and paraphrasing are presented in Section 2. The proposed method is described in Section 3. Experimental results for IWSLT and Europarl dataset are presented in Section 4. Finally, the paper is concluded with a summary and a few directions for future work in Section 5.

2 Related Work

Lattice decoding has been used to handle ambiguities of preprocessing. Bertoldi et al. (2007) employed a confusion network, which is a kind of lattice and represents speech recognition hypotheses in speech translation. Dyer (2009) also employed a segmentation lattice, which represents ambiguities of compound word segmentation in German, Hungarian and Turkish translation. However, to the best of our knowledge, there is no work which employed a lattice representing paraphrases of an input sentence.

On the other hand, paraphrasing has been used to enrich the SMT model. Callison-Burch et

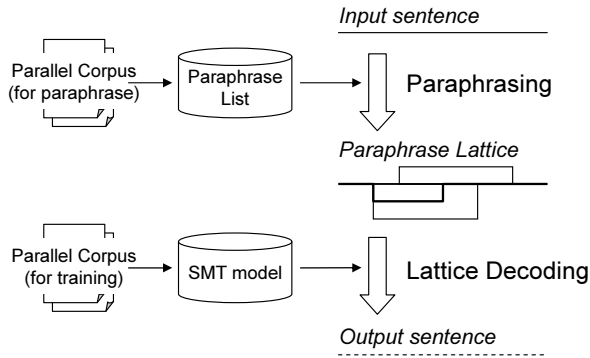


Figure 1: Overview of the proposed method.

al. (2006) and Marton et al. (2009) augmented the translation phrase table with paraphrases to translate unknown phrases. Bond et al. (2008) and Nakov (2008) augmented the training data by paraphrasing. However, there is no work which augments input sentences by paraphrasing and represents them in lattices.

3 Paraphrase Lattice for SMT

Overview of the proposed method is shown in Figure 1. In advance, we automatically acquire a paraphrase list from a parallel corpus. In order to acquire paraphrases of unknown phrases, this parallel corpus is different from the parallel corpus for training.

Given an input sentence, we build a lattice which represents paraphrases of the input sentence using the paraphrase list. We call this lattice a paraphrase lattice. Then, we give the paraphrase lattice to the lattice decoder.

3.1 Acquiring the paraphrase list

We acquire a paraphrase list using Bannard and Callison-Burch (2005)’s method. Their idea is, if two different phrases e_1 , e_2 in one language are aligned to the same phrase c in another language, they are hypothesized to be paraphrases of each other. Our paraphrase list is acquired in the same way.

The procedure is as follows:

1. Build a phrase table.
Build a phrase table from parallel corpus using standard SMT techniques.
2. Filter the phrase table by the sigtest-filter.
The phrase table built in 1 has many inappropriate phrase pairs. Therefore, we filter the

phrase table and keep only appropriate phrase pairs using the sigtest-filter (Johnson et al., 2007).

3. Calculate the paraphrase probability.
Calculate the paraphrase probability $p(e_2|e_1)$ if e_2 is hypothesized to be a paraphrase of e_1 .

$$p(e_2|e_1) = \sum_c P(c|e_1)P(e_2|c)$$

where $P(\cdot|\cdot)$ is phrase translation probability.

4. Acquire a paraphrase pair.
Acquire (e_1, e_2) as a paraphrase pair if $p(e_2|e_1) > p(e_1|e_1)$. The purpose of this threshold is to keep highly-accurate paraphrase pairs. In experiments, more than 80% of paraphrase pairs were eliminated by this threshold.

3.2 Building paraphrase lattice

An input sentence is paraphrased using the paraphrase list and transformed into a paraphrase lattice. The paraphrase lattice is a lattice which represents paraphrases of the input sentence. An example of a paraphrase lattice is shown in Figure 2. In this example, an input sentence is “*is there a beauty salon ?*”. This paraphrase lattice contains two paraphrase pairs “*beauty salon*” = “*beauty parlor*” and “*beauty salon*” = “*salon*”, and represents following three sentences.

- *is there a beauty salon ?*
- *is there a beauty parlor ?*
- *is there a salon ?*

In the paraphrase lattice, each node consists of a token, the distance to the next node and features for lattice decoding. We use following four features for lattice decoding.

- Paraphrase probability (p)
A paraphrase probability $p(e_2|e_1)$ calculated when acquiring the paraphrase.

$$h_p = p(e_2|e_1)$$

- Language model score (l)
A ratio between the language model probability of the paraphrased sentence (para) and that of the original sentence (orig).

$$h_l = \frac{lm(para)}{lm(orig)}$$

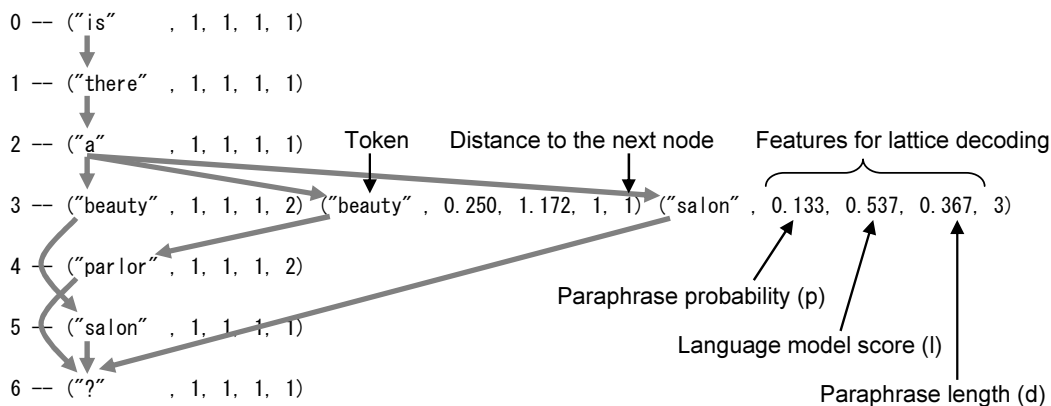


Figure 2: An example of a paraphrase lattice, which contains three features of (p, l, d).

- Normalized language model score (L)

A language model score where the language model probability is normalized by the sentence length. The sentence length is calculated as the number of tokens.

$$h_L = \frac{LM(para)}{LM(orig)},$$

where $LM(sent) = lm(sent)^{\frac{1}{length(sent)}}$

- Paraphrase length (d)

The difference between the original sentence length and the paraphrased sentence length.

$$h_d = \exp(length(para) - length(orig))$$

The values of these features are calculated only if the node is the first node of the paraphrase, for example the second “beauty” and “salon” in line 3 of Figure 2. In other nodes, for example “parlor” in line 4 and original nodes, we use 1 as the values of features.

The features related to the language model, such as (l) and (L), are affected by the context of source sentences even if the same paraphrase pair is applied. As these features can penalize paraphrases which are not appropriate to the context, appropriate paraphrases are chosen and appropriate translations are output in lattice decoding. The features related to the sentence length, such as (L) and (d), are added to penalize the language model score in case the paraphrased sentence length is shorter than the original sentence length and the language model score is unreasonably low.

In experiments, we use four combinations of these features, (p), (p, l), (p, L) and (p, l, d).

3.3 Lattice decoding

We use Moses (Koehn et al., 2007) as a decoder for lattice decoding. Moses is an open source

SMT system which allows lattice decoding. In lattice decoding, Moses selects the best path and the best translation according to features added in each node and other SMT features. These weights are optimized using Minimum Error Rate Training (MERT) (Och, 2003).

4 Experiments

In order to evaluate the proposed method, we conducted English-to-Japanese and English-to-Chinese translation experiments using IWSLT 2007 (Fordyce, 2007) dataset. This dataset contains EJ and EC parallel corpus for the travel domain and consists of 40k sentences for training and about 500 sentences sets (dev1, dev2 and dev3) for development and testing. We used the dev1 set for parameter tuning, the dev2 set for choosing the setting of the proposed method, which is described below, and the dev3 set for testing.

The English-English paraphrase list was acquired from the EC corpus for EJ translation and 53K pairs were acquired. Similarly, 47K pairs were acquired from the EJ corpus for EC translation.

4.1 Baseline

As baselines, we used Moses and Callison-Burch et al. (2006)’s method (hereafter CCB). In Moses, we used default settings without paraphrases. In CCB, we paraphrased the phrase table using the automatically acquired paraphrase list. Then, we augmented the phrase table with paraphrased phrases which were not found in the original phrase table. Moreover, we used an additional feature whose value was the paraphrase probability (p) if the entry was generated by paraphrasing and

	Moses (w/o Paraphrases)	CCB	Proposed Method
EJ	38.98	39.24 (+0.26)	40.34 (+1.36)
EC	25.11	26.14 (+1.03)	27.06 (+1.95)

Table 1: Experimental results for IWSLT (%BLEU).

1 if otherwise. Weights of the feature and other features in SMT were optimized using MERT.

4.2 Proposed method

In the proposed method, we conducted experiments with various settings for paraphrasing and lattice decoding. Then, we chose the best setting according to the result of the dev2 set.

4.2.1 Limitation of paraphrasing

As the paraphrase list was automatically acquired, there were many erroneous paraphrase pairs. Building paraphrase lattices with all erroneous paraphrase pairs and decoding these paraphrase lattices caused high computational complexity. Therefore, we limited the number of paraphrasing per phrase and per sentence. The number of paraphrasing per phrase was limited to three and the number of paraphrasing per sentence was limited to twice the size of the sentence length.

As a criterion for limiting the number of paraphrasing, we use three features (p), (l) and (L), which are same as the features described in Subsection 3.2. When building paraphrase lattices, we apply paraphrases in descending order of the value of the criterion.

4.2.2 Finding optimal settings

As previously mentioned, we have three choices for the criterion for building paraphrase lattices and four combinations of features for lattice decoding. Thus, there are $3 \times 4 = 12$ combinations of these settings. We conducted parameter tuning with the dev1 set for each setting and used as best the setting which got the highest BLEU score for the dev2 set.

4.3 Results

The experimental results are shown in Table 1. We used the case-insensitive BLEU metric for evaluation. In EJ translation, the proposed method obtained the highest score of 40.34%, which achieved an absolute improvement of 1.36 BLEU points over Moses and 1.10 BLEU points over CCB. In EC translation, the proposed method also obtained the highest score of 27.06% and achieved

an absolute improvement of 1.95 BLEU points over Moses and 0.92 BLEU points over CCB. As the relation of three systems is Moses < CCB < Proposed Method, paraphrasing is useful for SMT and using paraphrase lattices and lattice decoding is especially more useful than augmenting the phrase table. In Proposed Method, the criterion for building paraphrase lattices and the combination of features for lattice decoding were (p) and (p, L) in EJ translation and (L) and (p, l) in EC translation. Since features related to the source-side language model were chosen in each direction, using the source-side language model is useful for decoding paraphrase lattices.

We also tried a combination of Proposed Method and CCB, which is a method of decoding paraphrase lattices with an augmented phrase table. However, the result showed no significant improvements. This is because the proposed method includes the effect of augmenting the phrase table.

Moreover, we conducted German-English translation using the Europarl corpus (Koehn, 2005). We used the WMT08 dataset¹, which consists of 1M sentences for training and 2K sentences for development and testing. We acquired 5.3M pairs of German-German paraphrases from a 1M German-Spanish parallel corpus. We conducted experiments with various sizes of training corpus, using 10K, 20K, 40K, 80K, 160K and 1M. Figure 3 shows the proposed method consistently get higher score than Moses and CCB.

5 Conclusion

This paper has proposed a novel method for transforming a source sentence into a paraphrase lattice and applying lattice decoding. Since our method can employ source-side language models as a decoding feature, the decoder can choose proper paraphrases and translate properly. The experimental results showed significant gains for the IWSLT and Europarl dataset. In IWSLT dataset, we obtained 1.36 BLEU points over Moses in EJ translation and 1.95 BLEU points over Moses in

¹<http://www.statmt.org/wmt08/>

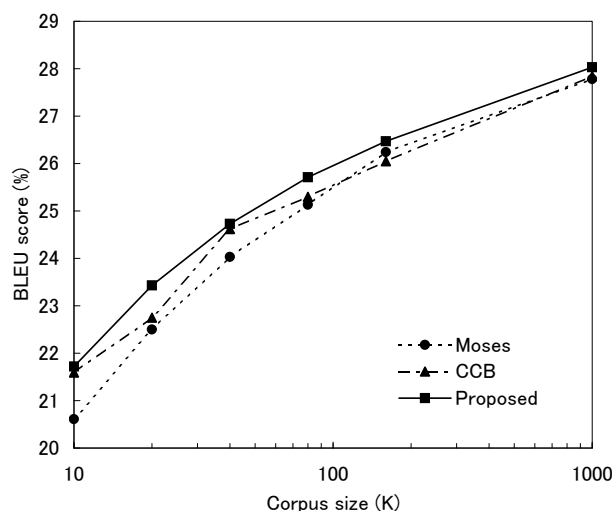


Figure 3: Effect of training corpus size.

EC translation. In Europarl dataset, the proposed method consistently get higher score than baselines.

In future work, we plan to apply this method with paraphrases derived from a massive corpus such as the Web corpus and apply this method to a hierarchical phrase based SMT.

References

- Colin Bannard and Chris Callison-Burch. 2005. Paraphrasing with Bilingual Parallel Corpora. In *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 597–604.
- Nicola Bertoldi, Richard Zens, and Marcello Federico. 2007. Speech translation by confusion network decoding. In *Proceedings of the International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, pages 1297–1300.
- Francis Bond, Eric Nichols, Darren Scott Appling, and Michael Paul. 2008. Improving Statistical Machine Translation by Paraphrasing the Training Data. In *Proceedings of the International Workshop on Spoken Language Translation (IWSLT)*, pages 150–157.
- Chris Callison-Burch, Philipp Koehn, and Miles Osborne. 2006. Improved Statistical Machine Translation Using Paraphrases. In *Proceedings of the Human Language Technology conference - North American chapter of the Association for Computational Linguistics (HLT-NAACL)*, pages 17–24.
- Chris Dyer. 2009. Using a maximum entropy model to build segmentation lattices for MT. In *Proceedings of the Human Language Technology conference - North American chapter of the Association for Computational Linguistics (HLT-NAACL)*, pages 406–414.
- Cameron S. Fordyce. 2007. Overview of the IWSLT 2007 Evaluation Campaign. In *Proceedings of the International Workshop on Spoken Language Translation (IWSLT)*, pages 1–12.
- J Howard Johnson, Joel Martin, George Foster, and Roland Kuhn. 2007. Improving Translation Quality by Discarding Most of the Phrasetable. In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 967–975.
- Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondrej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open Source Toolkit for Statistical Machine Translation. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 177–180.
- Philipp Koehn. 2005. Europarl: A Parallel Corpus for Statistical Machine Translation. In *Proceedings of the 10th Machine Translation Summit (MT Summit)*, pages 79–86.
- Yuval Marton, Chris Callison-Burch, and Philip Resnik. 2009. Improved Statistical Machine Translation Using Monolingually-Derived Paraphrases. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 381–390.
- Preslav Nakov. 2008. Improved Statistical Machine Translation Using Monolingual Paraphrases. In *Proceedings of the European Conference on Artificial Intelligence (ECAI)*, pages 338–342.
- Franz Josef Och. 2003. Minimum Error Rate Training in Statistical Machine Translation. In *Proceedings of the 41st Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 160–167.

A Joint Rule Selection Model for Hierarchical Phrase-based Translation*

Lei Cui[†], Dongdong Zhang[‡], Mu Li[‡], Ming Zhou[‡], and Tiejun Zhao[†]

[†]School of Computer Science and Technology
Harbin Institute of Technology, Harbin, China
{cuilei, tjzhao}@mtlab.hit.edu.cn

[‡]Microsoft Research Asia, Beijing, China
{dozhang, muli, mingzhou}@microsoft.com

Abstract

In hierarchical phrase-based SMT systems, statistical models are integrated to guide the hierarchical rule selection for better translation performance. Previous work mainly focused on the selection of either the source side of a hierarchical rule or the target side of a hierarchical rule rather than considering both of them simultaneously. This paper presents a joint model to predict the selection of hierarchical rules. The proposed model is estimated based on four sub-models where the rich context knowledge from both source and target sides is leveraged. Our method can be easily incorporated into the practical SMT systems with the log-linear model framework. The experimental results show that our method can yield significant improvements in performance.

1 Introduction

Hierarchical phrase-based model has strong expression capabilities of translation knowledge. It can not only maintain the strength of phrase translation in traditional phrase-based models (Koehn et al., 2003; Xiong et al., 2006), but also characterize the complicated long distance reordering similar to syntactic based statistical machine translation (SMT) models (Yamada and Knight, 2001; Quirk et al., 2005; Galley et al., 2006; Liu et al., 2006; Marcu et al., 2006; Mi et al., 2008; Shen et al., 2008).

In hierarchical phrase-based SMT systems, due to the flexibility of rule matching, a huge number of hierarchical rules could be automatically learnt from bilingual training corpus (Chiang, 2005). SMT decoders are forced to face the challenge of

proper rule selection for hypothesis generation, including both source-side rule selection and target-side rule selection where the source-side rule determines what part of source words to be translated and the target-side rule provides one of the candidate translations of the source-side rule. Improper rule selections may result in poor translations.

There is some related work about the hierarchical rule selection. In the original work (Chiang, 2005), the target-side rule selection is analogous to the model in traditional phrase-based SMT system such as Pharaoh (Koehn et al., 2003). Extending this work, (He et al., 2008; Liu et al., 2008) integrate rich context information of non-terminals to predict the target-side rule selection. Different from the above work where the probability distribution of source-side rule selection is uniform, (Setiawan et al., 2009) proposes to select source-side rules based on the captured function words which often play an important role in word reordering. There is also some work considering to involve more rich contexts to guide the source-side rule selection. (Marton and Resnik, 2008; Xiong et al., 2009) explore the source syntactic information to reward exact matching structure rules or punish crossing structure rules.

All the previous work mainly focused on either source-side rule selection task or target-side rule selection task rather than both of them together. The separation of these two tasks, however, weakens the high interrelation between them. In this paper, we propose to integrate both source-side and target-side rule selection in a unified model. The intuition is that the joint selection of source-side and target-side rules is more reliable as it conducts the search in a larger space than the single selection task does. It is expected that these two kinds of selection can help and affect each other, which may potentially lead to better hierarchical rule selections with a relative global optimum instead of a local optimum that might be reached in the pre-

This work was finished while the first author visited Microsoft Research Asia as an intern.

vious methods. Our proposed joint probability model is factored into four sub-models that can be further classified into source-side and target-side rule selection models or context-based and context-free selection models. The context-based models explore rich context features from both source and target sides, including function words, part-of-speech (POS) tags, syntactic structure information and so on. Our model can be easily incorporated as an independent feature into the practical hierarchical phrase-based systems with the log-linear model framework. The experimental results indicate our method can improve the system performance significantly.

2 Hierarchical Rule Selection Model

Following (Chiang, 2005), $\langle \alpha, \gamma \rangle$ is used to represent a synchronous context free grammar (SCFG) rule extracted from the training corpus, where α and γ are the source-side and target-side rule respectively. Let C be the context of $\langle \alpha, \gamma \rangle$. Formally, our joint probability model of hierarchical rule selection is described as follows:

$$P(\alpha, \gamma | C) = P(\alpha | C)P(\gamma | \alpha, C) \quad (1)$$

We decompose the joint probability model into two sub-models based on the Bayes formulation, where the first sub-model is *source-side rule selection model* and the second one is the *target-side rule selection model*.

For the source-side rule selection model, we further compute it by the interpolation of two sub-models:

$$\theta P_s(\alpha) + (1 - \theta)P_s(\alpha | C) \quad (2)$$

where $P_s(\alpha)$ is the *context-free source model* (CFSM) and $P_s(\alpha | C)$ is the *context-based source model* (CBSM), θ is the interpolation weight that can be optimized over the development data.

CFSM is the probability of source-side rule selection that can be estimated based on maximum likelihood estimation (MLE) method:

$$P_s(\alpha) = \frac{\sum_{\gamma} \text{Count}(\langle \alpha, \gamma \rangle)}{\text{Count}(\alpha)} \quad (3)$$

where the numerator is the total count of bilingual rule pairs with the same source-side rule that are extracted based on the extraction algorithm in (Chiang, 2005), and the denominator is the total amount of source-side rule patterns contained in

the monolingual source side of the training corpus. CFSM is used to capture how likely the source-side rule is linguistically motivated or has the corresponding target-side counterpart.

For CBSM, it can be naturally viewed as a classification problem where each distinct source-side rule is a single class. However, considering the huge number of classes may cause serious data sparseness problem and thereby degrade the classification accuracy, we approximate CBSM by a binary classification problem which can be solved by the maximum entropy (ME) approach (Berger et al., 1996) as follows:

$$\begin{aligned} P_s(\alpha | C) &\approx P_s(v | \alpha, C) \\ &= \frac{\exp[\sum_i \lambda_i h_i(v, \alpha, C)]}{\sum_{v'} \exp[\sum_i \lambda_i h_i(v', \alpha, C)]} \end{aligned} \quad (4)$$

where $v \in \{0, 1\}$ is the indicator whether the source-side rule is applied during decoding, $v = 1$ when the source-side rule is applied, otherwise $v = 0$; h_i is a feature function, λ_i is the weight of h_i . CBSM estimates the probability of the source-side rule being selected according to the rich context information coming from the surface strings and sub-phrases that will be reduced to non-terminals during decoding.

Analogously, we decompose the target-side rule selection model by the interpolation approach as well:

$$\varphi P_t(\gamma) + (1 - \varphi)P_t(\gamma | \alpha, C) \quad (5)$$

where $P_t(\gamma)$ is the *context-free target model* (CFTM) and $P_t(\gamma | \alpha, C)$ is the *context-based target model* (CBTM), φ is the interpolation weight that can be optimized over the development data.

In the similar way, we compute CFTM by the MLE approach and estimate CBTM by the ME approach. CFTM computes how likely the target-side rule is linguistically motivated, while CBTM predicts how likely the target-side rule is applied according to the clues from the rich context information.

3 Model Training of CBSM and CBTM

3.1 The acquisition of training instances

CBSM and CBTM are trained by ME approach for the binary classification, where a training instance consists of a label and the context related to SCFG rules. The context is divided into source context

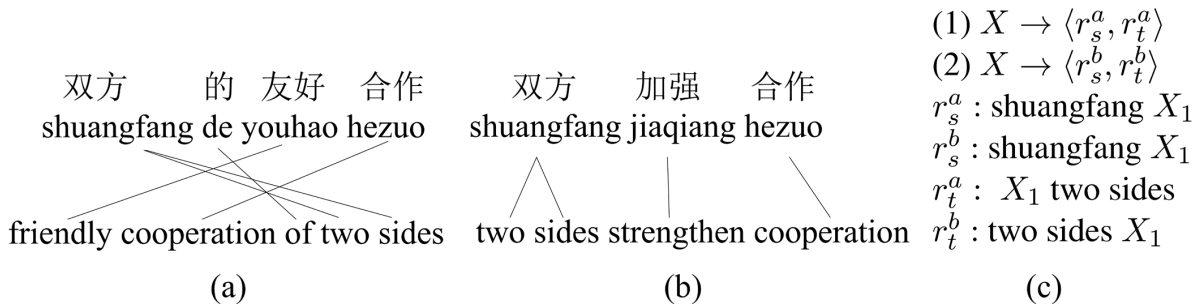


Figure 1: Example of training instances in CBSM and CBTM.

and target context. CBSM is trained only based on the source context while CBTM is trained over both the source and the target context. All the training instances are automatically constructed from the bilingual training corpus, which have labels of either positive (i.e., $v = 1$) or negative (i.e., $v = 0$). This section explains how the training instances are constructed for the training of CBSM and CBTM.

Let s and t be the source sentence and target sentence, W be the word alignment between them, r_s be a source-side rule that pattern-matches a sub-phrase of s , r_t be the target-side rule pattern-matching a sub-phrase of t and being aligned to r_s based on W , and $C(r)$ be the context features related to the rule r which will be explained in the following section.

For the training of CBSM, if the SCFG rule $\langle r_s, r_t \rangle$ can be extracted based on the rule extraction algorithm in (Chiang, 2005), $\langle v = 1, C(r_s) \rangle$ is constructed as a positive instance, otherwise $\langle v = 0, C(r_s) \rangle$ is constructed as a negative instance. For example in Figure 1(a), the context of source-side rule " X_1 hezuo" that pattern-matches the phrase "youhao hezuo" produces a positive instance, while the context of " X_1 youhao" that pattern-matches the source phrase "de youhao" or "shuangfang de youhao" will produce a negative instance as there are no corresponding plausible target-side rules that can be extracted legally¹.

For the training of CBTM, given r_s , suppose there is a SCFG rule set $\{\langle r_s, r_t^k \rangle | 1 \leq k \leq n\}$ extracted from multiple distinct sentence pairs in the bilingual training corpus, among which we assume $\langle r_s, r_t^i \rangle$ is extracted from the sentence pair $\langle s, t \rangle$. Then, we construct $\langle v = 1, C(r_s), C(r_t^i) \rangle$

¹Because the aligned target words are not contiguous and "cooperation" is aligned to the word outside the source-side rule.

as a positive instance, while the elements in $\{\langle v = 0, C(r_s), C(r_t^j) \rangle | j \neq i \wedge 1 \leq j \leq n\}$ are viewed as negative instances since they fail to be applied to the translation from s to t . For example in Figure 1(c), Rule (1) and Rule (2) are two different SCFG rules extracted from Figure 1(a) and Figure 1(b) respectively, where their source-side rules are the same. As Rule (1) cannot be applied to Figure 1(b) for the translation and Rule (2) cannot be applied to Figure 1(a) for the translation either, $\langle v = 1, C(r_s^a), C(r_t^a) \rangle$ and $\langle v = 1, C(r_s^b), C(r_t^b) \rangle$ are constructed as positive instances while $\langle v = 0, C(r_s^a), C(r_t^b) \rangle$ and $\langle v = 0, C(r_s^b), C(r_t^a) \rangle$ are viewed as negative instances. It is noticed that this instance construction method may lead to a large quantity of negative instances and choke the training procedure. In practice, to limit the size of the training set, the negative instances constructed based on low-frequency target-side rules are pruned.

3.2 Context-based features for ME training

ME approach has the merit of easily combining different features to predict the probability of each class. We incorporate into the ME based model the following informative context-based features to train CBSM and CBTM. These features are carefully designed to reduce the data sparseness problem and some of them are inspired by previous work (He et al., 2008; Gimpel and Smith, 2008; Marton and Resnik, 2008; Chiang et al., 2009; Setiawan et al., 2009; Shen et al., 2009; Xiong et al., 2009):

1. **Function word features**, which indicate whether the hierarchical source-side/target-side rule strings and sub-phrases covered by non-terminals contain function words that are often important clues of predicting syntactic structures.

2. **POS features**, which are POS tags of the boundary source words covered by non-terminals.
3. **Syntactic features**, which are the constituent constraints of hierarchical source-side rules exactly matching or crossing syntactic subtrees.
4. **Rule format features**, which are non-terminal positions and orders in source-side/target-side rules. This feature interacts between source and target components since it shows whether the translation ordering is affected.
5. **Length features**, which are the length of sub-phrases covered by source non-terminals.

4 Experiments

4.1 Experiment setting

We implement a hierarchical phrase-based system similar to the Hiero (Chiang, 2005) and evaluate our method on the Chinese-to-English translation task. Our bilingual training data comes from FBIS corpus, which consists of around 160K sentence pairs where the source data is parsed by the Berkeley parser (Petrov and Klein, 2007). The ME training toolkit, developed by (Zhang, 2006), is used to train our CBSM and CBTM. The training size of constructed positive instances for both CBSM and CBTM is 4.68M, while the training size of constructed negative instances is 3.74M and 3.03M respectively. Following (Setiawan et al., 2009), we identify function words as the 128 most frequent words in the corpus. The interpolation weights are set to $\theta = 0.75$ and $\varphi = 0.70$. The 5-gram language model is trained over the English portion of FBIS corpus plus Xinhua portion of the Gigaword corpus. The development data is from NIST 2005 evaluation data and the test data is from NIST 2006 and NIST 2008 evaluation data. The evaluation metric is the case-insensitive BLEU4 (Papineni et al., 2002). Statistical significance in BLEU score differences is tested by paired bootstrap re-sampling (Koehn, 2004).

4.2 Comparison with related work

Our baseline is the implemented Hiero-like SMT system where only the standard features are employed and the performance is state-of-the-art.

We compare our method with the baseline and some typical approaches listed in Table 1 where XP+ denotes the approach in (Marton and Resnik, 2008) and TOFW (topological ordering of function words) stands for the method in (Setiawan et al., 2009). As (Xiong et al., 2009)’s work is based on phrasal SMT system with bracketing transduction grammar rules (Wu, 1997) and (Shen et al., 2009)’s work is based on the string-to-dependency SMT model, we do not implement these two related work due to their different models from ours. We also do not compare with (He et al., 2008)’s work due to its less practicability of integrating numerous sub-models.

Methods	NIST 2006	NIST 2008
Baseline	0.3025	0.2200
XP+	0.3061	0.2254
TOFW	0.3089	0.2253
Our method	0.3141	0.2318

Table 1: Comparison results, our method is significantly better than the baseline, as well as the other two approaches ($p < 0.01$)

As shown in Table 1, all the methods outperform the baseline because they have extra models to guide the hierarchical rule selection in some ways which might lead to better translation. Apparently, our method also performs better than the other two approaches, indicating that our method is more effective in the hierarchical rule selection as both source-side and target-side rules are selected together.

4.3 Effect of sub-models

Due to the space limitation, we analyze the effect of sub-models upon the system performance, rather than that of ME features, part of which have been investigated in previous related work.

Settings	NIST 2006	NIST 2008
Baseline	0.3025	0.2200
Baseline+CFSM	0.3092*	0.2266*
Baseline+CBSM	0.3077*	0.2247*
Baseline+CFTM	0.3076*	0.2286*
Baseline+CBTM	0.3060	0.2255*
Baseline+CFSM+CFTM	0.3109*	0.2289*
Baseline+CFSM+CBSM	0.3104*	0.2282*
Baseline+CFTM+CBTM	0.3099*	0.2299*
Baseline+all sub-models	0.3141*	0.2318*

Table 2: Sub-model effect upon the performance, *: significantly better than baseline ($p < 0.01$)

As shown in Table 2, when sub-models are inte-

grated as independent features, the performance is improved compared to the baseline, which shows that each of the sub-models can improve the hierarchical rule selection. It is noticeable that the performance of the source-side rule selection model is comparable with that of the target-side rule selection model. Although CFSM and CFTM perform only slightly better than the others among the individual sub-models, the best performance is achieved when all the sub-models are integrated.

5 Conclusion

Hierarchical rule selection is an important and complicated task for hierarchical phrase-based SMT system. We propose a joint probability model for the hierarchical rule selection and the experimental results prove the effectiveness of our approach.

In the future work, we will explore more useful features and test our method over the large scale training corpus. A challenge might exist when running the ME training toolkit over a big size of training instances from the large scale training data.

Acknowledgments

We are especially grateful to the anonymous reviewers for their insightful comments. We also thank Hendra Setiawan, Yuval Marton, Chi-Ho Li, Shujie Liu and Nan Duan for helpful discussions.

References

- Adam L. Berger, Vincent J. Della Pietra, and Stephen A. Della Pietra. 1996. *A Maximum Entropy Approach to Natural Language Processing*. *Computational Linguistics*, 22(1): pages 39-72.
- David Chiang. 2005. *A Hierarchical Phrase-Based Model for Statistical Machine Translation*. In *Proc. ACL*, pages 263-270.
- David Chiang, Kevin Knight, and Wei Wang. 2009. *11,001 New Features for Statistical Machine Translation*. In *Proc. HLT-NAACL*, pages 218-226.
- Michel Galley, Jonathan Graehl, Kevin Knight, Daniel Marcu, Steve DeNeefe, Wei Wang, and Ignacio Thayer. 2006. *Scalable Inference and Training of Context-Rich Syntactic Translation Models*. In *Proc. ACL-Coling*, pages 961-968.
- Kevin Gimpel and Noah A. Smith. 2008. *Rich Source-Side Context for Statistical Machine Translation*. In *Proc. the Third Workshop on Statistical Machine Translation*, pages 9-17.
- Zhongjun He, Qun Liu, and Shouxun Lin. 2008. *Improving Statistical Machine Translation using Lexicalized Rule Selection*. In *Proc. Coling*, pages 321-328.
- Philipp Koehn. 2004. *Statistical Significance Tests for Machine Translation Evaluation*. In *Proc. EMNLP*.
- Philipp Koehn, Franz J. Och, and Daniel Marcu. 2003. *Statistical Phrase-Based Translation*. In *Proc. HLT-NAACL*, pages 127-133.
- Qun Liu, Zhongjun He, Yang Liu, and Shouxun Lin. 2008. *Maximum Entropy based Rule Selection Model for Syntax-based Statistical Machine Translation*. In *Proc. EMNLP*, pages 89-97.
- Yang Liu, Yun Huang, Qun Liu, and Shouxun Lin. 2007. *Forest-to-String Statistical Translation Rules*. In *Proc. ACL*, pages 704-711.
- Yang Liu, Qun Liu, and Shouxun Lin. 2006. *Tree-to-String Alignment Template for Statistical Machine Translation*. In *Proc. ACL-Coling*, pages 609-616.
- Daniel Marcu, Wei Wang, Abdessamad Echihabi, and Kevin Knight. 2006. *SPMT: Statistical Machine Translation with Syntactified Target Language Phrases*. In *Proc. EMNLP*, pages 44-52.
- Yuval Marton and Philip Resnik. 2008. *Soft Syntactic Constraints for Hierarchical Phrasal-Based Translation*. In *Proc. ACL*, pages 1003-1011.
- Haitao Mi, Liang Huang, and Qun Liu. 2008. *Forest-Based Translation*. In *Proc. ACL*, pages 192-199.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. *Bleu: a Method for Automatic Evaluation of Machine Translation*. In *Proc. ACL*, pages 311-318.
- Slav Petrov and Dan Klein. 2007. *Improved Inference for Unlexicalized Parsing*. In *Proc. HLT-NAACL*, pages 404-411.
- Chris Quirk, Arul Menezes, and Colin Cherry. 2005. *Dependency Treelet Translation: Syntactically Informed Phrasal SMT*. In *Proc. ACL*, pages 271-279.
- Libin Shen, Jinxi Xu, and Ralph Weischedel. 2008. *A New String-to-Dependency Machine Translation Algorithm with a Target Dependency Language Model*. In *Proc. ACL*, pages 577-585.
- Libin Shen, Jinxi Xu, Bing Zhang, Spyros Matsoukas, and Ralph Weischedel. 2009. *Effective Use of Linguistic and Contextual Information for Statistical Machine Translation*. In *Proc. EMNLP*, pages 72-80.
- Hendra Setiawan, Min Yen Kan, Haizhou Li, and Philip Resnik. 2009. *Topological Ordering of Function Words in Hierarchical Phrase-based Translation*. In *Proc. ACL*, pages 324-332.

- Dekai Wu. 1997. *Stochastic Inversion Transduction Grammars and Bilingual Parsing of Parallel Corpora*. *Computational Linguistics*, 23(3): pages 377-403.
- Deyi Xiong, Qun Liu, and Shouxun Lin. 2006. *Maximum Entropy Based Phrase Reordering Model for Statistical Machine Translation*. In *Proc. ACL-Coling*, pages 521-528.
- Deyi Xiong, Min Zhang, Aiti Aw, and Haizhou Li. 2009. *A Syntax-Driven Bracketing Model for Phrase-Based Translation*. In *Proc. ACL*, pages 315-323.
- Kenji Yamada and Kevin Knight. 2001. *A Syntax-based Statistical Translation Model*. In *Proc. ACL*, pages 523-530.
- Le Zhang. 2006. *Maximum entropy modeling toolkit for python and c++*. available at http://homepages.inf.ed.ac.uk/lzhang10/maxent_toolkit.html.

Learning Lexicalized Reordering Models from Reordering Graphs

Jinsong Su, Yang Liu, Yajuan Lü, Haitao Mi, Qun Liu

Key Laboratory of Intelligent Information Processing

Institute of Computing Technology

Chinese Academy of Sciences

P.O. Box 2704, Beijing 100190, China

{sujinsong, yliu, lvayajuan, htmi, liuqun}@ict.ac.cn

Abstract

Lexicalized reordering models play a crucial role in phrase-based translation systems. They are usually learned from the word-aligned bilingual corpus by examining the reordering relations of adjacent phrases. Instead of just checking whether there is one phrase adjacent to a given phrase, we argue that it is important to take the number of adjacent phrases into account for better estimations of reordering models. We propose to use a structure named *reordering graph*, which represents all phrase segmentations of a sentence pair, to learn lexicalized reordering models efficiently. Experimental results on the NIST Chinese-English test sets show that our approach significantly outperforms the baseline method.

1 Introduction

Phrase-based translation systems (Koehn et al., 2003; Och and Ney, 2004) prove to be the state-of-the-art as they have delivered translation performance in recent machine translation evaluations. While excelling at memorizing local translation and reordering, phrase-based systems have difficulties in modeling permutations among phrases. As a result, it is important to develop effective reordering models to capture such non-local reordering.

The early phrase-based paradigm (Koehn et al., 2003) applies a simple distance-based distortion penalty to model the phrase movements. More recently, many researchers have presented lexicalized reordering models that take advantage of lexical information to predict reordering (Tillmann, 2004; Xiong et al., 2006; Zens and Ney, 2006; Koehn et

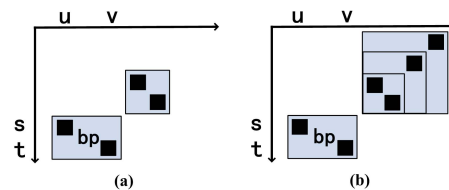


Figure 1: Occurrence of a swap with different numbers of adjacent bilingual phrases: only one phrase in (a) and three phrases in (b). Black squares denote word alignments and gray rectangles denote bilingual phrases. $[s,t]$ indicates the target-side span of bilingual phrase bp and $[u,v]$ represents the source-side span of bilingual phrase bp .

al., 2007; Galley and Manning, 2008). These models are learned from a word-aligned corpus to predict three orientations of a phrase pair with respect to the previous bilingual phrase: monotone (M), swap (S), and discontinuous (D). Take the bilingual phrase bp in Figure 1(a) for example. The word-based reordering model (Koehn et al., 2007) analyzes the word alignments at positions $(s-1, u-1)$ and $(s-1, v+1)$. The orientation of bp is set to D because the position $(s-1, v+1)$ contains no word alignment. The phrase-based reordering model (Tillmann, 2004) determines the presence of the adjacent bilingual phrase located in position $(s-1, v+1)$ and then treats the orientation of bp as S . Given no constraint on maximum phrase length, the hierarchical phrase reordering model (Galley and Manning, 2008) also analyzes the adjacent bilingual phrases for bp and identifies its orientation as S .

However, given a bilingual phrase, the above-mentioned models just consider the presence of an adjacent bilingual phrase rather than the number of adjacent bilingual phrases. See the examples in Fig-

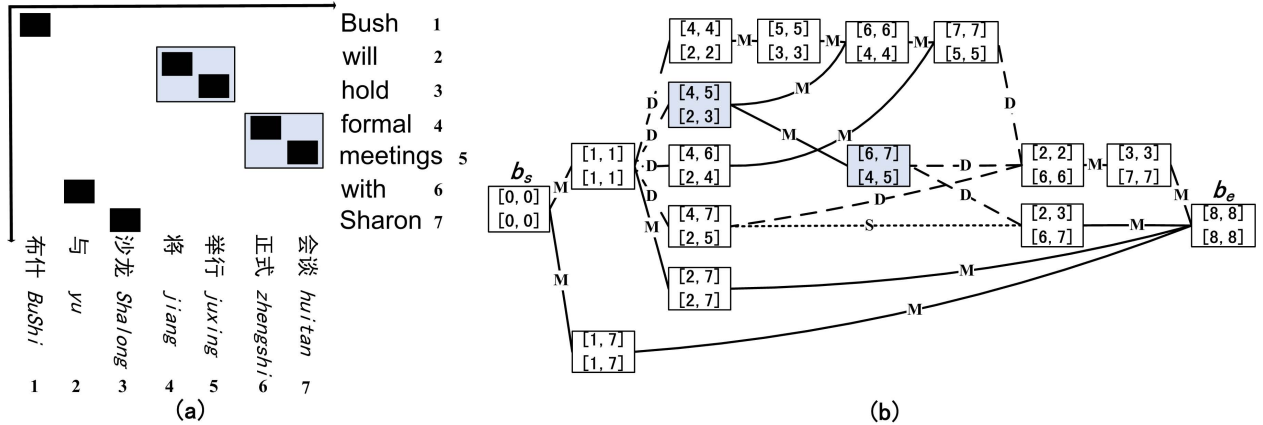


Figure 2: (a) A parallel Chinese-English sentence pair and (b) its corresponding reordering graph. In (b), we denote each bilingual phrase with a rectangle, where the upper and bottom numbers in the brackets represent the source and target spans of this bilingual phrase respectively. M = monotone (solid lines), S = swap (dotted line), and D = discontinuous (segmented lines). The bilingual phrases marked in the gray constitute a reordering example.

ure 1 for illustration. In Figure 1(a), bp is in a swap order with only one bilingual phrase. In Figure 1(b), bp swaps with three bilingual phrases. Lexicalized reordering models do not distinguish different numbers of adjacent phrase pairs, and just give bp the same count in the swap orientation.

In this paper, we propose a novel method to better estimate the reordering probabilities with the consideration of varying numbers of adjacent bilingual phrases. Our method uses reordering graphs to represent all phrase segmentations of parallel sentence pairs, and then gets the fractional counts of bilingual phrases for orientations from reordering graphs in an inside-outside fashion. Experimental results indicate that our method achieves significant improvements over the traditional lexicalized reordering model (Koehn et al., 2007).

This paper is organized as follows: in Section 2, we first give a brief introduction to the traditional lexicalized reordering model. Then we introduce our method to estimate the reordering probabilities from reordering graphs. The experimental results are reported in Section 3. Finally, we end with a conclusion and future work in Section 4.

2 Estimation of Reordering Probabilities Based on Reordering Graph

In this section, we first describe the traditional lexicalized reordering model, and then illustrate how to construct reordering graphs to estimate the reorder-

ing probabilities.

2.1 Lexicalized Reordering Model

Given a phrase pair $bp = (\bar{e}_i, \bar{f}_{a_i})$, where a_i defines that the source phrase \bar{f}_{a_i} is aligned to the target phrase \bar{e}_i , the traditional lexicalized reordering model computes the reordering count of bp in the orientation o based on the word alignments of boundary words. Specifically, the model collects bilingual phrases and distinguishes their orientations with respect to the previous bilingual phrase into three categories:

$$o = \begin{cases} M & a_i - a_{i-1} = 1 \\ S & a_i - a_{i-1} = -1 \\ D & |a_i - a_{i-1}| \neq 1 \end{cases} \quad (1)$$

Using the relative-frequency approach, the reordering probability regarding bp is

$$p(o|bp) = \frac{Count(o, bp)}{\sum_{o'} Count(o', bp)} \quad (2)$$

2.2 Reordering Graph

For a parallel sentence pair, its reordering graph indicates all possible translation derivations consisting of the extracted bilingual phrases. To construct a reordering graph, we first extract bilingual phrases using the way of (Och, 2003). Then, the adjacent

bilingual phrases are linked according to the target-side order. Some bilingual phrases, which have no adjacent bilingual phrases because of maximum length limitation, are linked to the nearest bilingual phrases in the target-side order.

Shown in Figure 2(b), the reordering graph for the parallel sentence pair (Figure 2(a)) can be represented as an undirected graph, where each rectangle corresponds to a phrase pair, each link is the orientation relationship between adjacent bilingual phrases, and two distinguished rectangles b_s and b_e indicate the beginning and ending of the parallel sentence pair, respectively. With the reordering graph, we can obtain all reordering examples containing the given bilingual phrase. For example, the bilingual phrase $\langle zhengshi\ huitan, \text{ formal meetings} \rangle$ (see Figure 2(a)), corresponding to the rectangle labeled with the source span [6,7] and the target span [4,5], is in a monotone order with one previous phrase and in a discontinuous order with two subsequent phrases (see Figure 2(b)).

2.3 Estimation of Reordering Probabilities

We estimate the reordering probabilities from reordering graphs. Given a parallel sentence pair, there are many translation derivations corresponding to different paths in its reordering graph. Assuming all derivations have a uniform probability, the fractional counts of bilingual phrases for orientations can be calculated by utilizing an algorithm in the inside-outside fashion.

Given a phrase pair bp in the reordering graph, we denote the number of paths from b_s to bp with $\alpha(bp)$. It can be computed in an iterative way $\alpha(bp) = \sum_{bp'} \alpha(bp')$, where bp' is one of the previous bilingual phrases of bp and $\alpha(b_s)=1$. In a similar way, the number of paths from b_e to bp , notated as $\beta(bp)$, is simply $\beta(bp) = \sum_{bp''} \beta(bp'')$, where bp'' is one of the subsequent bilingual phrases of bp and $\beta(b_e)=1$. Here, we show the α and β values of all bilingual phrases of Figure 2 in Table 1. Especially, for the reordering example consisting of the bilingual phrases $bp_1=\langle jiang\ juxing, \text{ will hold} \rangle$ and $bp_2=\langle zhengshi\ huitan, \text{ formal meetings} \rangle$, marked in the gray color in Figure 2, the α and β values can be calculated: $\alpha(bp_1) = 1$, $\beta(bp_2) = 1+1 = 2$, $\beta(b_s) = 8+1 = 9$.

Inspired by the parsing literature on pruning

<i>src span</i>	<i>trg span</i>	α	β
[0, 0]	[0, 0]	1	9
[1, 1]	[1, 1]	1	8
[1, 7]	[1, 7]	1	1
[4, 4]	[2, 2]	1	1
[4, 5]	[2, 3]	1	3
[4, 6]	[2, 4]	1	1
[4, 7]	[2, 5]	1	2
[2, 7]	[2, 7]	1	1
[5, 5]	[3, 3]	1	1
[6, 6]	[4, 4]	2	1
[6, 7]	[4, 5]	1	2
[7, 7]	[5, 5]	3	1
[2, 2]	[6, 6]	5	1
[2, 3]	[6, 7]	2	1
[3, 3]	[7, 7]	5	1
[8, 8]	[8, 8]	9	1

Table 1: The α and β values of the bilingual phrases shown in Figure 2.

(Charniak and Johnson, 2005; Huang, 2008), the fractional count of (o, bp', bp) is

$$Count(o, bp', bp) = \frac{\alpha(bp') \cdot \beta(bp)}{\beta(b_s)} \quad (3)$$

where the numerator indicates the number of paths containing the reordering example (o, bp', bp) and the denominator is the total number of paths in the reordering graph. Continuing with the reordering example described above, we obtain its fractional count using the formula (3): $Count(M, bp_1, bp_2) = (1 \times 2)/9 = 2/9$.

Then, the fractional count of bp in the orientation o is calculated as described below:

$$Count(o, bp) = \sum_{bp'} Count(o, bp', bp) \quad (4)$$

For example, we compute the fractional count of bp_2 in the monotone orientation by the formula (4): $Count(M, bp_2) = 2/9$.

As described in the lexicalized reordering model (Section 2.1), we apply the formula (2) to calculate the final reordering probabilities.

3 Experiments

We conduct experiments to investigate the effectiveness of our method on the **msd-fe** reordering model and the **msd-bidirectional-fe** reordering model. These two models are widely applied in

phrase-based system (Koehn et al., 2007). The msd-fe reordering model has three features, which represent the probabilities of bilingual phrases in three orientations: monotone, swap, or discontinuous. If a msd-bidirectional-fe model is used, then the number of features doubles: one for each direction.

3.1 Experiment Setup

Two different sizes of training corpora are used in our experiments: one is a small-scale corpus that comes from FBIS corpus consisting of 239K bilingual sentence pairs, the other is a large-scale corpus that includes 1.55M bilingual sentence pairs from LDC. The 2002 NIST MT evaluation test data is used as the development set and the 2003, 2004, 2005 NIST MT test data are the test sets. We choose the MOSES¹ (Koehn et al., 2007) as the experimental decoder. GIZA++ (Och and Ney, 2003) and the heuristics “grow-diag-final-and” are used to generate a word-aligned corpus, where we extract bilingual phrases with maximum length 7. We use SRILM Toolkits (Stolcke, 2002) to train a 4-gram language model on the Xinhua portion of Gigaword corpus.

In exception to the reordering probabilities, we use the same features in the comparative experiments. During decoding, we set $t_{table-limit} = 20$, $stack = 100$, and perform minimum-error-rate training (Och, 2003) to tune various feature weights. The translation quality is evaluated by case-insensitive BLEU-4 metric (Papineni et al., 2002). Finally, we conduct paired bootstrap sampling (Koehn, 2004) to test the significance in BLEU scores differences.

3.2 Experimental Results

Table 2 shows the results of experiments with the small training corpus. For the msd-fe model, the BLEU scores by our method are 30.51 32.78 and 29.50, achieving absolute improvements of **0.89**, **0.66** and **0.62** on the three test sets, respectively. For the msd-bidirectional-fe model, our method obtains BLEU scores of 30.49 32.73 and 29.24, with absolute improvements of **1.11**, **0.73** and **0.60** over the baseline method.

¹The phrase-based lexical reordering model (Tillmann, 2004) is also closely related to our model. However, due to the limit of time and space, we only use Moses-style reordering model (Koehn et al., 2007) as our baseline.

model	method	MT-03	MT-04	MT-05
m-f	baseline	29.62	32.12	28.88
	RG	30.51**	32.78**	29.50*
m-b-f	baseline	29.38	32.00	28.64
	RG	30.49**	32.73**	29.24*

Table 2: Experimental results with the **small-scale** corpus. m-f: msd-fe reordering model. m-b-f: msd-bidirectional-fe reordering model. RG: probabilities estimation based on Reordering Graph. * or **: significantly better than baseline ($p < 0.05$ or $p < 0.01$).

model	method	MT-03	MT-04	MT-05
m-f	baseline	31.58	32.39	31.49
	RG	32.44**	33.24**	31.64
m-b-f	baseline	32.43	33.07	31.69
	RG	33.29**	34.49**	32.79**

Table 3: Experimental results with the **large-scale** corpus.

Table 3 shows the results of experiments with the large training corpus. In the experiments of the msd-fe model, in exception to the MT-05 test set, our method is superior to the baseline method. The BLEU scores by our method are 32.44, 33.24 and 31.64, which obtain **0.86**, **0.85** and **0.15** gains on three test set, respectively. For the msd-bidirectional-fe model, the BLEU scores produced by our approach are 33.29, 34.49 and 32.79 on the three test sets, with **0.86**, **1.42** and **1.1** points higher than the baseline method, respectively.

4 Conclusion and Future Work

In this paper, we propose a method to improve the reordering model by considering the effect of the number of adjacent bilingual phrases on the reordering probabilities estimation. Experimental results on NIST Chinese-to-English tasks demonstrate the effectiveness of our method.

Our method is also general to other lexicalized reordering models. We plan to apply our method to the complex lexicalized reordering models, for example, the hierarchical reordering model (Galley and Manning, 2008) and the MEBTG reordering model (Xiong et al., 2006). In addition, how to further improve the reordering model by distinguishing the derivations with different probabilities will become another study emphasis in further research.

Acknowledgement

The authors were supported by National Natural Science Foundation of China, Contracts 60873167 and 60903138. We thank the anonymous reviewers for their insightful comments. We are also grateful to Hongmei Zhao and Shu Cai for their helpful feedback.

References

- Eugene Charniak and Mark Johnson. 2005. Coarse-to-fine n-best parsing and maxent discriminative reranking. In *Proc. of ACL 2005*, pages 173–180.
- Michel Galley and Christopher D. Manning. 2008. A simple and effective hierarchical phrase reordering model. In *Proc. of EMNLP 2008*, pages 848–856.
- Liang Huang. 2008. Forest reranking: Discriminative parsing with non-local features. In *Proc. of ACL 2008*, pages 586–594.
- Philipp Koehn, Franz Josef Och, and Daniel Marcu. 2003. Statistical phrase-based translation. In *Proc. of HLT-NAACL 2003*, pages 127–133.
- Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondrej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In *Proc. of ACL 2007, Demonstration Session*, pages 177–180.
- Philipp Koehn. 2004. Statistical significance tests for machine translation evaluation. In *Proc. of EMNLP 2004*, pages 388–395.
- Franz Josef Och and Hermann Ney. 2003. A systematic comparison of various statistical alignment models. *Computational Linguistics*, 29(1):19–51.
- Franz Joseph Och and Hermann Ney. 2004. The alignment template approach to statistical machine translation. *Computational Linguistics*, pages 417–449.
- Franz Josef Och. 2003. Minimum error rate training in statistical machine translation. In *Proc. of ACL 2003*, pages 160–167.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proc. of ACL 2002*, pages 311–318.
- Andreas Stolcke. 2002. Srilm - an extensible language modeling toolkit. In *Proc. of ICSLP 2002*, pages 901–904.
- Christoph Tillmann. 2004. A unigram orientation model for statistical machine translation. In *Proc. of HLT-ACL 2004, Short Papers*, pages 101–104.
- Deyi Xiong, Qun Liu, and Shouxun Lin. 2006. Maximum entropy based phrase reordering model for statistical machine translation. In *Proc. of ACL 2006*, pages 521–528.
- Richard Zens and Hermann Ney. 2006. Discriminative reordering models for statistical machine translation. In *Proc. of Workshop on Statistical Machine Translation 2006*, pages 521–528.

Filtering Syntactic Constraints for Statistical Machine Translation

Hailong Cao and Eiichiro Sumita

Language Translation Group, MASTAR Project
National Institute of Information and Communications Technology
3-5 Hikari-dai, Seika-cho, Soraku-gun, Kyoto, Japan, 619-0289
{hlcao, eiichiro.sumita}@nict.go.jp

Abstract

Source language parse trees offer very useful but imperfect reordering constraints for statistical machine translation. A lot of effort has been made for soft applications of syntactic constraints. We alternatively propose the selective use of syntactic constraints. A classifier is built automatically to decide whether a node in the parse trees should be used as a reordering constraint or not. Using this information yields a 0.8 BLEU point improvement over a full constraint-based system.

1 Introduction

In statistical machine translation (SMT), the search problem is NP-hard if arbitrary reordering is allowed (Knight, 1999). Therefore, we need to restrict the possible reordering in an appropriate way for both efficiency and translation quality. The most widely used reordering constraints are IBM constraints (Berger et al., 1996), ITG constraints (Wu, 1995) and syntactic constraints (Yamada et al., 2000; Galley et al., 2004; Liu et al., 2006; Marcu et al., 2006; Zollmann and Venugopal 2006; and numerous others). Syntactic constraints can be imposed from the source side or target side. This work will focus on syntactic constraints from source parse trees.

Linguistic parse trees can provide very useful reordering constraints for SMT. However, they are far from perfect because of both parsing errors and the crossing of the constituents and formal phrases extracted from parallel training data. The key challenge is how to take advantage of the prior knowledge in the linguistic parse trees without affecting the strengths of formal phrases. Recent efforts attack this problem by using the constraints softly (Cherry, 2008; Marton and Resnik, 2008). In their methods, a candidate

translation gets an extra credit if it respects the parse tree but may incur a cost if it violates a constituent boundary.

In this paper, we address this challenge from a less explored direction. Rather than use all constraints offered by the parse trees, we propose using them selectively. Based on parallel training data, a classifier is built automatically to decide whether a node in the parse trees should be used as a reordering constraint or not. As a result, we obtain a 0.8 BLEU point improvement over a full constraint-based system.

2 Reordering Constraints from Source Parse Trees

In this section we briefly review a constraint-based system named IST-ITG (Imposing Source Tree on Inversion Transduction Grammar, Yamamoto et al., 2008) upon which this work builds.

When using ITG constraints during decoding, the source-side parse tree structure is not considered. The reordering process can be more tightly constrained if constraints from the source parse tree are integrated with the ITG constraints. IST-ITG constraints directly apply source sentence tree structure to generate the target with the following constraint: the target sentence is obtained by rotating any node of the source sentence tree structure.

After parsing the source sentence, a bracketed sentence is obtained by removing the node syntactic labels; this bracketed sentence can then be directly expressed as a tree structure. For example¹, the parse tree “(S1 (S (NP (DT This)) (VP (AUX is) (NP (DT a) (NN pen))))))” is obtained from the source sentence “This is a pen”, which consists of four words. By removing

¹ We use English examples for the sake of readability.

the node syntactic labels, the bracketed sentence “((This) ((is) ((a) (pen))))” is obtained. Such a bracketed sentence can be used to produce constraints.

For example, for the source-side bracketed tree “((f1 f2) (f3 f4))”, eight target sequences [e1, e2, e3, e4], [e2, e1, e3, e4], [e1, e2, e4, e3], [e2, e1, e4, e3], [e3, e4, e1, e2], [e3, e4, e2, e1], [e4, e3, e1, e2], and [e4, e3, e2, e1] are possible. For the source-side bracketed tree “(((f1f2) f3) f4),” eight sequences [e1, e2, e3, e4], [e2, e1, e3, e4], [e3, e1, e2, e4], [e3, e2, e1, e4], [e4, e1, e2, e3], [e4, e2, e1, e3], [e4, e3, e1, e2], and [e4, e3, e2, e1] are possible. When the source sentence tree structure is a binary tree, the number of word orderings is reduced to 2^{N-1} where N is the length of the source sentence.

The parsing results sometimes do not produce binary trees. In this case, some subtrees have more than two child nodes. For a non-binary subtree, any reordering of child nodes is allowed. For example, if a subtree has three child nodes, six reorderings of the nodes are possible.

3 Learning to Classify Parse Tree Nodes

In IST-ITG and many other methods which use syntactic constraints, all of the nodes in the parse trees are utilized. Though many nodes in the parse trees are useful, we would argue that some nodes are not trustworthy. For example, if we constrain the translation of “f1 f2 f3 f4” with node N2 illustrated in Figure 1, then word “e1” will never be put in the middle the other three words. If we want to obtain the translation “e2 e1 e4 e3”, node N3 can offer a good constraint while node N2 should be filtered out. In real corpora, cases such as node N2 are frequent enough to be noticeable (see Fox (2002) or section 4.1 in this paper).

Therefore, we use the definitions in Galley et al. (2004) to classify the nodes in parse trees into two types: frontier nodes and interior nodes. Though the definitions were originally made for target language parse trees, they can be straightforwardly applied to the source side. A node which satisfies both of the following two conditions is referred as a frontier node:

- All the words covered by the node can be translated separately. That is to say, these words do not share a translation with any word outside the coverage of the node.

- All the words covered by the node remain contiguous after translation.

Otherwise the node is an interior node.

For example, in Figure 1, both node N1 and node N3 are frontier nodes. Node N2 is an interior node because the source words f2, f3 and f4 are translated into e2, e3 and e4, which are not contiguous in the target side.

Clearly, only frontier nodes should be used as reordering constraints while interior nodes are not suitable for this. However, little work has been done on how to explicitly distinguish these two kinds of nodes in the source parse trees. In this section, we will explore building a classifier which can label the nodes in the parse trees as frontier nodes or interior nodes.

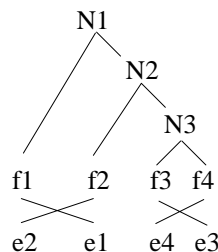


Figure 1: An example parse tree and alignments

3.1 Training

Ideally, we would have a human-annotated corpus in which each sentence is parsed and each node in the parse trees is labeled as a frontier node or an interior node. But such a target language specific corpus is hard to come by, and never in the quantity we would like.

Instead, we generate such a corpus automatically. We begin with a parallel corpus which will be used to train our SMT model. In our case, it is the FBIS Chinese-English corpus.

Firstly, the Chinese sentences are segmented, POS tagged and parsed by the tools described in Kruengkrai et al. (2009) and Cao et al. (2007), both of which are trained on the Penn Chinese Treebank 6.0.

Secondly, we use GIZA++ to align the sentences in both the Chinese-English and English-Chinese directions. We combine the alignments using the “grow-diag-final-and” procedure provided with MOSES (Koehn, 2007). Because there are many errors in the alignment, we remove the links if the alignment count is less than three for the source or the target word. Additionally, we also remove notoriously bad links in

{de, le} × {the, a, an} following Fossum and Knight (2008).

Thirdly, given the parse trees and the alignment information, we label each node as a frontier node or an interior node according to the definition introduced in this section. Using the labeled nodes as training data, we can build a classifier. In theory, a broad class of machine learning tools can be used; however, due to the scale of the task (see section 4), we utilize the Pegasos² which is a very fast SVM solver (Shalev-Shwartz et al, 2007).

3.2 Features

For each node in the parse trees, we use the following feature templates:

- A context-free grammar rule which rewrites the current node (In this and all the following grammar based features, a mark is used to indicate which non terminal is the current node.)
- A context-free grammar rule which rewrites the current node’s father
- The combination of the above two rules
- A lexicalized context-free grammar rule which rewrites the current node
- A lexicalized context-free grammar rule which rewrites the current node’s father
- Syntactic label, head word, and head POS tag of the current node
- Syntactic label, head word, and head POS tag of the current node’s left child
- Syntactic label, head word, and head POS tag of the current node’s right child
- Syntactic label, head word, and head POS tag of the current node’s left brother
- Syntactic label, head word, and head POS tag of the current node’s right brother
- Syntactic label, head word, and head POS tag of the current node’s father
- The leftmost word covered by the current node and the word before it
- The rightmost word covered by the current node and the word after it

4 Experiments

Our SMT system is based on a fairly typical phrase-based model (Finch and Sumita, 2008). For the training of our SMT model, we use a modified training toolkit adapted from the

² <http://www.cs.huji.ac.il/~shais/code/index.html>

MOSES decoder. Our decoder can operate on the same principles as the MOSES decoder. Minimum error rate training (MERT) with respect to BLEU score is used to tune the decoder’s parameters, and it is performed using the standard technique of Och (2003). A lexical reordering model was used in our experiments.

The translation model was created from the FBIS corpus. We used a 5-gram language model trained with modified Knesser-Ney smoothing. The language model was trained on the target side of FBIS corpus and the Xinhua news in GIGAWORD corpus. The development and test sets are from NIST MT08 evaluation campaign. Table 1 shows the statistics of the corpora used in our experiments.

Data	Sentences	Chinese words	English words
Training set	243,698	7,933,133	10,343,140
Development set	1664	38,779	46,387
Test set	1357	32377	42,444
GIGAWORD	19,049,757	-	306,221,306

Table 1: Corpora statistics

4.1 Experiments on Nodes Classification

We extracted about 3.9 million example nodes from the training data, i.e. the FBIS corpus. There were 2.37 million frontier nodes and 1.59 million interior nodes in these examples, give rise to about 4.4 million features. To test the performance of our classifier, we simply use the last ten thousand examples as a test set, and the rest being used as Pegasos training data. All the parameters in Pegasos were set as default values. In this way, the accuracy of the classifier was 71.59%.

Then we retrained our classifier by using all of the examples. The nodes in the automatically parsed NIST MT08 test set were labeled by the classifier. As a result, 17,240 nodes were labeled as frontier nodes and 5,736 nodes were labeled as interior nodes.

4.2 Experiments on Chinese-English SMT

In order to confirm that it is advantageous to distinguish between frontier nodes and interior nodes, we performed four translation experiments.

The first one was a typical beam search decoding without any syntactic constraints.

All the other three experiments were based on the IST-ITG method which makes use of syntac-

tic constraints. The difference between these three experiments lies in what constraints are used. In detail, the second one used all nodes recognized by the parser; the third one only used frontier nodes labeled by the classifier; the fourth one only used interior nodes labeled by the classifier.

With the exception of the above differences, all the other settings were the same in the four experiments. Table 2 summarizes the SMT performance.

Syntactic Constraints	BLEU
none	17.26
all nodes	16.83
frontier nodes	17.63
interior nodes	16.59

Table 2: Comparison of different constraints by SMT quality

Clearly, we obtain the best performance if we constrain the search with only frontier nodes. Using just frontier yields a 0.8 BLEU point improvement over the baseline constraint-based system which uses all the constraints.

On the other hand, constraints from interior nodes result in the worst performance. This comparison shows it is necessary to explicitly distinguish nodes in the source parse trees when they are used as reordering constraints.

The improvement over the system without constraints is only modest. It may be too coarse to use parse trees as hard constraints. We believe a greater improvement can be expected if we apply our idea to finer-grained approaches that use constraints softly (Marton and Resnik (2008) and Cherry (2008)).

5 Conclusion and Future Work

We propose a selectively approach to syntactic constraints during decoding. A classifier is built automatically to decide whether a node in the parse trees should be used as a reordering constraint or not. Preliminary results show that it is not only advantageous but necessary to explicitly distinguish between frontier nodes and interior nodes.

The idea of selecting syntactic constraints is compatible with the idea of using constraints softly; we plan to combine the two ideas and obtain further improvements in future work.

Acknowledgments

We would like to thank Taro Watanabe and Andrew Finch for insightful discussions. We also would like to thank the anonymous reviewers for their constructive comments.

Reference

- A.L. Berger, P.F. Brown, S.A.D. Pietra, V.J.D. Pietra, J.R. Gillett, A.S. Kehler, and R.L. Mercer. 1996. Language translation apparatus and method of using context-based translation models. United States patent, patent number 5510981, April.
- Hailong Cao, Yujie Zhang and Hitoshi Isahara. Empirical study on parsing Chinese based on Collins' model. 2007. In *PACLING*.
- Colin Cherry. 2008. Cohesive phrase-Based decoding for statistical machine translation. In *ACL- HLT*.
- Andrew Finch and Eiichiro Sumita. 2008. Dynamic model interpolation for statistical machine translation. In *SMT Workshop*.
- Victoria Fossum and Kevin Knight. 2008. Using bilingual Chinese-English word alignments to resolve PP attachment ambiguity in English. In *AMTA Student Workshop*.
- Heidi J. Fox. 2002. Phrasal cohesion and statistical machine translation. In *EMNLP*.
- Michel Galley, Mark Hopkins, Kevin Knight, and Daniel Marcu. 2004. What's in a translation rule? In *HLT-NAACL*.
- Kevin Knight. 1999. Decoding complexity in word replacement translation models. *Computational Linguistics*, 25(4):607–615.
- Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondrej Bojar, Alexandra Constantin, Evan Herbst. 2007. Moses: Open Source Toolkit for Statistical Machine Translation. In *ACL demo and poster sessions*.
- Canasai Kruengkrai, Kiyotaka Uchimoto, Jun'ichi Kazama, Yiou Wang, Kentaro Torisawa and Hitoshi Isahara. 2009. An error-driven word-character hybrid model for joint Chinese word segmentation and POS tagging. In *ACL-IJCNLP*.
- Yang Liu, Qun Liu, Shouxun Lin. 2006. Tree-to-string alignment template for statistical machine translation. In *ACL-COLING*.
- Daniel Marcu, Wei Wang, Abdessamad Echihabi, and Kevin Knight. 2006. SPMT: Statistical machine translation with syntactified target language phrases. In *EMNLP*.

- Yuval Marton and Philip Resnik. 2008. Soft syntactic constraints for hierarchical phrased-based translation. In *ACL-HLT*.
- Franz Och. 2003. Minimum error rate training in statistical machine translation. In *ACL*.
- Shai Shalev-Shwartz, Yoram Singer and Nathan Srebro. 2007. Pegasos: Primal estimated sub-gradient solver for SVM. In *ICML*.
- Dekai Wu. 1995. Stochastic inversion transduction grammars with application to segmentation, bracketing, and alignment of parallel corpora. In *IJCAI*.
- Kenji Yamada and Kevin Knight. 2000. A syntax-based statistical translation model. In *ACL*.
- Hirofumi Yamamoto, Hideo Okuma and Eiichiro Sumita. 2008. Imposing constraints from the source tree on ITG constraints for SMT. In *Workshop on syntax and structure in statistical translation*.
- Andreas Zollmann and Ashish Venugopal. 2006. Syntax augmented machine translation via chart parsing. In *SMT Workshop, HLT-NAACL*.

Diversify and Combine: Improving Word Alignment for Machine Translation on Low-Resource Languages

Bing Xiang, Yonggang Deng, and Bowen Zhou

IBM T. J. Watson Research Center

Yorktown Heights, NY 10598

{bxiang, ydeng, zhou}@us.ibm.com

Abstract

We present a novel method to improve word alignment quality and eventually the translation performance by producing and combining complementary word alignments for low-resource languages. Instead of focusing on the improvement of a single set of word alignments, we generate multiple sets of diversified alignments based on different motivations, such as linguistic knowledge, morphology and heuristics. We demonstrate this approach on an English-to-Pashto translation task by combining the alignments obtained from syntactic reordering, stemming, and partial words. The combined alignment outperforms the baseline alignment, with significantly higher F-scores and better translation performance.

1 Introduction

Word alignment usually serves as the starting point and foundation for a statistical machine translation (SMT) system. It has received a significant amount of research over the years, notably in (Brown et al., 1993; Ittycheriah and Roukos, 2005; Fraser and Marcu, 2007; Hermjakob, 2009). They all focused on the improvement of word alignment models. In this work, we leverage existing aligners and generate multiple sets of word alignments based on complementary information, then combine them to get the final alignment for phrase training. The resource required for this approach is little, compared to what is needed to build a reasonable discriminative alignment model, for example. This makes the approach especially appealing for SMT on low-resource languages.

Most of the research on alignment combination in the past has focused on how to combine the alignments from two different directions, source-to-target and target-to-source. Usually people start from the intersection of two sets of alignments, and gradually add links in the union based on certain heuristics, as in (Koehn et al., 2003), to achieve a better balance compared to using either intersection (high precision) or union (high recall). In (Ayan and Dorr, 2006) a maximum entropy approach was proposed to combine multiple alignments based on a set of linguistic and alignment features. A different approach was presented in (Deng and Zhou, 2009), which again concentrated on the combination of two sets of alignments, but with a different criterion. It tries to maximize the number of phrases that can be extracted in the combined alignments. A greedy search method was utilized and it achieved higher translation performance than the baseline.

More recently, an alignment selection approach was proposed in (Huang, 2009), which computes confidence scores for each link and prunes the links from multiple sets of alignments using a hand-picked threshold. The alignments used in that work were generated from different aligners (HMM, block model, and maximum entropy model). In this work, we use soft voting with weighted confidence scores, where the weights can be tuned with a specific objective function. There is no need for a pre-determined threshold as used in (Huang, 2009). Also, we utilize various knowledge sources to enrich the alignments instead of using different aligners. Our strategy is to diversify and then combine in order to catch any complementary information captured in the word alignments for low-resource languages.

The rest of the paper is organized as follows.

We present three different sets of alignments in Section 2 for an English-to-Pashto MT task. In Section 3, we propose the alignment combination algorithm. The experimental results are reported in Section 4. We conclude the paper in Section 5.

2 Diversified Word Alignments

We take an English-to-Pashto MT task as an example and create three sets of additional alignments on top of the baseline alignment.

2.1 Syntactic Reordering

Pashto is a subject-object-verb (SOV) language, which puts verbs after objects. People have proposed different syntactic rules to pre-reorder SOV languages, either based on a constituent parse tree (Drábek and Yarowsky, 2004; Wang et al., 2007) or dependency parse tree (Xu et al., 2009). In this work, we apply syntactic reordering for verb phrases (VP) based on the English constituent parse. The VP-based reordering rule we apply in the work is:

- $VP(VB^*, *) \rightarrow VP(*, VB^*)$

where VB^* represents VB, VBD, VBG, VBN, VBP and VBZ .

In Figure 1, we show the reference alignment between an English sentence and the corresponding Pashto translation, where E is the original English sentence, P is the Pashto sentence (in romanized text), and E' is the English sentence after reordering. As we can see, after the VP-based reordering, the alignment between the two sentences becomes monotone, which makes it easier for the aligner to get the alignment correct. During the reordering of English sentences, we store the index changes for the English words. After getting the alignment trained on the reordered English and original Pashto sentence pairs, we map the English words back to the original order, along with the learned alignment links. In this way, the alignment is ready to be combined with the baseline alignment and any other alternatives.

2.2 Stemming

Pashto is one of the morphologically rich languages. In addition to the linguistic knowledge applied in the syntactic reordering described above, we also utilize morphological analysis by applying stemming on both the English and Pashto sides. For English, we use Porter stemming (Porter,

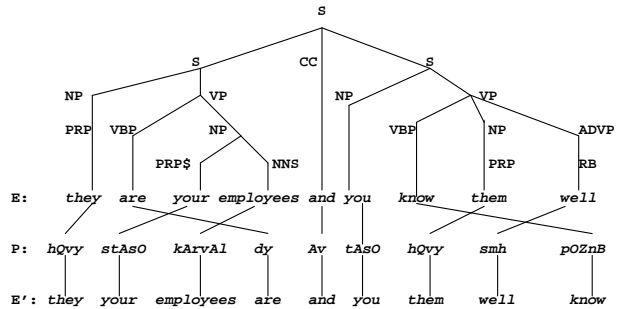


Figure 1: Alignment before/after VP-based reordering.

1980), a widely applied algorithm to remove the common morphological and inflexional endings from words in English. For Pashto, we utilize a morphological decomposition algorithm that has been shown to be effective for Arabic speech recognition (Xiang et al., 2006). We start from a fixed set of affixes with 8 prefixes and 21 suffixes. The prefixes and suffixes are stripped off from the Pashto words under the two constraints:(1) Longest matched affixes first; (2) Remaining stem must be at least two characters long.

2.3 Partial Word

For low-resource languages, we usually suffer from the data sparsity issue. Recently, a simple method was presented in (Chiang et al., 2009), which keeps partial English and Urdu words in the training data for alignment training. This is similar to the stemming method, but is more heuristics-based, and does not rely on a set of available affixes. With the same motivation, we keep the first 4 characters of each English and Pashto word to generate one more alternative for the word alignment.

3 Confidence-Based Alignment Combination

Now we describe the algorithm to combine multiple sets of word alignments based on weighted confidence scores. Suppose a_{ijk} is an alignment link in the i -th set of alignments between the j -th source word and k -th target word in sentence pair (S, T) . Similar to (Huang, 2009), we define the confidence of a_{ijk} as

$$c(a_{ijk}|S, T) = \sqrt{q_{s2t}(a_{ijk}|S, T)q_{t2s}(a_{ijk}|T, S)}, \quad (1)$$

where the source-to-target link posterior probability

$$q_{s2t}(a_{ijk}|S, T) = \frac{p_i(t_k|s_j)}{\sum_{k'=1}^K p_i(t_{k'}|s_j)}, \quad (2)$$

and the target-to-source link posterior probability $q_{t2s}(a_{ijk}|T, S)$ is defined similarly. $p_i(t_k|s_j)$ is the lexical translation probability between source word s_j and target word t_k in the i -th set of alignments.

Our alignment combination algorithm is as follows.

1. Each candidate link a_{ijk} gets soft votes from N sets of alignments via weighted confidence scores:

$$v(a_{ijk}|S, T) = \sum_{i=1}^N w_i * c(a_{ijk}|S, T), \quad (3)$$

where the weight w_i for each set of alignment can be optimized under various criteria. In this work, we tune it on a hand-aligned development set to maximize the alignment F-score.

2. All candidates are sorted by soft votes in descending order and evaluated sequentially. A candidate link a_{ijk} is included if one of the following is true:
 - Neither s_j nor t_k is aligned so far;
 - s_j is not aligned and its left or right neighboring word is aligned to t_k so far;
 - t_k is not aligned and its left or right neighboring word is aligned to s_j so far.
3. Repeat scanning all candidate links until no more links can be added.

In this way, those alignment links with higher confidence scores have higher priority to be included in the combined alignment.

4 Experiments

4.1 Baseline

Our training data contains around 70K English-Pashto sentence pairs released under the DARPA TRANSTAC project, with about 900K words on the English side. The baseline is a phrase-based MT system similar to (Koehn et al., 2003). We use GIZA++ (Och and Ney, 2000) to generate the baseline alignment for each direction and then

apply grow-diagonal-final (*gdf*). The decoding weights are optimized with minimum error rate training (MERT) (Och, 2003) to maximize BLEU scores (Papineni et al., 2002). There are 2028 sentences in the tuning set and 1019 sentences in the test set, both with one reference. We use another 150 sentence pairs as a heldout hand-aligned set to measure the word alignment quality. The three sets of alignments described in Section 2 are generated on the same training data separately with GIZA++ and enhanced by *gdf* as for the baseline alignment. The English parse tree used for the syntactic reordering was produced by a maximum entropy based parser (Ratnaparkhi, 1997).

4.2 Improvement in Word Alignment

In Table 1 we show the precision, recall and F-score of each set of word alignments for the 150-sentence set. Using partial word provides the highest F-score among all individual alignments. The F-score is 5% higher than for the baseline alignment. The VP-based reordering itself does not improve the F-score, which could be due to the parse errors on the conversational training data. We experiment with three options (c_0, c_1, c_2) when combining the baseline and reordering-based alignments. In c_0 , the weights w_i and confidence scores $c(a_{ijk}|S, T)$ in Eq. (3) are all set to 1. In c_1 , we set confidence scores to 1, while tuning the weights with hill climbing to maximize the F-score on a hand-aligned tuning set. In c_2 , we compute the confidence scores as in Eq. (1) and tune the weights as in c_1 . The numbers in Table 1 show the effectiveness of having both weights and confidence scores during the combination.

Similarly, we combine the baseline with each of the other sets of alignments using c_2 . They all result in significantly higher F-scores. We also generate alignments on VP-reordered partial words (X in Table 1) and compared $B + X$ and $B + V + P$. The better results with $B + V + P$ show the benefit of keeping the alignments as diversified as possible before the combination. Finally, we compare the proposed alignment combination c_2 with the heuristics-based method (*gdf*), where the latter starts from the intersection of all 4 sets of alignments and then applies grow-diagonal-final (Koehn et al., 2003) based on the links in the union. The proposed combination approach on $B + V + S + P$ results in close to 7% higher F-scores than the baseline and also 2% higher than

gdf. We also notice that its higher F-score is mainly due to the higher precision, which should result from the consideration of confidence scores.

Alignment	Comb	P	R	F
Baseline		0.6923	0.6414	0.6659
V		0.6934	0.6388	0.6650
S		0.7376	0.6495	0.6907
P		0.7665	0.6643	0.7118
X		0.7615	0.6641	0.7095
B+V	c_0	0.7639	0.6312	0.6913
B+V	c_1	0.7645	0.6373	0.6951
B+V	c_2	0.7895	0.6505	0.7133
B+S	c_2	0.7942	0.6553	0.7181
B+P	c_2	0.8006	0.6612	0.7242
B+X	c_2	0.7827	0.6670	0.7202
B+V+P	c_2	0.7912	0.6755	0.7288
B+V+S+P	<i>gdf</i>	0.7238	0.7042	0.7138
B+V+S+P	c_2	0.7906	0.6852	0.7342

Table 1: Alignment precision, recall and F-score (B: baseline; V: VP-based reordering; S: stemming; P: partial word; X: VP-reordered partial word).

4.3 Improvement in MT Performance

In Table 2, we show the corresponding BLEU scores on the test set for the systems built on each set of word alignment in Table 1. Similar to the observation from Table 1, c_2 outperforms c_0 and c_1 , and $B + V + S + P$ with c_2 outperforms $B + V + S + P$ with *gdf*. We also ran one experiment in which we concatenated all 4 sets of alignments into one big set (shown as *cat*). Overall, the BLEU score with confidence-based combination was increased by 1 point compared to the baseline, 0.6 compared to *gdf*, and 0.7 compared to *cat*. All results are statistically significant with $p < 0.05$ using the sign-test described in (Collins et al., 2005).

5 Conclusions

In this work, we have presented a word alignment combination method that improves both the alignment quality and the translation performance. We generated multiple sets of diversified alignments based on linguistics, morphology, and heuristics, and demonstrated the effectiveness of combination on the English-to-Pashto translation task. We showed that the combined alignment significantly outperforms the baseline alignment with

Alignment	Comb	Links	Phrase	BLEU
Baseline		963K	565K	12.67
V		965K	624K	12.82
S		915K	692K	13.04
P		906K	716K	13.30
X		911K	689K	13.00
B+V	c_0	870K	890K	13.20
B+V	c_1	865K	899K	13.32
B+V	c_2	874K	879K	13.60
B+S	c_2	864K	948K	13.41
B+P	c_2	863K	942K	13.40
B+X	c_2	871K	905K	13.37
B+V+P	c_2	880K	914K	13.60
B+V+S+P	<i>cat</i>	3749K	1258K	13.01
B+V+S+P	<i>gdf</i>	1021K	653K	13.14
B+V+S+P	c_2	907K	771K	13.73

Table 2: Improvement in BLEU scores (B: baseline; V: VP-based reordering; S: stemming; P: partial word; X: VP-reordered partial word).

both higher F-score and higher BLEU score. The combination approach itself is not limited to any specific alignment. It provides a general framework that can take advantage of as many alignments as possible, which could differ in preprocessing, alignment modeling, or any other aspect.

Acknowledgments

This work was supported by the DARPA TRANSTAC program. We would like to thank Upendra Chaudhari, Sameer Maskey and Xiaoqiang Luo for providing useful resources and the anonymous reviewers for their constructive comments.

References

- Necip Fazil Ayan and Bonnie J. Dorr. 2006. A maximum entropy approach to combining word alignments. In *Proc. HLT/NAACL*, June.
- Peter Brown, Vincent Della Pietra, Stephen Della Pietra, and Robert Mercer. 1993. The mathematics of statistical machine translation: parameter estimation. *Computational Linguistics*, 19(2):263–311.
- David Chiang, Kevin Knight, Samad Echiabi, et al. 2009. Isi/language weaver mist 2009 systems. In *Presentation at NIST MT 2009 Workshop*, August.
- Michael Collins, Philipp Koehn, and Ivona Kučerová. 2005. Clause restructuring for statistical machine translation. In *Proc. of ACL*, pages 531–540.

- Yonggang Deng and Bowen Zhou. 2009. Optimizing word alignment combination for phrase table training. In *Proc. ACL*, pages 229–232, August.
- Elliott Franco Drábek and David Yarowsky. 2004. Improving bitext word alignments via syntax-based reordering of english. In *Proc. ACL*.
- Alexander Fraser and Daniel Marcu. 2007. Getting the structure right for word alignment: Leaf. In *Proc. of EMNLP*, pages 51–60, June.
- Ulf Hermjakob. 2009. Improved word alignment with statistics and linguistic heuristics. In *Proc. EMNLP*, pages 229–237, August.
- Fei Huang. 2009. Confidence measure for word alignment. In *Proc. ACL*, pages 932–940, August.
- Abraham Ittycheriah and Salim Roukos. 2005. A maximum entropy word aligner for arabic-english machine translation. In *Proc. of HLT/EMNLP*, pages 89–96, October.
- Philipp Koehn, Franz Och, and Daniel Marcu. 2003. Statistical phrase-based translation. In *Proc. NAACL/HLT*.
- Franz Josef Och and Hermann Ney. 2000. Improved statistical alignment models. In *Proc. of ACL*, pages 440–447, Hong Kong, China, October.
- Franz Josef Och. 2003. Minimum error rate training in statistical machine translation. In *Proc. of ACL*, pages 160–167.
- Kishore Papineni, Salim Roukos, Todd Ward, and Weijing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proc. of ACL*, pages 311–318.
- Martin Porter. 1980. An algorithm for suffix stripping. In *Program*, volume 14, pages 130–137.
- Adwait Ratnaparkhi. 1997. A linear observed time statistical parser based on maximum entropy models. In *Proc. of EMNLP*, pages 1–10.
- Chao Wang, Michael Collins, and Philipp Koehn. 2007. Chinese syntactic reordering for statistical machine translation. In *Proc. EMNLP*, pages 737–745.
- Bing Xiang, Kham Nguyen, Long Nguyen, Richard Schwartz, and John Makhoul. 2006. Morphological decomposition for arabic broadcast news transcription. In *Proc. ICASSP*.
- Peng Xu, Jaeho Kang, Michael Ringgaard, and Franz Och. 2009. Using a dependency parser to improve smt for subject-object-verb languages. In *Proc. NAACL/HLT*, pages 245–253, June.

Efficient Path Counting Transducers for Minimum Bayes-Risk Decoding of Statistical Machine Translation Lattices

Graeme Blackwood, Adrià de Gispert, William Byrne

Machine Intelligence Laboratory
Cambridge University Engineering Department
Trumpington Street, CB2 1PZ, U.K.
{gwb24 | ad465 | wjb31}@cam.ac.uk

Abstract

This paper presents an efficient implementation of linearised lattice minimum Bayes-risk decoding using weighted finite state transducers. We introduce transducers to efficiently count lattice paths containing n -grams and use these to gather the required statistics. We show that these procedures can be implemented exactly through simple transformations of word sequences to sequences of n -grams. This yields a novel implementation of lattice minimum Bayes-risk decoding which is fast and exact even for very large lattices.

1 Introduction

This paper focuses on an exact implementation of the linearised form of lattice minimum Bayes-risk (LMBR) decoding using general purpose weighted finite state transducer (WFST) operations¹. The LMBR decision rule in Tromble et al. (2008) has the form

$$\hat{E} = \operatorname{argmax}_{E' \in \mathcal{E}} \left\{ \theta_0 |E'| + \sum_{u \in \mathcal{N}} \theta_u \#_u(E') p(u|\mathcal{E}) \right\} \quad (1)$$

where \mathcal{E} is a lattice of translation hypotheses, \mathcal{N} is the set of all n -grams in the lattice (typically, $n = 1 \dots 4$), and the parameters θ are constants estimated on held-out data. The quantity $p(u|\mathcal{E})$ we refer to as the path posterior probability of the n -gram u . This particular posterior is defined as

$$p(u|\mathcal{E}) = p(\mathcal{E}_u|\mathcal{E}) = \sum_{E \in \mathcal{E}_u} P(E|F), \quad (2)$$

where $\mathcal{E}_u = \{E \in \mathcal{E} : \#_u(E) > 0\}$ is the subset of lattice paths containing the n -gram u at least

¹We omit an introduction to WFSTs for space reasons. See Mohri et al. (2008) for details of the general purpose WFST operations used in this paper.

once. It is the efficient computation of these path posterior n -gram probabilities that is the primary focus of this paper. We will show how general purpose WFST algorithms can be employed to efficiently compute $p(u|\mathcal{E})$ for all $u \in \mathcal{N}$.

Tromble et al. (2008) use Equation (1) as an approximation to the general form of statistical machine translation MBR decoder (Kumar and Byrne, 2004):

$$\hat{E} = \operatorname{argmin}_{E' \in \mathcal{E}} \sum_{E \in \mathcal{E}} L(E, E') P(E|F) \quad (3)$$

The approximation replaces the sum over all paths in the lattice by a sum over lattice n -grams. Even though a lattice may have many n -grams, it is possible to extract and enumerate them exactly whereas this is often impossible for individual paths. Therefore, while the Tromble et al. (2008) linearisation of the gain function in the decision rule is an approximation, Equation (1) can be computed exactly even over very large lattices. The challenge is to do so efficiently.

If the quantity $p(u|\mathcal{E})$ had the form of a conditional expected count

$$c(u|\mathcal{E}) = \sum_{E \in \mathcal{E}} \#_u(E) P(E|F), \quad (4)$$

it could be computed efficiently using counting transducers (Allauzen et al., 2003). The statistic $c(u|\mathcal{E})$ counts the number of times an n -gram occurs on each path, accumulating the weighted count over all paths. By contrast, what is needed by the approximation in Equation (1) is to identify all paths containing an n -gram and accumulate their probabilities. The accumulation of probabilities at the path level, rather than the n -gram level, makes the exact computation of $p(u|\mathcal{E})$ hard.

Tromble et al. (2008) approach this problem by building a separate word sequence acceptor for each n -gram in \mathcal{N} and intersecting this acceptor

with the lattice to discard all paths that do not contain the n -gram; they then sum the probabilities of all paths in the filtered lattice. We refer to this as the *sequential method*, since $p(u|\mathcal{E})$ is calculated separately for each u in sequence.

Allauzen et al. (2010) introduce a transducer for simultaneous calculation of $p(u|\mathcal{E})$ for all unigrams $u \in \mathcal{N}_1$ in a lattice. This transducer is effective for finding path posterior probabilities of unigrams because there are relatively few unique unigrams in the lattice. As we will show, however, it is less efficient for higher-order n -grams.

Allauzen et al. (2010) use exact statistics for the unigram path posterior probabilities in Equation (1), but use the conditional expected counts of Equation (4) for higher-order n -grams. Their hybrid MBR decoder has the form

$$\begin{aligned} \hat{E} = & \operatorname{argmax}_{E' \in \mathcal{E}} \left\{ \theta_0 |E'| \right. \\ & + \sum_{u \in \mathcal{N}: 1 \leq |u| \leq k} \theta_u \#_u(E') p(u|\mathcal{E}) \\ & \left. + \sum_{u \in \mathcal{N}: k < |u| \leq 4} \theta_u \#_u(E') c(u|\mathcal{E}) \right\}, \quad (5) \end{aligned}$$

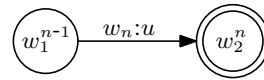
where k determines the range of n -gram orders at which the path posterior probabilities $p(u|\mathcal{E})$ of Equation (2) and conditional expected counts $c(u|\mathcal{E})$ of Equation (4) are used to compute the expected gain. For $k < 4$, Equation (5) is thus an approximation to the approximation. In many cases it will be perfectly fine, depending on how closely $p(u|\mathcal{E})$ and $c(u|\mathcal{E})$ agree for higher-order n -grams. Experimentally, Allauzen et al. (2010) find this approximation works well at $k = 1$ for MBR decoding of statistical machine translation lattices. However, there may be scenarios in which $p(u|\mathcal{E})$ and $c(u|\mathcal{E})$ differ so that Equation (5) is no longer useful in place of the original Tromble et al. (2008) approximation.

In the following sections, we present an efficient method for simultaneous calculation of $p(u|\mathcal{E})$ for n -grams of a fixed order. While other fast MBR approximations are possible (Kumar et al., 2009), we show how the exact path posterior probabilities can be calculated and applied in the implementation of Equation (1) for efficient MBR decoding over lattices.

2 N -gram Mapping Transducer

We make use of a trick to count higher-order n -grams. We build transducer Φ_n to map word se-

quences to n -gram sequences of order n . Φ_n has a similar form to the WFST implementation of an n -gram language model (Allauzen et al., 2003). Φ_n includes for each n -gram $u = w_1^n$ arcs of the form:



The n -gram lattice of order n is called \mathcal{E}_n and is found by composing $\mathcal{E} \circ \Phi_n$, projecting on the output, removing ϵ -arcs, determinizing, and minimising. The construction of \mathcal{E}_n is fast even for large lattices and is memory efficient. \mathcal{E}_n itself may have more states than \mathcal{E} due to the association of distinct n -gram histories with states. However, the counting transducer for unigrams is simpler than the corresponding counting transducer for higher-order n -grams. As a result, counting unigrams in \mathcal{E}_n is easier than counting n -grams in \mathcal{E} .

3 Efficient Path Counting

Associated with each \mathcal{E}_n we have a transducer Ψ_n which can be used to calculate the path posterior probabilities $p(u|\mathcal{E})$ for all $u \in \mathcal{N}_n$. In Figures 1 and 2 we give two possible forms² of Ψ_n that can be used to compute path posterior probabilities over n -grams $u_{1,2} \in \mathcal{N}_n$ for some n . No modification to the ρ -arc matching mechanism is required even in counting higher-order n -grams since all n -grams are represented as individual symbols after application of the mapping transducer Φ_n .

Transducer Ψ_n^L is used by Allauzen et al. (2010) to compute the exact unigram contribution to the conditional expected gain in Equation (5). For example, in counting paths that contain u_1 , Ψ_n^L retains the *first* occurrence of u_1 and maps every other symbol to ϵ . This ensures that in any path containing a given u , only the first u is counted, avoiding multiple counting of paths.

We introduce an alternative path counting transducer Ψ_n^R that effectively deletes all symbols except the *last* occurrence of u on any path by ensuring that any paths in composition which count earlier instances of u do not end in a final state. Multiple counting is avoided by counting only the last occurrence of each symbol u on a path.

We note that initial $\epsilon:\epsilon$ arcs in Ψ_n^L effectively create $|\mathcal{N}_n|$ copies of \mathcal{E}_n in composition while searching for the first occurrence of each u . Com-

²The special composition symbol σ matches any arc; ρ matches any arc other than those with an explicit transition. See the OpenFst documentation: <http://openfst.org>

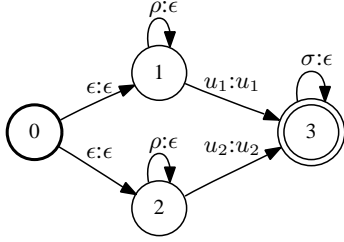


Figure 1: Path counting transducer Ψ_n^L matching first (left-most) occurrence of each $u \in \mathcal{N}_n$.

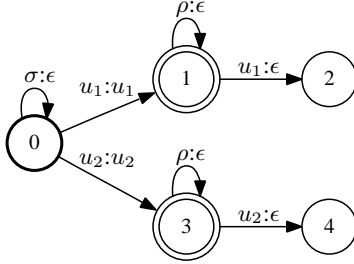
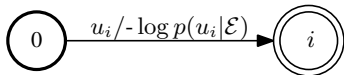


Figure 2: Path counting transducer Ψ_n^R matching last (right-most) occurrence of each $u \in \mathcal{N}_n$.

posing with Ψ_n^R creates a single copy of \mathcal{E}_n while searching for the last occurrence of u ; we find this to be much more efficient for large \mathcal{N}_n .

Path posterior probabilities are calculated over each \mathcal{E}_n by composing with Ψ_n in the log semiring, projecting on the output, removing ϵ -arcs, determinizing, minimising, and pushing weights to the initial state (Allauzen et al., 2010). Using either Ψ_n^L or Ψ_n^R , the resulting counts acceptor is \mathcal{X}_n . It has a compact form with one arc from the start state for each $u_i \in \mathcal{N}_n$:



3.1 Efficient Path Posterior Calculation

Although \mathcal{X}_n has a convenient and elegant form, it can be difficult to build for large \mathcal{N}_n because the composition $\mathcal{E}_n \circ \Psi_n$ results in millions of states and arcs. The log semiring ϵ -removal and determinization required to sum the probabilities of paths labelled with each u can be slow.

However, if we use the proposed Ψ_n^R , then each path in $\mathcal{E}_n \circ \Psi_n^R$ has only one non- ϵ output label u and all paths leading to a given final state share the same u . A modified forward algorithm can be used to calculate $p(u|\mathcal{E})$ without the costly ϵ -removal and determinization. The modification simply requires keeping track of which symbol u is encountered along each path to a final state.

More than one final state may gather probabilities for the same u ; to compute $p(u|\mathcal{E})$ these probabilities are added. The forward algorithm requires that $\mathcal{E}_n \circ \Psi_n^R$ be topologically sorted; although sorting can be slow, it is still quicker than log semiring ϵ -removal and determinization.

The statistics gathered by the forward algorithm could also be gathered under the expectation semiring (Eisner, 2002) with suitably defined features. We take the view that the full complexity of that approach is not needed here, since only one symbol is introduced per path and per exit state.

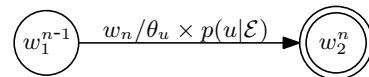
Unlike $\mathcal{E}_n \circ \Psi_n^R$, the composition $\mathcal{E}_n \circ \Psi_n^L$ does not segregate paths by u such that there is a direct association between final states and symbols. The forward algorithm does not readily yield the per-symbol probabilities, although an arc weight vector indexed by symbols could be used to correctly aggregate the required statistics (Riley et al., 2009). For large \mathcal{N}_n this would be memory intensive. The association between final states and symbols could also be found by label pushing, but we find this slow for large $\mathcal{E}_n \circ \Psi_n$.

4 Efficient Decoder Implementation

In contrast to Equation (5), we use the exact values of $p(u|\mathcal{E})$ for all $u \in \mathcal{N}_n$ at orders $n = 1 \dots 4$ to compute

$$\hat{E} = \operatorname{argmin}_{E' \in \mathcal{E}} \left\{ \theta_0 |E'| + \sum_{n=1}^4 g_n(E, E') \right\}, \quad (6)$$

where $g_n(E, E') = \sum_{u \in \mathcal{N}_n} \theta_u \#_u(E') p(u|\mathcal{E})$ using the exact path posterior probabilities at each order. We make acceptors Ω_n such that $\mathcal{E} \circ \Omega_n$ assigns order n partial gain $g_n(E, E')$ to all paths $E \in \mathcal{E}$. Ω_n is derived from Φ_n directly by assigning arc weight $\theta_u \times p(u|\mathcal{E})$ to arcs with output label u and then projecting on the input labels. For each n -gram $u = w_1^n$ in \mathcal{N}_n arcs of Ω_n have the form:



To apply θ_0 we make a copy of \mathcal{E} , called \mathcal{E}_0 , with fixed weight θ_0 on all arcs. The decoder is formed as the composition $\mathcal{E}_0 \circ \Omega_1 \circ \Omega_2 \circ \Omega_3 \circ \Omega_4$ and \hat{E} is extracted as the maximum cost string.

5 Lattice Generation for LMBR

Lattice MBR decoding performance and efficiency is evaluated in the context of the NIST

		mt0205tune	mt0205test	mt08nw	mt08ng
ML		54.2	53.8	51.4	36.3
k	0	52.6	52.3	49.8	34.5
	1	54.8	54.4	52.2	36.6
	2	54.9	54.5	52.4	36.8
	3	54.9	54.5	52.4	36.8
LMBR		55.0	54.6	52.4	36.8

Table 1: BLEU scores for Arabic→English maximum likelihood translation (ML), MBR decoding using the hybrid decision rule of Equation (5) at $0 \leq k \leq 3$, and regular linearised lattice MBR (LMBR).

		mt0205tune	mt0205test	mt08nw	mt08ng
Posteriors	sequential	3160	3306	2090	3791
	Ψ_n^L	6880	7387	4201	8796
	Ψ_n^R	1746	1789	1182	2787
Decoding	sequential	4340	4530	2225	4104
	Ψ_n	284	319	118	197
Total	sequential	7711	8065	4437	8085
	Ψ_n^L	7458	8075	4495	9199
	Ψ_n^R	2321	2348	1468	3149

Table 2: Time in seconds required for path posterior n -gram probability calculation and LMBR decoding using sequential method and left-most (Ψ_n^L) or right-most (Ψ_n^R) counting transducer implementations.

Arabic→English machine translation task³. The development set mt0205tune is formed from the odd numbered sentences of the NIST MT02–MT05 testsets; the even numbered sentences form the validation set mt0205test. Performance on NIST MT08 newswire (mt08nw) and newsgroup (mt08ng) data is also reported.

First-pass translation is performed using HiFST (Iglesias et al., 2009), a hierarchical phrase-based decoder. Word alignments are generated using MTTK (Deng and Byrne, 2008) over 150M words of parallel text for the constrained NIST MT08 Arabic→English track. In decoding, a Shallow-1 grammar with a single level of rule nesting is used and no pruning is performed in generating first-pass lattices (Iglesias et al., 2009).

The first-pass language model is a modified Kneser-Ney (Kneser and Ney, 1995) 4-gram estimated over the English parallel text and an 881M word subset of the GigaWord Third Edition (Graff et al., 2007). Prior to LMBR, the lattices are rescored with large stupid-backoff 5-gram language models (Brants et al., 2007) estimated over more than 6 billion words of English text.

The n -gram factors $\theta_0, \dots, \theta_4$ are set according to Tromble et al. (2008) using unigram precision

$p = 0.85$ and average recall ratio $r = 0.74$. Our translation decoder and MBR procedures are implemented using OpenFst (Allauzen et al., 2007).

6 LMBR Speed and Performance

Lattice MBR decoding performance is shown in Table 1. Compared to the maximum likelihood translation hypotheses (row ML), LMBR gives gains of +0.8 to +1.0 BLEU for newswire data and +0.5 BLEU for newsgroup data (row LMBR).

The other rows of Table 1 show the performance of LMBR decoding using the hybrid decision rule of Equation (5) for $0 \leq k \leq 3$. When the conditional expected counts $c(u|\mathcal{E})$ are used at all orders (i.e. $k = 0$), the hybrid decoder BLEU scores are considerably lower than even the ML scores. This poor performance is because there are many unigrams u for which $c(u|\mathcal{E})$ is much greater than $p(u|\mathcal{E})$. The consensus translation maximising the conditional expected gain is then dominated by unigram matches, significantly degrading LMBR decoding performance. Table 1 shows that for these lattices the hybrid decision rule is an accurate approximation to Equation (1) only when $k \geq 2$ and the exact contribution to the gain function is computed using the path posterior probabilities at orders $n = 1$ and $n = 2$.

³<http://www.itl.nist.gov/iad/mig/tests/mt>

We now analyse the efficiency of lattice MBR decoding using the exact path posterior probabilities of Equation (2) at all orders. We note that the sequential method and both simultaneous implementations using path counting transducers Ψ_n^L and Ψ_n^R yield the same hypotheses (allowing for numerical accuracy); they differ only in speed and memory usage.

Posteriors Efficiency Computation times for the steps in LMBR are given in Table 2. In calculating path posterior n -gram probabilities $p(u|\mathcal{E})$, we find that the use of Ψ_n^L is more than twice as slow as the sequential method. This is due to the difficulty of counting higher-order n -grams in large lattices. Ψ_n^L is effective for counting unigrams, however, since there are far fewer of them. Using Ψ_n^R is almost twice as fast as the sequential method. This speed difference is due to the simple forward algorithm. We also observe that for higher-order n , the composition $\mathcal{E}_n \circ \Psi_n^R$ requires less memory and produces a smaller machine than $\mathcal{E}_n \circ \Psi_n^L$. It is easier to count paths by the final occurrence of a symbol than by the first.

Decoding Efficiency Decoding times are significantly faster using Ω_n than the sequential method; average decoding time is around 0.1 seconds per sentence. The total time required for lattice MBR is dominated by the calculation of the path posterior n -gram probabilities, and this is a function of the number of n -grams in the lattice $|\mathcal{N}|$. For each sentence in mt0205tune, Figure 3 plots the total LMBR time for the sequential method (marked ‘o’) and for probabilities computed using Ψ_n^R (marked ‘+’). This compares the two techniques on a sentence-by-sentence basis. As $|\mathcal{N}|$ grows, the simultaneous path counting transducer is found to be much more efficient.

7 Conclusion

We have described an efficient and exact implementation of the linear approximation to LMBR using general WFST operations. A simple transducer was used to map words to sequences of n -grams in order to simplify the extraction of higher-order statistics. We presented a counting transducer Ψ_n^R that extracts the statistics required for all n -grams of order n in a single composition and allows path posterior probabilities to be computed efficiently using a modified forward procedure.

We take the view that even approximate search

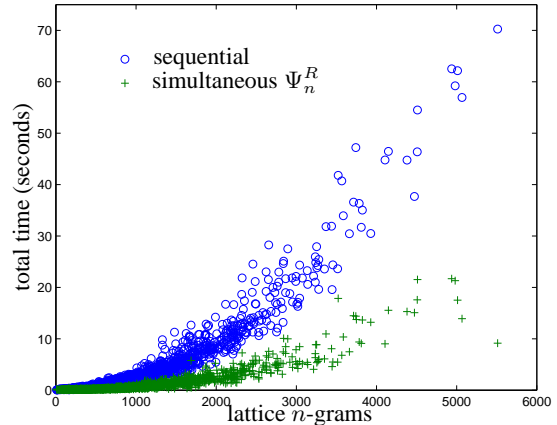


Figure 3: Total time in seconds versus $|\mathcal{N}|$.

criteria should be implemented exactly where possible, so that it is clear exactly what the system is doing. For machine translation lattices, conflating the values of $p(u|\mathcal{E})$ and $c(u|\mathcal{E})$ for higher-order n -grams might not be a serious problem, but in other scenarios – especially where symbol sequences are repeated multiple times on the same path – it may be a poor approximation.

We note that since much of the time in calculation is spent dealing with ϵ -arcs that are ultimately removed, an optimised composition algorithm that skips over such redundant structure may lead to further improvements in time efficiency.

Acknowledgments

This work was supported in part under the GALE program of the Defense Advanced Research Projects Agency, Contract No. HR0011-06-C-0022.

References

- Cyril Allauzen, Mehryar Mohri, and Brian Roark. 2003. Generalized algorithms for constructing statistical language models. In *Proceedings of the 41st Meeting of the Association for Computational Linguistics*, pages 557–564.
- Cyril Allauzen, Michael Riley, Johan Schalkwyk, Wojciech Skut, and Mehryar Mohri. 2007. OpenFst: a general and efficient weighted finite-state transducer library. In *Proceedings of the 9th International Conference on Implementation and Application of Automata*, pages 11–23. Springer.
- Cyril Allauzen, Shankar Kumar, Wolfgang Macherey, Mehryar Mohri, and Michael Riley. 2010. Expected

- sequence similarity maximization. In *Human Language Technologies 2010: The 11th Annual Conference of the North American Chapter of the Association for Computational Linguistics*, Los Angeles, California, June.
- Thorsten Brants, Ashok C. Popat, Peng Xu, Franz J. Och, and Jeffrey Dean. 2007. Large language models in machine translation. In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 858–867.
- Yonggang Deng and William Byrne. 2008. HMM word and phrase alignment for statistical machine translation. *IEEE Transactions on Audio, Speech, and Language Processing*, 16(3):494–507.
- Jason Eisner. 2002. Parameter estimation for probabilistic finite-state transducers. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 1–8, Philadelphia, July.
- David Graff, Junbo Kong, Ke Chen, and Kazuaki Maeda. 2007. English Gigaword Third Edition.
- Gonzalo Iglesias, Adrià de Gispert, Eduardo R. Banga, and William Byrne. 2009. Hierarchical phrase-based translation with weighted finite state transducers. In *Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 433–441, Boulder, Colorado, June. Association for Computational Linguistics.
- R. Kneser and H. Ney. 1995. Improved backing-off for m-gram language modeling. In *Acoustics, Speech, and Signal Processing*, pages 181–184.
- Shankar Kumar and William Byrne. 2004. Minimum Bayes-risk decoding for statistical machine translation. In *Proceedings of Human Language Technologies: The 2004 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 169–176.
- Shankar Kumar, Wolfgang Macherey, Chris Dyer, and Franz Och. 2009. Efficient minimum error rate training and minimum bayes-risk decoding for translation hypergraphs and lattices. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the Association for Computational Linguistics and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, pages 163–171, Suntec, Singapore, August. Association for Computational Linguistics.
- M. Mohri, F.C.N. Pereira, and M. Riley. 2008. Speech recognition with weighted finite-state transducers. *Handbook on Speech Processing and Speech Communication*.
- Michael Riley, Cyril Allauzen, and Martin Jansche. 2009. OpenFst: An Open-Source, Weighted Finite-State Transducer Library and its Applications to Speech and Language. In *Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics, Companion Volume: Tutorial Abstracts*, pages 9–10, Boulder, Colorado, May. Association for Computational Linguistics.
- Roy Tromble, Shankar Kumar, Franz Och, and Wolfgang Macherey. 2008. Lattice Minimum Bayes-Risk decoding for statistical machine translation. In *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*, pages 620–629, Honolulu, Hawaii, October. Association for Computational Linguistics.

The Same-head Heuristic for Coreference

Micha Elsner and Eugene Charniak

Brown Laboratory for Linguistic Information Processing (BLLIP)

Brown University

Providence, RI 02912

{melsner, ec}@cs.brown.edu

Abstract

We investigate coreference relationships between NPs with the same head noun. It is relatively common in unsupervised work to assume that such pairs are coreferent— but this is not always true, especially if realistic mention detection is used. We describe the distribution of non-coreferent same-head pairs in news text, and present an unsupervised generative model which learns not to link some same-head NPs using syntactic features, improving precision.

1 Introduction

Full NP coreference, the task of discovering which non-pronominal NPs in a discourse refer to the same entity, is widely known to be challenging. In practice, however, most work focuses on the subtask of linking NPs with different head words. Decisions involving NPs with the same head word have not attracted nearly as much attention, and many systems, especially unsupervised ones, operate under the assumption that all same-head pairs corefer. This is by no means always the case—there are several systematic exceptions to the rule. In this paper, we show that these exceptions are fairly common, and describe an unsupervised system which learns to distinguish them from coreferent same-head pairs.

There are several reasons why relatively little attention has been paid to same-head pairs. Primarily, this is because they are a comparatively easy subtask in a notoriously difficult area; Stoyanov et al. (2009) shows that, among NPs headed by common nouns, those which have an exact match earlier in the document are the easiest to

resolve (variant MUC score .82 on MUC-6) and while those with partial matches are quite a bit harder (.53), by far the worst performance is on those without any match at all (.27). This effect is magnified by most popular metrics for coreference, which reward finding links within large clusters more than they punish proposing spurious links, making it hard to improve performance by linking conservatively. Systems that use gold mention boundaries (the locations of NPs marked by annotators)¹ have even less need to worry about same-head relationships, since most NPs which disobey the conventional assumption are not marked as mentions.

In this paper, we count how often same-head pairs fail to corefer in the MUC-6 corpus, showing that gold mention detection hides most such pairs, but more realistic detection finds large numbers. We also present an unsupervised generative model which learns to make certain same-head pairs non-coreferent. The model is based on the idea that pronoun referents are likely to be salient noun phrases in the discourse, so we can learn about NP antecedents using pronominal antecedents as a starting point. Pronoun anaphora, in turn, is learnable from raw data (Cherry and Bergsma, 2005; Charniak and Elsner, 2009). Since our model links fewer NPs than the baseline, it improves precision but decreases recall. This tradeoff is favorable for CEAF, but not for b^3 .

2 Related work

Unsupervised systems specify the assumption of same-head coreference in several ways: by as-

¹Gold mention detection means something slightly different in the ACE corpus, where the system input contains every NP annotated with an entity type.

sumption (Haghighi and Klein, 2009), using a head-prediction clause (Poon and Domingos, 2008), and using a sparse Dirichlet prior on word emissions (Haghighi and Klein, 2007). (These three systems, perhaps not coincidentally, use gold mentions.) An exception is Ng (2008), who points out that head identity is not an entirely reliable cue and instead uses exact string match (minus determiners) for common NPs and an alias detection system for proper NPs. This work uses mentions extracted with an NP chunker. No specific results are reported for same-head NPs. However, while using exact string match raises precision, many non-matching phrases are still coreferent, so this approach cannot be considered a full solution to the problem.

Supervised systems do better on the task, but not perfectly. Recent work (Stoyanov et al., 2009) attempts to determine the contributions of various categories of NP to coreference scores, and shows (as stated above) that common NPs which partially match an earlier mention are not well resolved by the state-of-the-art RECONCILE system, which uses pairwise classification. They also show that using gold mention boundaries makes the coreference task substantially easier, and argue that this experimental setting is “rather unrealistic”.

3 Descriptive study: MUC-6

We begin by examining how often non-same-head pairs appear in the MUC-6 coreference dataset. To do so, we compare two artificial coreference systems: the **link-all** strategy links all, and only, full (non-pronominal) NP pairs with the same head which occur within 10 sentences of one another. The **oracle** strategy links NP pairs with the same head which occur within 10 sentences, but only if they are actually coreferent (according to the gold annotation)² The link-all system, in other words, does what most existing unsupervised systems do on the same-head subset of NPs, while the oracle system performs perfectly.

We compare our results to the gold standard using two metrics. b^3 (Bagga and Baldwin, 1998) is a standard metric which calculates a precision and recall for each mention. The mention CEAF (Luo, 2005) constructs a maximum-weight bipar-

²The choice of 10 sentences as the window size captures most, but not all, of the available recall. Using *nouns* mention detection, it misses 117 possible same-head links, or about 10%. However, precision drops further as the window size increases.

tite matching between gold and proposed clusters, then gives the percentage of entities whose gold label and proposed label match. b^3 gives more weight to errors involving larger clusters (since these lower scores for several mentions at once); for mention CEAF, all mentions are weighted equally.

We annotate the data with the self-trained Charniak parser (McClosky et al., 2006), then extract mentions using three different methods. The **gold mentions** method takes only mentions marked by annotators. The **nps** method takes all base noun phrases detected by the parser. Finally, the **nouns** method takes all nouns, even those that do not head NPs; this method maximizes recall, since it does not exclude prenominals in phrases like “a **Bush** spokesman”. (High-precision models of the internal structure of flat Penn Treebank-style NPs were investigated by Vadas and Curran (2007).) For each experimental setting, we show the number of **mentions** detected, and how many of them are **linked** to some antecedent by the system.

The data is shown in Table 1. b^3 shows a large drop in precision when all same-head pairs are linked; in fact, in the *nps* and *nouns* settings, only about half the same-headed NPs are actually coreferent (864 real links, 1592 pairs for *nps*). This demonstrates that non-coreferent same-head pairs not only occur, but are actually rather common in the dataset. The drop in precision is much less obvious in the *gold mentions* setting, however; most unlinked same-head pairs are not annotated as mentions in the gold data, which is one reason why systems run in this experimental setting can afford to ignore them.

Improperly linking same-head pairs causes a loss in precision, but scores are dominated by recall³. Thus, reporting b^3 helps to mask the impact of these pairs when examining the final f-score.

We roughly characterize what sort of same-headed NPs are non-coreferent by hand-examining 100 randomly selected pairs. 39 pairs denoted different entities (“recent employees” vs “employees who have worked for longer”) disambiguated by modifiers or sometimes by discourse position. The next largest group (24) consists of time and measure phrases like “ten miles”. 12 pairs refer to parts or quantities

³This bias is exaggerated for systems which only link same-head pairs, but continues to apply to real systems; for instance (Haghighi and Klein, 2009) has a b^3 precision of 84 and recall of 67.

	Mentions	Linked	b^3 pr	rec	F	mention CEAF
Gold mentions						
Oracle	1929	1164	100	32.3	48.8	54.4
Link all	1929	1182	80.6	31.7	45.5	53.8
Alignment	1929	495	93.7	22.1	35.8	40.5
NPs						
Oracle	3993	864	100	30.6	46.9	73.4
Link all	3993	1592	<u>67.2</u>	29.5	41.0	<u>62.2</u>
Alignment	3993	518	<u>87.2</u>	24.7	38.5	<u>67.0</u>
Nouns						
Oracle	5435	1127	100	41.5	58.6	83.5
Link all	5435	2541	<u>56.6</u>	40.9	45.7	<u>67.0</u>
Alignment	5435	935	<u>83.0</u>	32.8	47.1	<u>74.4</u>

Table 1: Oracle, system and baseline scores on MUC-6 test data. **Gold mentions leave little room for improvement** between baseline and oracle; **detecting more mentions widens the gap between them**. With realistic mention detection, precision and CEAF scores improve over baselines, while recall and f-scores drop.

(“members of...”), and 12 contained a generic (“In a corporate campaign, a union tries...”). 9 contained an annotator error. The remaining 4 were mistakes involving proper noun phrases headed by *Inc.* and other abbreviations; this case is easy to handle, but apparently not the primary cause of errors.

4 System

Our system is a version of the popular IBM model 2 for machine translation. To define our generative model, we assume that the parse trees for the entire document D are given, *except* for the subtrees with root nonterminal NP, denoted n_i , which our system will generate. These subtrees are related by a hidden set of alignments, a_i , which link each NP to another NP (which we call a *generator*) appearing somewhere before it in the document, or to a null antecedent. The set of potential generators G (which plays the same role as the source-language text in MT) is taken to be all the NPs occurring within 10 sentences of the target, plus a special null antecedent which plays the same role as the null word in machine translation—it serves as a dummy generator for NPs which are unrelated to any real NP in G .

The generative process fills in all the NP nodes in order, from left to right. This process ensures that, when generating node n_i , we have already filled in all the NPs in the set G (since these all precede n_i). When deciding on a generator for NP n_i , we can extract features characterizing its

relationship to a potential generator g_j . These features, which we denote $f(n_i, g_j, D)$, may depend on their relative position in the document D , and on any features of g_j , since we have already generated its tree. However, we cannot extract features from the subtree under n_i , since we have yet to generate it!

As usual for IBM models, we learn using EM, and we need to start our alignment function off with a good initial set of parameters. Since antecedents of NPs and pronouns (both salient NPs) often occur in similar syntactic environments, we use an alignment function for pronoun coreference as a starting point. This alignment can be learned from raw data, making our approach unsupervised.

We take the pronoun model of Charniak and Elsnier (2009)⁴ as our starting point. We re-express it in the IBM framework, using a log-linear model for our alignment. Then our alignment (parameterized by feature weights w) is:

$$p(a_i = j | G, D) \propto \exp(f(n_i, g_j, D) \bullet w)$$

The weights w are learned by gradient descent on the log-likelihood. To use this model within EM, we alternate an E-step where we calculate the expected alignments $E[a_i = j]$, then an M-step where we run gradient descent. (We have also had some success with stepwise EM as in (Liang and Klein, 2009), but this requires some tuning to work properly.)

⁴Downloaded from <http://bllip.cs.brown.edu>.

As features, we take the same features as Charniak and Elsnar (2009): sentence and word-count distance between n_i and g_j , sentence position of each, syntactic role of each, and head type of g_j (proper, common or pronoun). We add binary features for the nonterminal directly over g_j (NP, VP, PP, any S type, or other), the type of phrases modifying g_j (proper nouns, phrasals (except QP and PP), QP, PP-of, PP-other, other modifiers, or nothing), and the type of determiner of g_j (possessive, definite, indefinite, deictic, other, or nothing). We designed this feature set to distinguish prominent NPs in the discourse, and also to be able to detect abstract or partitive phrases by examining modifiers and determiners.

To produce full NPs and learn same-head coreference, we focus on learning a good alignment using the pronoun model as a starting point. For translation, we use a trivial model, $p(n_i|g_{a_i}) = 1$ if the two have the same head, and 0 otherwise, except for the null antecedent, which draws heads from a multinomial distribution over words.

While we could learn an alignment and then treat all generators as antecedents, so that only NPs aligned to the null antecedent were not labeled coreferent, in practice this model would align nearly all the same-head pairs. This is true because many words are “bursty”; the probability of a second occurrence given the first is higher than the a priori probability of occurrence (Church, 2000). Therefore, our model is actually a mixture of two IBM models, p_C and p_N , where p_C produces NPs with antecedents and p_N produces pairs that share a head, but are not coreferent. To break the symmetry, we allow p_C to use any parameters w , while p_N uses a uniform alignment, $w \equiv \vec{0}$. We interpolate between these two models with a constant λ , the single manually set parameter of our system, which we fixed at .9.

The full model, therefore, is:

$$\begin{aligned}
 p(n_i|G, D) &= \lambda p_T(n_i|G, D) \\
 &\quad + (1 - \lambda) p_N(n_i|G, D) \\
 p_T(n_i|G, D) &= \frac{1}{Z} \sum_{j \in G} \exp(f(n_i, g_j, D) \bullet w) \\
 &\quad \times \mathbb{I}\{\text{head}(n_i) = \text{head}(j)\} \\
 p_N(n_i|G, D) &= \sum_{j \in G} \frac{1}{|G|} \mathbb{I}\{\text{head}(n_i) = \text{head}(g_j)\}
 \end{aligned}$$

NPs for which the maximum-likelihood gener-

ator (the largest term in either of the sums) is from p_T and is not the null antecedent are marked as coreferent to the generator. Other NPs are marked not coreferent.

5 Results

Our results on the MUC-6 formal test set are shown in Table 1. In all experimental settings, the model improves precision over the baseline while decreasing recall— that is, it misses some legitimate coreferent pairs while correctly excluding many of the spurious ones. Because of the precision-recall tradeoff at which the systems operate, this results in reduced b^3 and link F. However, for the *nps* and *nouns* settings, where the parser is responsible for finding mentions, the tradeoff is positive for the CEAF metrics. For instance, in the *nps* setting, it improves over baseline by 57%.

As expected, the model does poorly in the gold mentions setting, doing worse than baseline on both metrics. Although it is possible to get very high precision in this setting, the model is far too conservative, linking less than half of the available mentions to anything, when in fact about 60% of them are coreferent. As we explain above, this experimental setting makes it mostly unnecessary to worry about non-coreferent same-head pairs because the MUC-6 annotators don’t often mark them.

6 Conclusions

While same-head pairs are easier to resolve than same-other pairs, they are still non-trivial and deserve further attention in coreference research. To effectively measure their effect on performance, researchers should report multiple metrics, since under b^3 the link-all heuristic is extremely difficult to beat. It is also important to report results using a realistic mention detector as well as gold mentions.

Acknowledgements

We thank Jean Carletta for the SWITCHBOARD annotations, and Dan Jurafsky and eight anonymous reviewers for their comments and suggestions. This work was funded by a Google graduate fellowship.

References

- Amit Bagga and Breck Baldwin. 1998. Algorithms for scoring coreference chains. In *LREC Workshop on Linguistics Coreference*, pages 563–566.
- Eugene Charniak and Micha Elsner. 2009. EM works for pronoun anaphora resolution. In *Proceedings of EACL*, Athens, Greece.
- Colin Cherry and Shane Bergsma. 2005. An Expectation Maximization approach to pronoun resolution. In *Proceedings of CoNLL*, pages 88–95, Ann Arbor, Michigan.
- Kenneth W. Church. 2000. Empirical estimates of adaptation: the chance of two Noriegas is closer to $p/2$ than p^2 . In *Proceedings of ACL*, pages 180–186.
- Aria Haghighi and Dan Klein. 2007. Unsupervised coreference resolution in a nonparametric Bayesian model. In *Proceedings of ACL*, pages 848–855.
- Aria Haghighi and Dan Klein. 2009. Simple coreference resolution with rich syntactic and semantic features. In *Proceedings of EMNLP*, pages 1152–1161.
- Percy Liang and Dan Klein. 2009. Online EM for unsupervised models. In *HLT-NAACL*.
- Xiaoqiang Luo. 2005. On coreference resolution performance metrics. In *Proceedings of HLT-EMNLP*, pages 25–32, Morristown, NJ, USA. Association for Computational Linguistics.
- David McClosky, Eugene Charniak, and Mark Johnson. 2006. Effective self-training for parsing. In *Proceedings of HLT-NAACL*, pages 152–159.
- Vincent Ng. 2008. Unsupervised models for coreference resolution. In *Proceedings of EMNLP*, pages 640–649, Honolulu, Hawaii. Association for Computational Linguistics.
- Hoifung Poon and Pedro Domingos. 2008. Joint unsupervised coreference resolution with Markov Logic. In *Proceedings of EMNLP*, pages 650–659, Honolulu, Hawaii, October. Association for Computational Linguistics.
- Veselin Stoyanov, Nathan Gilbert, Claire Cardie, and Ellen Riloff. 2009. Conundrums in noun phrase coreference resolution: Making sense of the state-of-the-art. In *Proceedings of ACL-IJCNLP*, pages 656–664, Suntec, Singapore, August. Association for Computational Linguistics.
- David Vadas and James Curran. 2007. Adding noun phrase structure to the penn treebank. In *Proceedings of ACL*, pages 240–247, Prague, Czech Republic, June. Association for Computational Linguistics.

Authorship Attribution Using Probabilistic Context-Free Grammars

Sindhu Raghavan Adriana Kovashka Raymond Mooney

Department of Computer Science

The University of Texas at Austin

1 University Station C0500

Austin, TX 78712-0233, USA

{sindhu, adriana, mooney}@cs.utexas.edu

Abstract

In this paper, we present a novel approach for authorship attribution, the task of identifying the author of a document, using probabilistic context-free grammars. Our approach involves building a probabilistic context-free grammar for each author and using this grammar as a language model for classification. We evaluate the performance of our method on a wide range of datasets to demonstrate its efficacy.

1 Introduction

Natural language processing allows us to build language models, and these models can be used to distinguish between languages. In the context of written text, such as newspaper articles or short stories, the author’s style could be considered a distinct “language.” Authorship attribution, also referred to as authorship identification or prediction, studies strategies for discriminating between the styles of different authors. These strategies have numerous applications, including settling disputes regarding the authorship of old and historically important documents (Mosteller and Wallace, 1984), automatic plagiarism detection, determination of document authenticity in court (Juola and Sofko, 2004), cyber crime investigation (Zheng et al., 2009), and forensics (Luyckx and Daelemans, 2008).

The general approach to authorship attribution is to extract a number of style markers from the text and use these style markers as features to train a classifier (Burrows, 1987; Binongo and Smith, 1999; Diederich et al., 2000; Holmes and Forsyth, 1995; Joachims, 1998; Mosteller and Wallace, 1984). These style markers could include the frequencies of certain characters, function words, phrases or sentences. Peng et al. (2003) build a character-level n -gram model for each author. Stamatatos et al. (1999) and Luyckx and Daelemans

(2008) use a combination of word-level statistics and part-of-speech counts or n -grams. Baayen et al. (1996) demonstrate that the use of syntactic features from parse trees can improve the accuracy of authorship attribution. While there have been several approaches proposed for authorship attribution, it is not clear if the performance of one is better than the other. Further, it is difficult to compare the performance of these algorithms because they were primarily evaluated on different datasets. For more information on the current state of the art for authorship attribution, we refer the reader to a detailed survey by Stamatatos (2009).

We further investigate the use of syntactic information by building complete models of each author’s syntax to distinguish between authors. Our approach involves building a probabilistic context-free grammar (PCFG) for each author and using this grammar as a language model for classification. Experiments on a variety of corpora including poetry and newspaper articles on a number of topics demonstrate that our PCFG approach performs fairly well, but it only outperforms a bigram language model on a couple of datasets (e.g. poetry). However, combining our approach with other methods results in an ensemble that performs the best on most datasets.

2 Authorship Attribution using PCFG

We now describe our approach to authorship attribution. Given a training set of documents from different authors, we build a PCFG for each author based on the documents they have written. Given a test document, we parse it using each author’s grammar and assign it to the author whose PCFG produced the highest likelihood for the document.

In order to build a PCFG, a standard statistical parser takes a corpus of parse trees of sentences as training input. Since we do not have access to authors’ documents annotated with parse trees, we use a statistical parser trained on a generic

corpus like the Wall Street Journal (WSJ) or Brown corpus from the Penn Treebank (<http://www.cis.upenn.edu/~treebank/>) to automatically annotate (i.e. treebank) the training documents for each author. In our experiments, we used the Stanford Parser (Klein and Manning, 2003b; Klein and Manning, 2003a) and the OpenNLP sentence segmenter (<http://opennlp.sourceforge.net/>). Our approach is summarized below:

Input – A training set of documents labeled with author names and a test set of documents with unknown authors.

1. Train a statistical parser on a generic corpus like the WSJ or Brown corpus.
2. Treebank each training document using the parser trained in Step 1.
3. Train a PCFG G_i for each author A_i using the treebanked documents for that author.
4. For each test document, compute its likelihood for each grammar G_i by multiplying the probability of the top PCFG parse for each sentence.
5. For each test document, find the author A_i whose grammar G_i results in the highest likelihood score.

Output – A label (author name) for each document in the test set.

3 Experimental Comparison

This section describes experiments evaluating our approach on several real-world datasets.

3.1 Data

We collected a variety of documents with known authors including news articles on a wide range of topics and literary works like poetry. We downloaded all texts from the Internet and manually removed extraneous information as well as titles, author names, and chapter headings. We collected several news articles from the New York Times online journal (<http://global.nytimes.com/>) on topics related to business, travel, and football. We also collected news articles on cricket from the ESPN cricinfo website (<http://www.cricinfo.com>).

In addition, we collected poems from the Project Gutenberg website (http://www.gutenberg.org/wiki/Main_Page). We attempted to collect sets of documents on a shared topic written by multiple authors. This was done to ensure that the datasets truly tested authorship attribution as opposed to topic identification. However, since it is very difficult to find authors that write literary works on the same topic, the Poetry dataset exhibits higher topic variability than our news datasets. We had 5 different datasets in total – Football, Business, Travel, Cricket, and Poetry. The number of authors in our datasets ranged from 3 to 6.

For each dataset, we split the documents into training and test sets. Previous studies (Stamatatos et al., 1999) have observed that having unequal number of words per author in the training set leads to poor performance for the authors with fewer words. Therefore, we ensured that, in the training set, the total number of words per author was roughly the same. We would like to note that we could have also selected the training set such that the total number of sentences per author was roughly the same. However, since we would like to compare the performance of the PCFG-based approach with a bag-of-words baseline, we decided to normalize the training set based on the number of words, rather than sentences. For testing, we used 15 documents per author for datasets with news articles and 5 or 10 documents per author for the Poetry dataset. More details about the datasets can be found in Table 1.

Dataset	# authors	# words/auth	# docs/auth	# sent/auth
Football	3	14374.67	17.3	786.3
Business	6	11215.5	14.16	543.6
Travel	4	23765.75	28	1086
Cricket	4	23357.25	24.5	1189.5
Poetry	6	7261.83	24.16	329

Table 1: Statistics for the training datasets used in our experiments. The numbers in columns 3, 4 and 5 are averages.

3.2 Methodology

We evaluated our approach to authorship prediction on the five datasets described above. For news articles, we used the first 10 sections of the WSJ corpus, which consists of annotated news articles on finance, to build the initial statistical parser in

Step 1. For Poetry, we used 7 sections of the Brown corpus which consists of annotated documents from different areas of literature.

In the basic approach, we trained a PCFG model for each author based solely on the documents written by that author. However, since the number of documents per author is relatively low, this leads to very sparse training data. Therefore, we also augmented the training data by adding one, two or three sections of the WSJ or Brown corpus to each training set, and up-sampling (replicating) the data from the original author. We refer to this model as “PCFG-*I*”, where *I* stands for *interpolation* since this effectively exploits linear interpolation with the base corpus to smooth parameters. Based on our preliminary experiments, we replicated the original data three or four times.

We compared the performance of our approach to bag-of-words classification and n -gram language models. When using bag-of-words, one generally removes commonly occurring “stop words.” However, for the task of authorship prediction, we hypothesized that the frequency of specific stop words could provide useful information about the author’s writing style. Preliminary experiments verified that eliminating stop words degraded performance; therefore, we did not remove them. We used the Maximum Entropy (MaxEnt) and Naive Bayes classifiers in the MALLET software package (McCallum, 2002) as initial baselines. We surmised that a discriminative classifier like MaxEnt might perform better than a generative classifier like Naive Bayes. However, when sufficient training data is not available, generative models are known to perform better than discriminative models (Ng and Jordan, 2001). Hence, we chose to compare our method to both Naive Bayes and MaxEnt.

We also compared the performance of the PCFG approach against n -gram language models. Specifically, we tried unigram, bigram and trigram models. We used the same background corpus mixing method used for the PCFG-*I* model to effectively smooth the n -gram models. Since a generative model like Naive Bayes that uses n -gram frequencies is equivalent to an n -gram language model, we also used the Naive Bayes classifier in MALLET to implement the n -gram models. Note that a Naive-Bayes bag-of-words model is equivalent to a unigram language model.

While the PCFG model captures the author’s

writing style at the syntactic level, it may not accurately capture lexical information. Since both syntactic and lexical information is presumably useful in capturing the author’s overall writing style, we also developed an ensemble using a PCFG model, the bag-of-words MaxEnt classifier, and an n -gram language model. We linearly combined the confidence scores assigned by each model to each author, and used the combined score for the final classification. We refer to this model as “PCFG-*E*”, where *E* stands for *ensemble*. We also developed another ensemble based on MaxEnt and n -gram language models to demonstrate the contribution of the PCFG model to the overall performance of PCFG-*E*. For each dataset, we report *accuracy*, the fraction of the test documents whose authors were correctly identified.

3.3 Results and Discussion

Table 2 shows the accuracy of authorship prediction on different datasets. For the n -gram models, we only report the results for the bigram model with smoothing (Bigram-*I*) as it was the best performing model for most datasets (except for Cricket and Poetry). For the Cricket dataset, the trigram-*I* model was the best performing n -gram model with an accuracy of 98.34%. Generally, a higher order n -gram model ($n = 3$ or higher) performs poorly as it requires a fair amount of smoothing due to the exponential increase in all possible n -gram combinations. Hence, the superior performance of the trigram-*I* model on the Cricket dataset was a surprising result. For the Poetry dataset, the unigram-*I* model performed best among the smoothed n -gram models at 81.8% accuracy. This is unsurprising because as mentioned above, topic information is strongest in the Poetry dataset, and it is captured well in the unigram model. For bag-of-words methods, we find that the generatively trained Naive Bayes model (unigram language model) performs better than or equal to the discriminatively trained MaxEnt model on most datasets (except for Business). This result is not surprising since our datasets are limited in size, and generative models tend to perform better than discriminative methods when there is very little training data available. Amongst the different baseline models (MaxEnt, Naive Bayes, Bigram-*I*), we find Bigram-*I* to be the best performing model (except for Cricket and Poetry). For both Cricket and Poetry, Naive Bayes

Dataset	MaxEnt	Naive Bayes	Bigram- <i>I</i>	PCFG	PCFG- <i>I</i>	PCFG- <i>E</i>	MaxEnt+Bigram- <i>I</i>
Football	84.45	86.67	86.67	93.34	80	91.11	86.67
Business	83.34	77.78	90.00	77.78	85.56	91.11	92.22
Travel	83.34	83.34	91.67	81.67	86.67	91.67	90.00
Cricket	91.67	95.00	91.67	86.67	91.67	95.00	93.34
Poetry	56.36	78.18	70.90	78.18	83.63	87.27	76.36

Table 2: Accuracy in % for authorship prediction on different datasets. Bigram-*I* refers to the bigram language model with smoothing. PCFG-*E* refers to the ensemble based on MaxEnt, Bigram-*I*, and PCFG-*I*. MaxEnt+Bigram-*I* refers to the ensemble based on MaxEnt and Bigram-*I*.

is the best performing baseline model. While discussing the performance of the PCFG model and its variants, we consider the best performing baseline model.

We observe that the basic PCFG model and the PCFG-*I* model do not usually outperform the best baseline method (except for Football and Poetry, as discussed below). For Football, the basic PCFG model outperforms the best baseline, while for Poetry, the PCFG-*I* model outperforms the best baseline. Further, the performance of the basic PCFG model is inferior to that of PCFG-*I* for most datasets, likely due to the insufficient training data used in the basic model. Ideally one would use more training documents, but in many domains it is impossible to obtain a large corpus of documents written by a single author. For example, as Luyckx and Daelemans (2008) argue, in forensics one would like to identify the authorship of documents based on a limited number of documents written by the author. Hence, we investigated smoothing techniques to improve the performance of the basic PCFG model. We found that the interpolation approach resulted in a substantial improvement in the performance of the PCFG model for all but the Football dataset (discussed below). However, for some datasets, even this improvement was not sufficient to outperform the best baseline.

The results for PCFG and PCFG-*I* demonstrate that syntactic information alone is generally a bit less accurate than using *n*-grams. In order to utilize *both* syntactic and lexical information, we developed PCFG-*E* as described above. We combined the best *n*-gram model (Bigram-*I*) and PCFG model (PCFG-*I*) with MaxEnt to build PCFG-*E*. For the Travel dataset, we find that the performance of the PCFG-*E* model is equal to that of the best constituent model (Bigram-*I*). For the remaining datasets, the performance of PCFG-*E*

is better than the best constituent model. Furthermore, for the Football, Cricket and Poetry datasets this improvement is quite substantial. We now find that the performance of some variant of PCFG is always better than or equal to that of the best baseline. While the basic PCFG model outperforms the baseline for the Football dataset, PCFG-*E* outperforms the best baseline for the Poetry and Business datasets. For the Cricket and Travel datasets, the performance of the PCFG-*E* model equals that of the best baseline. In order to assess the statistical significance of any performance difference between the best PCFG model and the best baseline, we performed the McNemar’s test, a non-parametric test for binomial variables (Rosner, 2005). We found that the difference in the performance of the two methods was not statistically significant at .05 significance level for any of the datasets, probably due to the small number of test samples.

The performance of PCFG and PCFG-*I* is particularly impressive on the Football and Poetry datasets. For the Football dataset, the basic PCFG model is the best performing PCFG model and it performs much better than other methods. It is surprising that smoothing using PCFG-*I* actually results in a drop in performance on this dataset. We hypothesize that the authors in the Football dataset may have very different syntactic writing styles that are effectively captured by the basic PCFG model. Smoothing the data apparently weakens this signal, hence causing a drop in performance. For Poetry, PCFG-*I* achieves much higher accuracy than the baselines. This is impressive given the much looser syntactic structure of poetry compared to news articles, and it indicates the value of syntactic information for distinguishing between literary authors.

Finally, we consider the specific contribution of the PCFG-*I* model towards the performance of

the PCFG-*E* ensemble. Based on comparing the results for PCFG-*E* and MaxEnt+Bigram-*I*, we find that there is a drop in performance for most datasets when removing PCFG-*I* from the ensemble. This drop is quite substantial for the Football and Poetry datasets. This indicates that PCFG-*I* is contributing substantially to the performance of PCFG-*E*. Thus, it further illustrates the importance of broader syntactic information for the task of authorship attribution.

4 Future Work and Conclusions

In this paper, we have presented our ongoing work on authorship attribution, describing a novel approach that uses probabilistic context-free grammars. We have demonstrated that both syntactic and lexical information are useful in effectively capturing authors' overall writing style. To this end, we have developed an ensemble approach that performs better than the baseline models on several datasets. An interesting extension of our current approach is to consider discriminative training of PCFGs for each author. Finally, we would like to compare the performance of our method to other state-of-the-art approaches to authorship prediction.

Acknowledgments

Experiments were run on the Mastodon Cluster, provided by NSF Grant EIA-0303609.

References

- H. Baayen, H. van Halteren, and F. Tweedie. 1996. Outside the cave of shadows: using syntactic annotation to enhance authorship attribution. *Literary and Linguistic Computing*, 11(3):121–132, September.
- Binongo and Smith. 1999. A Study of Oscar Wilde's Writings. *Journal of Applied Statistics*, 26:781.
- J Burrows. 1987. Word-patterns and Story-shapes: The Statistical Analysis of Narrative Style.
- Joachim Diederich, Jörg Kindermann, Edda Leopold, and Gerhard Paass. 2000. Authorship Attribution with Support Vector Machines. *Applied Intelligence*, 19:2003.
- D. I. Holmes and R. S. Forsyth. 1995. The Federalist Revisited: New Directions in Authorship Attribution. *Literary and Linguistic Computing*, 10:111–127.
- Thorsten Joachims. 1998. Text categorization with Support Vector Machines: Learning with many relevant features. In *Proceedings of the 10th European Conference on Machine Learning (ECML)*, pages 137–142, Berlin, Heidelberg. Springer-Verlag.
- Patrick Juola and John Sofko. 2004. Proving and Improving Authorship Attribution Technologies. In *Proceedings of Canadian Symposium for Text Analysis (CaSTA)*.
- Dan Klein and Christopher D. Manning. 2003a. Accurate unlexicalized parsing. In *Proceedings of the 41st Annual Meeting on Association for Computational Linguistics (ACL)*, pages 423–430, Morristown, NJ, USA. Association for Computational Linguistics.
- Dan Klein and Christopher D. Manning. 2003b. Fast Exact Inference with a Factored Model for Natural Language Parsing. In *Advances in Neural Information Processing Systems 15 (NIPS)*, pages 3–10. MIT Press.
- Kim Luyckx and Walter Daelemans. 2008. Authorship Attribution and Verification with Many Authors and Limited Data. In *Proceedings of the 22nd International Conference on Computational Linguistics (COLING)*, pages 513–520, August.
- Andrew Kachites McCallum. 2002. MALLET: A Machine Learning for Language Toolkit. <http://mallet.cs.umass.edu>.
- Frederick Mosteller and David L. Wallace. 1984. *Applied Bayesian and Classical Inference: The Case of the Federalist Papers*. Springer-Verlag.
- Andrew Y. Ng and Michael I. Jordan. 2001. On Discriminative vs. Generative classifiers: A comparison of logistic regression and naive Bayes. In *Advances in Neural Information Processing Systems 14 (NIPS)*, pages 841–848.
- Fuchun Peng, Dale Schuurmans, Viado Keselj, and Shaojun Wang. 2003. Language Independent Authorship Attribution using Character Level Language Models. In *Proceedings of the 10th Conference of the European Chapter of the Association for Computational Linguistics (EACL)*.
- Bernard Rosner. 2005. *Fundamentals of Biostatistics*. Duxbury Press.
- E. Stamatatos, N. Fakotakis, and G. Kokkinakis. 1999. Automatic Authorship Attribution. In *Proceedings of the 9th Conference of the European Chapter of the Association for Computational Linguistics (EACL)*, pages 158–164, Morristown, NJ, USA. Association for Computational Linguistics.
- E. Stamatatos. 2009. A Survey of Modern Authorship Attribution Methods. *Journal of the American Society for Information Science and Technology*, 60(3):538–556.
- Rong Zheng, Yi Qin, Zan Huang, and Hsinchun Chen. 2009. Authorship Analysis in Cybercrime Investigation. *Lecture Notes in Computer Science*, 2665/2009:959.

The impact of interpretation problems on tutorial dialogue

Myroslava O. Dzikovska and Johanna D. Moore

School of Informatics, University of Edinburgh, Edinburgh, United Kingdom
{m.dzikovska, j.moore}@ed.ac.uk

Natalie Steinhauser and Gwendolyn Campbell

Naval Air Warfare Center Training Systems Division, Orlando, FL, USA
{natalie.steihauser, gwendolyn.campbell}@navy.mil

Abstract

Supporting natural language input may improve learning in intelligent tutoring systems. However, interpretation errors are unavoidable and require an effective recovery policy. We describe an evaluation of an error recovery policy in the BEE-TLE II tutorial dialogue system and discuss how different types of interpretation problems affect learning gain and user satisfaction. In particular, the problems arising from student use of non-standard terminology appear to have negative consequences. We argue that existing strategies for dealing with terminology problems are insufficient and that improving such strategies is important in future ITS research.

1 Introduction

There is a mounting body of evidence that student self-explanation and contentful talk in human-human tutorial dialogue are correlated with increased learning gain (Chi et al., 1994; Purandare and Litman, 2008; Litman et al., 2009). Thus, computer tutors that understand student explanations have the potential to improve student learning (Graesser et al., 1999; Jordan et al., 2006; Alevan et al., 2001; Dzikovska et al., 2008). However, understanding and correctly assessing the student's contributions is a difficult problem due to the wide range of variation observed in student input, and especially due to students' sometimes vague and incorrect use of domain terminology.

Many tutorial dialogue systems limit the range of student input by asking short-answer questions. This provides a measure of robustness, and previous evaluations of ASR in spoken tutorial dialogue systems indicate that neither word error rate nor concept error rate in such systems affect learning gain (Litman and Forbes-Riley, 2005; Pon-Barry

et al., 2004). However, limiting the range of possible input limits the contentful talk that the students are expected to produce, and therefore may limit the overall effectiveness of the system.

Most of the existing tutoring systems that accept unrestricted language input use classifiers based on statistical text similarity measures to match student answers to open-ended questions with pre-authored anticipated answers (Graesser et al., 1999; Jordan et al., 2004; McCarthy et al., 2008). While such systems are robust to unexpected terminology, they provide only a very coarse-grained assessment of student answers. Recent research aims to develop methods that produce detailed analyses of student input, including correct, incorrect and missing parts (Nielsen et al., 2008; Dzikovska et al., 2008), because the more detailed assessments can help tailor tutoring to the needs of individual students.

While the detailed assessments of answers to open-ended questions are intended to improve potential learning, they also increase the probability of misunderstandings, which negatively impact tutoring and therefore negatively impact student learning (Jordan et al., 2009). Thus, appropriate error recovery strategies are crucially important for tutorial dialogue applications. We describe an evaluation of an implemented tutorial dialogue system which aims to accept unrestricted student input and limit misunderstandings by rejecting low confidence interpretations and employing a range of error recovery strategies depending on the cause of interpretation failure.

By comparing two different system policies, we demonstrate that with less restricted language input the rate of non-understanding errors impacts both learning gain and user satisfaction, and that problems arising from incorrect use of terminology have a particularly negative impact. A more detailed analysis of the results indicates that, even though we based our policy on an approach ef-

fective in task-oriented dialogue (Hockey et al., 2003), many of our strategies were not successful in improving learning gain. At the same time, students appear to be aware that the system does not fully understand them even if it accepts their input without indicating that it is having interpretation problems, and this is reflected in decreased user satisfaction. We argue that this indicates that we need better strategies for dealing with terminology problems, and that accepting non-standard terminology without explicitly addressing the difference in acceptable phrasing may not be sufficient for effective tutoring.

In Section 2 we describe our tutoring system, and the two tutoring policies implemented for the experiment. In Section 3 we present experimental results and an analysis of correlations between different types of interpretation problems, learning gain and user satisfaction. Finally, in Section 4 we discuss the implications of our results for error recovery policies in tutorial dialogue systems.

2 Tutorial Dialogue System and Error Recovery Policies

This work is based on evaluation of BEETLE II (Dzikovska et al., 2010), a tutorial dialogue system which provides tutoring in basic electricity and electronics. Students read pre-authored materials, experiment with a circuit simulator, and then are asked to explain their observations. BEETLE II uses a deep parser together with a domain-specific diagnoser to process student input, and a deep generator to produce tutorial feedback automatically depending on the current tutorial policy. It also implements an error recovery policy to deal with interpretation problems.

Students currently communicate with the system via a typed chat interface. While typing removes the uncertainty and errors involved in speech recognition, expected student answers are considerably more complex and varied than in a typical spoken dialogue system. Therefore, a significant number of interpretation errors arise, primarily during the semantic interpretation process. These errors can lead to *non-understandings*, when the system cannot produce a syntactic parse (or a reasonable fragmentary parse), or when it does not know how to interpret an out-of-domain word; and *misunderstandings*, where a system arrives at an incorrect interpretation, due to either an incorrect attachment in the parse, an incorrect

word sense assigned to an ambiguous word, or an incorrectly resolved referential expression.

Our approach to selecting an error recovery policy is to prefer non-understandings to misunderstandings. There is a known trade-off in spoken dialogue systems between allowing misunderstandings, i.e., cases in which a system accepts and acts on an incorrect interpretation of an utterance, and non-understandings, i.e., cases in which a system rejects an utterance as uninterpretable (Bohus and Rudnicky, 2005). Since misunderstandings on the part of a computer tutor are known to negatively impact student learning, and since in human-human tutorial dialogue the majority of student responses using unexpected terminology are classified as incorrect (Jordan et al., 2009), it would be a reasonable approach for a tutorial dialogue system to deal with potential interpretation problems by treating low-confidence interpretations as non-understandings and focusing on an effective non-understanding recovery policy.¹

We implemented two different policies for comparison. Our baseline policy does not attempt any remediation or error recovery. All student utterances are passed through the standard interpretation pipeline, so that the results can be analyzed later. However, the system does not attempt to address the student content. Instead, regardless of the answer analysis, the system always uses a neutral acceptance and bottom out strategy, giving the student the correct answer every time, e.g., “OK. One way to phrase the correct answer is: the open switch creates a gap in the circuit”. Thus, the students are never given any indication of whether they have been understood or not.

The full policy acts differently depending on the analysis of the student answer. For correct answers, it acknowledges the answer as correct and optionally restates it (see (Dzikovska et al., 2008) for details). For incorrect answers, it restates the correct portion of the answer (if any) and provides a hint to guide the student towards the completely correct answer. If the student’s utterance cannot be interpreted, the system responds with a help message indicating the cause of the problem together with a hint. In both cases, after 3 unsuccessful attempts to address the problem the system uses the bottom out strategy and gives away the answer.

¹While there is no confidence score from a speech recognizer, our system uses a combination of a parse quality score assigned by the parser and a set of consistency checks to determine whether an interpretation is sufficiently reliable.

The content of the bottom out is the same as in the baseline, except that the full system indicates clearly that the answer was incorrect or was not understood, e.g., “Not quite. Here is the answer: the open switch creates a gap in the circuit”.

The help messages are based on the Targeted-Help approach successfully used in spoken dialogue (Hockey et al., 2003), together with the error classification we developed for tutorial dialogue (Dzikovska et al., 2009). There are 9 different error types, each associated with a different targeted help message. The goal of the help messages is to give the student as much information as possible as to why the system failed to understand them but without giving away the answer.

In comparing the two policies, we would expect that the students in both conditions would learn something, but that the learning gain and user satisfaction would be affected by the difference in policies. We hypothesized that students who receive feedback on their errors in the full condition would learn more compared to those in the baseline condition.

3 Evaluation

We collected data from 76 subjects interacting with the system. The subjects were randomly assigned to either the baseline (BASE) or the full (FULL) policy condition. Each subject took a pre-test, then worked through a lesson with the system, and then took a post-test and filled in a user satisfaction survey. Each session lasted approximately 4 hours, with 232 student language turns in FULL ($SD = 25.6$) and 156 in BASE ($SD = 2.02$). Additional time was taken by reading and interacting with the simulation environment. The students had little prior knowledge of the domain. The survey consisted of 63 questions on the 5-point Likert scale covering the lesson content, the graphical user interface, and tutor’s understanding and feedback. For purposes of this study, we are using an averaged tutor score.

The average learning gain was 0.57 ($SD = 0.23$) in FULL, and 0.63 ($SD = 0.26$) in BASE. There was no significant difference in learning gain between conditions. Students liked BASE better: the average tutor evaluation score for FULL was 2.56 out of 5 ($SD = 0.65$), compared to 3.32 ($SD = 0.65$) in BASE. These results are significantly different (t -test, $p < 0.05$). In informal comments after the session many students said that

they were frustrated when the system said that it did not understand them. However, some students in BASE also mentioned that they sometimes were not sure if the system’s answer was correcting a problem with their answer, or simply phrasing it in a different way.

We used mean frequency of non-interpretable utterances (out of all student utterances in each session) to evaluate the effectiveness of the two different policies. On average, 14% of utterances in both conditions resulted in non-understandings.² The frequency of non-understandings was negatively correlated with learning gain in FULL: $r = -0.47, p < 0.005$, but not significantly correlated with learning gain in BASE: $r = -0.09, p = 0.59$. However, in both conditions the frequency of non-understandings was negatively correlated with user satisfaction: FULL $r = -0.36, p = 0.03$, BASE $r = -0.4, p = 0.01$. Thus, even though in BASE the system did not indicate non-understanding, students were negatively affected. That is, they were not satisfied with the policy that did not directly address the interpretation problems. We discuss possible reasons for this below.

We investigated the effect of different types of interpretation errors using two criteria. First, we checked whether the mean frequency of errors was reduced between BASE and FULL for each individual strategy. The reduced frequency means that the recovery strategy for this particular error type is effective in reducing the error frequency. Second, we looked for the cases where the frequency of a given error type is negatively correlated with either learning gain or user satisfaction. This provides evidence that such errors are negatively impacting the learning process, and therefore improving recovery strategies for those error types is likely to improve overall system effectiveness,

The results, shown in Table 1, indicate that the majority of interpretation problems are not significantly correlated with learning gain. However, several types of problems appear to be particularly significant, and are all related to improper use of domain terminology. These were *irrelevant_answer*, *no_appr_terms*, *selective_restriction_failure* and *program_error*.

An *irrelevant_answer* error occurs when the student makes a statement that uses domain termi-

²We do not know the percentage of misunderstandings or concept error rate as yet. We are currently annotating the data with the goal to evaluate interpretation correctness.

error type	full			baseline		
	mean freq. (std. dev)	satisfac- tion r	gain r	mean freq. (std. dev)	satisfac- tion r	gain r
irrelevant_answer	0.008 (0.01)	-0.08	-0.19	0.012 (0.01)	-0.07	-0.47**
no_appr_terms	0.005 (0.01)	-0.57**	-0.42**	0.003 (0.01)	-0.38**	-0.01
selectional_restr_failure	0.032 (0.02)	-0.12	-0.55**	0.040 (0.03)	0.13	0.26*
program_error	0.002 (0.003)	0.02	0.26	0.003 (0.003)	0	-0.35**
unknown_word	0.023 (0.01)	0.05	-0.21	0.024 (0.02)	-0.15	-0.09
disambiguation_failure	0.013 (0.01)	-0.04	0.02	0.007 (0.01)	-0.18	0.19
no_parse	0.019 (0.01)	-0.14	-0.08	0.022(0.02)	-0.3*	0.01
partial_interpretation	0.004 (0.004)	-0.11	-0.01	0.004 (0.005)	-0.19	0.22
reference_failure	0.012 (0.02)	-0.31*	-0.09	0.017 (0.01)	-0.15	-0.23
Overall	0.134 (0.05)	-0.36**	-0.47**	0.139 (0.04)	-0.4**	-0.09

Table 1: Correlations between frequency of different error types and student learning gain and satisfaction. ** - correlation is significant with $p < 0.05$, * - with $p \leq 0.1$.

nology but does not appear to answer the system’s question directly. For example, the expected answer to “In circuit 1, which components are in a closed path?” is “the bulb”. Some students misread the question and say “Circuit 1 is closed.” If that happens, in FULL the system says “Sorry, this isn’t the form of answer that I expected. I am looking for a component”, pointing out to the student the kind of information it is looking for. The BASE system for this error, and for all other errors discussed below, gives away the correct answer without indicating that there was a problem with interpreting the student’s utterance, e.g., “OK, the correct answer is the bulb.”

The *no_appr_terms* error happens when the student is using terminology inappropriate for the lesson in general. Students are expected to learn to explain everything in terms of connections and terminal states. For example, the expected answer to “What is voltage?” is “the difference in states between two terminals”. If instead the student says “Voltage is electricity”, FULL responds with “I am sorry, I am having trouble understanding. I see no domain concepts in your answer. Here’s a hint: your answer should mention a terminal.” The motivation behind this strategy is that in general, it is very difficult to reason about vaguely used domain terminology. We had hoped that by telling the student that the content of their utterance is outside the domain as understood by the system, and hinting at the correct terms to use, the system would guide students towards a better answer.

Selectional_restr_failure errors are typically due to incorrect terminology, when the students phrased answers in a way that contradicted the sys-

tem’s domain knowledge. For example, the system can reason about damaged bulbs and batteries, and open and closed paths. So if the student says “The path is damaged”, the FULL system would respond with “I am sorry, I am having trouble understanding. Paths cannot be damaged. Only bulbs and batteries can be damaged.”

Program_error were caused by faults in the underlying network software, but usually occurred when the student was using extremely long and complicated utterances.

Out of the four important error types described above, only the strategy for *irrelevant_answer* was effective: the frequency of *irrelevant_answer* errors is significantly higher in BASE (t -test, $p < 0.05$), and it is negatively correlated with learning gain in BASE. The frequencies of other error types did not significantly differ between conditions.

However, one other finding is particularly interesting: the frequency of *no_appr_terms* errors is negatively correlated with user satisfaction in BASE. This indicates that simply accepting the student’s answer when they are using incorrect terminology and exposing them to the correct answer is not the best strategy, possibly because the students are noticing the unexplained lack of alignment between their utterance and the system’s answer.

4 Discussion and Future Work

As discussed in Section 1, previous studies of short-answer tutorial dialogue systems produced a counter-intuitive result: measures of interpretation accuracy were not correlated with learning gain. With less restricted language, misunderstandings

negatively affected learning. Our study provides further evidence that interpretation quality significantly affects learning gain in tutorial dialogue. Moreover, while it has long been known that user satisfaction is negatively correlated with interpretation error rates in spoken dialogue, this is the first attempt to evaluate the impact of different types of interpretation errors on task success and usability of a tutoring system.

Our results demonstrate that different types of errors may matter to a different degree. In our system, all of the error types negatively correlated with learning gain stem from the same underlying problem: the use of incorrect or vague terminology by the student. With the exception of the *irrelevant_answer* strategy, the targeted help strategies we implemented were not effective in reducing error frequency or improving learning gain. Additional research is needed to understand why. One possibility is that *irrelevant_answer* was easier to remediate compared to other error types. It usually happened in situations where there was a clear expectation of the answer type (e.g., a list of component names, a yes/no answer). Therefore, it was easier to design an effective prompt. Help messages for other error types were more frequent when the expected answer was a complex sentence, and multiple possible ways of phrasing the correct answer were acceptable. Therefore, it was more difficult to formulate a prompt that would clearly describe the problem in all contexts.

One way to improve the help messages may be to have the system indicate more clearly when user terminology is a problem. Our system apologized each time there was a non-understanding, leading students to believe that they may be answering correctly but the answer is not being understood. A different approach would be to say something like “I am sorry, you are not using the correct terminology in your answer. Here’s a hint: your answer should mention a terminal”. Together with an appropriate mechanism to detect paraphrases of correct answers (as opposed to vague answers whose correctness is difficult to determine), this approach could be more beneficial in helping students learn. We are considering implementing and evaluating this as part of our future work.

Some of the errors, in particular instances of *no_appr_terms* and *selectional_restr_failure*, also stemmed from unrecognized paraphrases with non-standard terminology. Those answers could

conceivably be accepted by a system using semantic similarity as a metric (e.g., using LSA with pre-authored answers). However, our results also indicate that simply accepting the incorrect terminology may not be the best strategy. Users appear to be sensitive when the system’s language does not align with their terminology, as reflected in the decreased satisfaction ratings associated with higher rates of incorrect terminology problems in BASE. Moreover, prior analysis of human-human data indicates that tutors use different restate strategies depending on the “quality” of the student answers, even if they are accepting them as correct (Dzikovska et al., 2008). Together, these point at an important unaddressed issue: existing systems are often built on the assumption that only incorrect and missing parts of the student answer should be remediated, and a wide range of terminology should be accepted (Graesser et al., 1999; Jordan et al., 2006). While it is obviously important for the system to accept a range of different phrasings, our analysis indicates that this may not be sufficient by itself, and students could potentially benefit from addressing the terminology issues with a specifically devised strategy.

Finally, it could also be possible that some differences between strategy effectiveness were caused by incorrect error type classification. Manual examination of several dialogues suggests that most of the errors are assigned to the appropriate type, though in some cases incorrect syntactic parses resulted in unexpected interpretation errors, causing the system to give a confusing help message. These misclassifications appear to be evenly split between different error types, though a more formal evaluation is planned in the future. However from our initial examination, we believe that the differences in strategy effectiveness that we observed are due to the actual differences in the help messages. Therefore, designing better prompts would be the key factor in improving learning and user satisfaction.

Acknowledgments

This work has been supported in part by US Office of Naval Research grants N000140810043 and N0001410WX20278. We thank Katherine Harrison, Leanne Taylor, Charles Callaway, and Elaine Farrow for help with setting up the system and running the evaluation. We would like to thank anonymous reviewers for their detailed feedback.

References

- V. Aleven, O. Popescu, and K. R. Koedinger. 2001. Towards tutorial dialog to support self-explanation: Adding natural language understanding to a cognitive tutor. In *Proceedings of the 10th International Conference on Artificial Intelligence in Education (AIED '01)*.
- Dan Bohus and Alexander Rudnicky. 2005. Sorry, I didn't catch that! - An investigation of non-understanding errors and recovery strategies. In *Proceedings of SIGdial-2005*, Lisbon, Portugal.
- Michelene T. H. Chi, Nicholas de Leeuw, Mei-Hung Chiu, and Christian LaVancher. 1994. Eliciting self-explanations improves understanding. *Cognitive Science*, 18(3):439–477.
- Myroslava O. Dzikovska, Gwendolyn E. Campbell, Charles B. Callaway, Natalie B. Steinhauer, Elaine Farrow, Johanna D. Moore, Leslie A. Butler, and Colin Matheson. 2008. Diagnosing natural language answers to support adaptive tutoring. In *Proceedings 21st International FLAIRS Conference*, Coconut Grove, Florida, May.
- Myroslava O. Dzikovska, Charles B. Callaway, Elaine Farrow, Johanna D. Moore, Natalie B. Steinhauer, and Gwendolyn C. Campbell. 2009. Dealing with interpretation errors in tutorial dialogue. In *Proceedings of SIGDIAL-09*, London, UK, Sep.
- Myroslava O. Dzikovska, Johanna D. Moore, Natalie Steinhauer, Gwendolyn Campbell, Elaine Farrow, and Charles B. Callaway. 2010. Beetle II: a system for tutoring and computational linguistics experimentation. In *Proceedings of ACL-2010 demo session*.
- A. C. Graesser, P. Wiemer-Hastings, P. Wiemer-Hastings, and R. Kreuz. 1999. Autotutor: A simulation of a human tutor. *Cognitive Systems Research*, 1:35–51.
- Beth Ann Hockey, Oliver Lemon, Ellen Campana, Laura Hiatt, Gregory Aist, James Hieronymus, Alexander Gruenstein, and John Dowding. 2003. Targeted help for spoken dialogue systems: intelligent feedback improves naive users' performance. In *Proceedings of the tenth conference on European chapter of the Association for Computational Linguistics*, pages 147–154, Morristown, NJ, USA.
- Pamela W. Jordan, Maxim Makatchev, and Kurt VanLehn. 2004. Combining competing language understanding approaches in an intelligent tutoring system. In James C. Lester, Rosa Maria Vicari, and Fábio Paraguaçu, editors, *Intelligent Tutoring Systems*, volume 3220 of *Lecture Notes in Computer Science*, pages 346–357. Springer.
- Pamela Jordan, Maxim Makatchev, Umarani Pappuswamy, Kurt VanLehn, and Patricia Albacete. 2006. A natural language tutorial dialogue system for physics. In *Proceedings of the 19th International FLAIRS conference*.
- Pamela Jordan, Diane Litman, Michael Lipschultz, and Joanna Drummond. 2009. Evidence of misunderstandings in tutorial dialogue and their impact on learning. In *Proceedings of the 14th International Conference on Artificial Intelligence in Education (AIED)*, Brighton, UK, July.
- Diane Litman and Kate Forbes-Riley. 2005. Speech recognition performance and learning in spoken dialogue tutoring. In *Proceedings of EUROSPEECH-2005*, page 1427.
- Diane Litman, Johanna Moore, Myroslava Dzikovska, and Elaine Farrow. 2009. Generalizing tutorial dialogue results. In *Proceedings of 14th International Conference on Artificial Intelligence in Education (AIED)*, Brighton, UK, July.
- Philip M. McCarthy, Vasile Rus, Scott Crossley, Arthur C. Graesser, and Danielle S. McNamara. 2008. Assessing forward-, reverse-, and average-entailment indices on natural language input from the intelligent tutoring system, iSTART. In *Proceedings of the 21st International FLAIRS conference*, pages 165–170.
- Rodney D. Nielsen, Wayne Ward, and James H. Martin. 2008. Learning to assess low-level conceptual understanding. In *Proceedings 21st International FLAIRS Conference*, Coconut Grove, Florida, May.
- Heather Pon-Barry, Brady Clark, Elizabeth Owen Bratt, Karl Schultz, and Stanley Peters. 2004. Evaluating the effectiveness of SCoT: A spoken conversational tutor. In J. Mostow and P. Tedesco, editors, *Proceedings of the ITS 2004 Workshop on Dialog-based Intelligent Tutoring Systems*, pages 23–32.
- Amruta Purandare and Diane Litman. 2008. Content-learning correlations in spoken tutoring dialogs at word, turn and discourse levels. In *Proceedings 21st International FLAIRS Conference*, Coconut Grove, Florida, May.

The Prevalence of Descriptive Referring Expressions in News and Narrative

Raquel Hervás

Departamento de Ingeniería
del Software e Inteligencia Artificial
Universidad Complutense de Madrid
Madrid, 28040 Spain
raquelhb@fdi.ucm.es

Mark Alan Finlayson

Computer Science and
Artificial Intelligence Laboratory
Massachusetts Institute of Technology
Cambridge, MA, 02139 USA
markaf@mit.edu

Abstract

Generating referring expressions is a key step in Natural Language Generation. Researchers have focused almost exclusively on generating *distinctive* referring expressions, that is, referring expressions that uniquely identify their intended referent. While undoubtedly one of their most important functions, referring expressions can be more than distinctive. In particular, *descriptive* referring expressions – those that provide additional information not required for distinction – are critical to fluent, efficient, well-written text. We present a corpus analysis in which approximately one-fifth of 7,207 referring expressions in 24,422 words of news and narrative are descriptive. These data show that if we are ever to fully master natural language generation, especially for the genres of news and narrative, researchers will need to devote more attention to understanding how to generate descriptive, and not just distinctive, referring expressions.

1 A Distinctive Focus

Generating referring expressions is a key step in Natural Language Generation (NLG). From early treatments in seminal papers by Appelt (1985) and Reiter and Dale (1992) to the recent set of Referring Expression Generation (REG) Challenges (Gatt et al., 2009) through different corpora available for the community (Eugenio et al., 1998; van Deemter et al., 2006; Viethen and Dale, 2008), generating referring expressions has become one of the most studied areas of NLG.

Researchers studying this area have, almost without exception, focused exclusively on how to generate *distinctive* referring expressions, that is, referring expressions that unambiguously iden-

tify their intended referent. Referring expressions, however, may be more than distinctive. It is widely acknowledged that they can be used to achieve multiple goals, above and beyond distinction. Here we focus on *descriptive* referring expressions, that is, referring expressions that are not only distinctive, but provide additional information not required for identifying their intended referent. Consider the following text, in which some of the referring expressions have been underlined:

Once upon a time there was a man, who had three daughters. They lived in a house and their dresses were made of fabric.

While a bit strange, the text is perfectly well-formed. All the referring expressions are distinctive, in that we can properly identify the referents of each expression. But the real text, the opening lines to the folktale *The Beauty and the Beast*, is actually much more lyrical:

*Once upon a time there was **a rich merchant**, who had three daughters. They lived in **a very fine house** and their **gowns** were made of **the richest fabric sewn with jewels**.*

All the boldfaced portions – namely, the choice of head nouns, the addition of adjectives, the use of appositive phrases – serve to perform a descriptive function, and, importantly, are all unnecessary for distinction! In all of these cases, the author is using the referring expressions as a vehicle for communicating information about the referents. This descriptive information is sometimes new, sometimes necessary for understanding the text, and sometimes just for added flavor. But when the expression is *descriptive*, as opposed to *distinctive*, this additional information is not required for identifying the referent of the expression, and it is these sorts of referring expressions that we will be concerned with here.

Although these sorts of referring expression have been mostly ignored by researchers in this area¹, we show in this corpus study that descriptive expressions are in fact quite prevalent: nearly one-fifth of referring expressions in news and narrative are descriptive. In particular, our data, the trained judgments of native English speakers, show that 18% of all distinctive referring expressions in news and 17% of those in narrative folktales are descriptive. With this as motivation, we argue that descriptive referring expressions must be studied more carefully, especially as the field progresses from referring in a physical, immediate context (like that in the REG Challenges) to generating more literary forms of text.

2 Corpus Annotation

This is a corpus study; our procedure was therefore to define our annotation guidelines (Section 2.1), select texts to annotate (2.2), create an annotation tool for our annotators (2.3), and, finally, train annotators, have them annotate referring expressions' constituents and function, and then adjudicate the double-annotated texts into a gold standard (2.4).

2.1 Definitions

We wrote an annotation guide explaining the difference between distinctive and descriptive referring expressions. We used the guide when training annotators, and it was available to them while annotating. With limited space here we can only give an outline of what is contained in the guide; for full details see (Finlayson and Hervás, 2010a).

Referring Expressions We defined referring expressions as referential noun phrases and their coreferential expressions, e.g., “John kissed Mary. She blushed.”. This included referring expressions to generics (e.g., “Lions are fierce”), dates, times, and numbers, as well as events if they were referred to using a noun phrase. We included in each referring expression all the determiners, quantifiers, adjectives, appositives, and prepositional phrases that syntactically attached to that expression. When referring expressions were nested, all the nested referring expressions were also marked separately.

Nuclei vs. Modifiers In the only previous corpus study of descriptive referring expressions, on

¹With the exception of a small amount of work, discussed in Section 4.

museum labels, Cheng et al. (2001) noted that descriptive information is often integrated into referring expressions using modifiers to the head noun. To study this, and to allow our results to be more closely compared with Cheng's, we had our annotators split referring expressions into their constituents, portions called either *nuclei* or *modifiers*. The nuclei were the portions of the referring expression that performed the ‘core’ referring function; the modifiers were those portions that could be varied, syntactically speaking, independently of the nuclei. Annotators then assigned a distinctive or descriptive function to each constituent, rather than the referring expression as a whole.

Normally, the nuclei corresponded to the head of the noun phrase. In (1), the nucleus is the token *king*, which we have here surrounded with square brackets. The modifiers, surrounded by parentheses, are *The* and *old*.

(1) *(The) (old) [king] was wise.*

Phrasal modifiers were marked as single modifiers, for example, in (2).

(2) *(The) [roof] (of the house) collapsed.*

It is significant that we had our annotators mark and tag the nuclei of referring expressions. Cheng and colleagues only mentioned the possibility that additional information could be introduced in the modifiers. However, O'Donnell et al. (1998) observed that often the choice of head noun can also influence the function of a referring expression. Consider (3), in which the word *villain* is used to refer to the King.

(3) *The King assumed the throne today.
I don't trust (that) [villain] one bit.*

The speaker could have merely used *him* to refer to the King—the choice of that particular head noun *villain* gives us additional information about the disposition of the speaker. Thus *villain* is descriptive.

Function: Distinctive vs. Descriptive As already noted, instead of tagging the whole referring expression, annotators tagged each constituent (nuclei and modifiers) as distinctive or descriptive.

The two main tests for determining descriptiveness were (a) if presence of the constituent was unnecessary for identifying the referent, or (b) if

the constituent was expressed using unusual or ostentatious word choice. If either was true, the constituent was considered descriptive; otherwise, it was tagged as distinctive. In cases where the constituent was completely irrelevant to identifying the referent, it was tagged as descriptive. For example, in the folktale *The Princess and the Pea*, from which (1) was extracted, there is only one king in the entire story. Thus, in that story, *the king* is sufficient for identification, and therefore the modifier *old* is descriptive. This points out the importance of context in determining distinctiveness or descriptiveness; if there had been a roomful of kings, the tags on those modifiers would have been reversed.

There is some question as to whether copular predicates, such as *the plumber* in (4), are actually referring expressions.

(4) *John is the plumber*

Our annotators marked and tagged these constructions as normal referring expressions, but they added an additional flag to identify them as copular predicates. We then excluded these constructions from our final analysis. Note that copular predicates were treated differently from appositives: in appositives the predicate was included in the referring expression, and in most cases (again, depending on context) was marked descriptive (e.g., *John, the plumber, slept*).

2.2 Text Selection

Our corpus comprised 62 texts, all originally written in English, from two different genres, news and folktales. We began with 30 folktales of different sizes, totaling 12,050 words. These texts were used in a previous work on the influence of dialogues on anaphora resolution algorithms (Aggarwal et al., 2009); they were assembled with an eye toward including different styles, different authors, and different time periods. Following this, we matched, approximately, the number of words in the folktales by selecting 32 texts from Wall Street Journal section of the Penn Treebank (Marcus et al., 1993). These texts were selected at random from the first 200 texts in the corpus.

2.3 The Story Workbench

We used the Story Workbench application (Finlayson, 2008) to actually perform the annotation. The Story Workbench is a semantic annotation

program that, among other things, includes the ability to annotate referring expressions and coreferential relationships. We added the ability to annotate nuclei, modifiers, and their functions by writing a workbench “plugin” in Java that could be installed in the application.

The Story Workbench is not yet available to the public at large, being in a limited distribution beta testing phase. The developers plan to release it as free software within the next year. At that time, we also plan to release our plugin as free, downloadable software.

2.4 Annotation & Adjudication

The main task of the study was the annotation of the constituents of each referring expression, as well as the function (distinctive or descriptive) of each constituent. The system generated a first pass of constituent analysis, but did not mark functions. We hired two native English annotators, neither of whom had any linguistics background, who corrected these automatically-generated constituent analyses, and tagged each constituent as descriptive or distinctive. Every text was annotated by both annotators. Adjudication of the differences was conducted by discussion between the two annotators; the second author moderated these discussions and settled irreconcilable disagreements. We followed a “train-as-you-go” paradigm, where there was no distinct training period, but rather adjudication proceeded in step with annotation, and annotators received feedback during those sessions.

We calculated two measures of inter-annotator agreement: a kappa statistic and an f-measure, shown in Table 1. All of our f-measures indicated that annotators agreed almost perfectly on the location of referring expressions and their breakdown into constituents. These agreement calculations were performed on the annotators’ original corrected texts.

All the kappa statistics were calculated for two tags (nuclei vs. modifier for the constituents, and distinctive vs. descriptive for the functions) over both each token assigned to a nucleus or modifier and each referring expression pair. Our kappas indicate moderate to good agreement, especially for the folktales. These results are expected because of the inherent subjectivity of language. During the adjudication sessions it became clear that different people do not consider the same information

as obvious or descriptive for the same concepts, and even the contexts deduced by each annotators from the texts were sometimes substantially different.

	Tales	Articles	Total
Ref. Exp. (F_1)	1.00	0.99	0.99
Constituents (F_1)	0.99	0.98	0.98
Nuc./Mod. (κ)	0.97	0.95	0.96
Const. Func. (κ)	0.61	0.48	0.54
Ref. Exp. Func. (κ)	0.65	0.54	0.59

Table 1: Inter-annotator agreement measures

3 Results

Table 2 lists the primary results of the study. We considered a referring expression descriptive if any of its constituents were descriptive. Thus, 18% of the referring expressions in the corpus added additional information beyond what was required to unambiguously identify their referent. The results were similar in both genres.

	Tales	Articles	Total
Texts	30	32	62
Words	12,050	12,372	24,422
Sentences	904	571	1,475
Ref. Exp.	3,681	3,526	7,207
Dist. Ref. Exp.	3,057	2,830	5,887
Desc. Ref. Exp.	609	672	1,281
% Dist. Ref.	83%	81%	82%
% Desc. Ref.	17%	19%	18%

Table 2: Primary results.

Table 3 contains the percentages of descriptive and distinctive tags broken down by constituent. Like Cheng’s results, our analysis shows that descriptive referring expressions make up a significant fraction of all referring expressions. Although Cheng did not examine nuclei, our results show that the use of descriptive nuclei is small but not negligible.

4 Relation to the Field

Researchers working on generating referring expressions typically acknowledge that referring expressions can perform functions other than distinction. Despite this widespread acknowledgment, researchers have, for the most part, explicitly ignored these functions. Exceptions to this trend

	Tales	Articles	Total
Nuclei	3,666	3,502	7,168
Max. Nuc/Ref	1	1	1
Dist. Nuc.	95%	97%	96%
Desc. Nuc.	5%	3%	4%
Modifiers	2,277	3,627	5,904
Avg. Mod/Ref	0.6	1.0	0.8
Max. Mod/Ref	4	6	6
Dist. Mod.	78%	81%	80%
Desc. Mod.	22%	19%	20%

Table 3: Breakdown of Constituent Tags

are three. First is the general study of *aggregation* in the process of referring expression generation. Second and third are corpus studies by Cheng et al. (2001) and Jordan (2000a) that bear on the prevalence of descriptive referring expressions.

The NLG subtask of aggregation can be used to imbue referring expressions with a descriptive function (Reiter and Dale, 2000, §5.3). There is a specific kind of aggregation called *embedding* that moves information from one clause to another inside the structure of a separate noun phrase. This type of aggregation can be used to transform two sentences such as “*The princess lived in a castle. She was pretty*” into “*The pretty princess lived in a castle*”. The adjective *pretty*, previously a copular predicate, becomes a descriptive modifier of the reference to the princess, making the second text more natural and fluent. This kind of aggregation is widely used by humans for making the discourse more compact and efficient. In order to create NLG systems with this ability, we must take into account the caveat, noted by Cheng (1998), that any non-distinctive information in a referring expression must not lead to confusion about the distinctive function of the referring expression. This is by no means a trivial problem – this sort of aggregation interferes with referring and coherence planning at both a local and global level (Cheng and Mellish, 2000; Cheng et al., 2001). It is clear, from the current state of the art of NLG, that we have not yet obtained a deep enough understanding of aggregation to enable us to handle these interactions. More research on the topic is needed.

Two previous corpus studies have looked at the use of descriptive referring expressions. The first showed explicitly that people craft descriptive referring expressions to accomplish different

goals. Jordan and colleagues (Jordan, 2000b; Jordan, 2000a) examined the use of referring expressions using the COCONUT corpus (Eugenio et al., 1998). They tested how domain and discourse goals can influence the content of non-pronominal referring expressions in a dialogue context, checking whether or not a subject’s goals led them to include non-referring information in a referring expression. Their results are intriguing because they point toward heretofore unexamined constraints, utilities and expectations (possibly genre- or style-dependent) that may underlie the use of descriptive information to perform different functions, and are not yet captured by aggregation modules in particular or NLG systems in general.

In the other corpus study, which partially inspired this work, Cheng and colleagues analyzed a set of museum descriptions, the GNOME corpus (Poesio, 2004), for the pragmatic functions of referring expressions. They had three functions in their study, in contrast to our two. Their first function (marked by their `uniq` tag) was equivalent to our distinctive function. The other two were specializations of our descriptive tag, where they differentiated between additional information that helped to understand the text (`int`), or additional information not necessary for understanding (`attr`). Despite their annotators seeming to have trouble distinguishing between the latter two tags, they did achieve good overall inter-annotator agreement. They identified 1,863 modifiers to referring expressions in their corpus, of which 47.3% fulfilled a descriptive (`attr` or `int`) function. This is supportive of our main assertion, namely, that descriptive referring expressions, not only crucial for efficient and fluent text, are actually a significant phenomenon. It is interesting, though, that Cheng’s fraction of descriptive referring expression was so much higher than ours (47.3% versus our 18%). We attribute this substantial difference to genre, in that Cheng studied museum labels, in which the writer is space-constrained, having to pack a lot of information into a small label. The issue bears further study, and perhaps will lead to insights into differences in writing style that may be attributed to author or genre.

5 Contributions

We make two contributions in this paper.

First, we assembled, double-annotated, and ad-

judicated into a gold-standard a corpus of 24,422 words. We marked all referring expressions, coreferential relations, and referring expression constituents, and tagged each constituent as having a descriptive or distinctive function. We wrote an annotation guide and created software that allows the annotation of this information in free text. The corpus and the guide are available on-line in a permanent digital archive (Finlayson and Hervás, 2010a; Finlayson and Hervás, 2010b). The software will also be released in the same archive when the Story Workbench annotation application is released to the public. This corpus will be useful for the automatic generation and analysis of both descriptive and distinctive referring expressions. Any kind of system intended to generate text as humans do must take into account that identification is not the only function of referring expressions. Many analysis applications would benefit from the automatic recognition of descriptive referring expressions.

Second, we demonstrated that descriptive referring expressions comprise a substantial fraction (18%) of the referring expressions in news and narrative. Along with museum descriptions, studied by Cheng, it seems that news and narrative are genres where authors naturally use a large number of descriptive referring expressions. Given that so little work has been done on descriptive referring expressions, this indicates that the field would be well served by focusing more attention on this phenomenon.

Acknowledgments

This work was supported in part by the Air Force Office of Scientific Research under grant number A9550-05-1-0321, as well as by the Office of Naval Research under award number N00014091059. Any opinions, findings, and conclusions or recommendations expressed in this paper are those of the authors and do not necessarily reflect the views of the Office of Naval Research. This research is also partially funded the Spanish Ministry of Education and Science (TIN2009-14659-C03-01) and Universidad Complutense de Madrid (GR58/08). We also thank Whitman Richards, Ozlem Uzuner, Peter Szolovits, Patrick Winston, Pablo Gervás, and Mark Seifter for their helpful comments and discussion, and thank our annotators Saam Batmanghelidj and Geneva Trotter.

References

- Alaukik Aggarwal, Pablo Gervás, and Raquel Hervás. 2009. Measuring the influence of errors induced by the presence of dialogues in reference clustering of narrative text. In *Proceedings of ICON-2009: 7th International Conference on Natural Language Processing*, India. Macmillan Publishers.
- Douglas E. Appelt. 1985. Planning English referring expressions. *Artificial Intelligence*, 26:1–33.
- Hua Cheng and Chris Mellish. 2000. Capturing the interaction between aggregation and text planning in two generation systems. In *INLG '00: First international conference on Natural Language Generation 2000*, pages 186–193, Morristown, NJ, USA. Association for Computational Linguistics.
- Hua Cheng, Massimo Poesio, Renate Henschel, and Chris Mellish. 2001. Corpus-based np modifier generation. In *NAACL '01: Second meeting of the North American Chapter of the Association for Computational Linguistics on Language technologies 2001*, pages 1–8, Morristown, NJ, USA. Association for Computational Linguistics.
- Hua Cheng. 1998. Embedding new information into referring expressions. In *ACL-36: Proceedings of the 36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics*, pages 1478–1480, Morristown, NJ, USA. Association for Computational Linguistics.
- Barbara Di Eugenio, Johanna D. Moore, Pamela W. Jordan, and Richmond H. Thomason. 1998. An empirical investigation of proposals in collaborative dialogues. In *Proceedings of the 17th international conference on Computational linguistics*, pages 325–329, Morristown, NJ, USA. Association for Computational Linguistics.
- Mark A. Finlayson and Raquel Hervás. 2010a. Annotation guide for the UCM/MIT indications, referring expressions, and coreference corpus (UMIREC corpus). Technical Report MIT-CSAIL-TR-2010-025, MIT Computer Science and Artificial Intelligence Laboratory. <http://hdl.handle.net/1721.1/54765>.
- Mark A. Finlayson and Raquel Hervás. 2010b. UCM/MIT indications, referring expressions, and coreference corpus (UMIREC corpus). Work product, MIT Computer Science and Artificial Intelligence Laboratory. <http://hdl.handle.net/1721.1/54766>.
- Mark A. Finlayson. 2008. Collecting semantics in the wild: The Story Workbench. In *Proceedings of the AAAI Fall Symposium on Naturally-Inspired Artificial Intelligence*, pages 46–53, Menlo Park, CA, USA. AAAI Press.
- Albert Gatt, Anja Belz, and Eric Kow. 2009. The TUNA-REG challenge 2009: overview and evaluation results. In *ENLG '09: Proceedings of the 12th European Workshop on Natural Language Generation*, pages 174–182, Morristown, NJ, USA. Association for Computational Linguistics.
- Pamela W. Jordan. 2000a. Can nominal expressions achieve multiple goals?: an empirical study. In *ACL '00: Proceedings of the 38th Annual Meeting on Association for Computational Linguistics*, pages 142–149, Morristown, NJ, USA. Association for Computational Linguistics.
- Pamela W. Jordan. 2000b. Influences on attribute selection in redescription: A corpus study. In *Proceedings of CogSci2000*, pages 250–255.
- Mitchell P. Marcus, Mary Ann Marcinkiewicz, and Beatrice Santorini. 1993. Building a large annotated corpus of english: the penn treebank. *Computational Linguistics*, 19(2):313–330.
- Michael O'Donnell, Hua Cheng, and Janet Hitzeman. 1998. Integrating referring and informing in NP planning. In *Proceedings of COLING-ACL'98 Workshop on the Computational Treatment of Nominals*, pages 46–56.
- Massimo Poesio. 2004. Discourse annotation and semantic annotation in the GNOME corpus. In *DiscAnnotation '04: Proceedings of the 2004 ACL Workshop on Discourse Annotation*, pages 72–79, Morristown, NJ, USA. Association for Computational Linguistics.
- Ehud Reiter and Robert Dale. 1992. A fast algorithm for the generation of referring expressions. In *Proceedings of the 14th conference on Computational linguistics*, Nantes, France.
- Ehud Reiter and Robert Dale. 2000. *Building Natural Language Generation Systems*. Cambridge University Press.
- Kees van Deemter, Ielka van der Sluis, and Albert Gatt. 2006. Building a semantically transparent corpus for the generation of referring expressions. In *Proceedings of the 4th International Conference on Natural Language Generation (Special Session on Data Sharing and Evaluation)*, INLG-06.
- Jette Viethen and Robert Dale. 2008. The use of spatial relations in referring expressions. In *Proceedings of the 5th International Conference on Natural Language Generation*.

Preferences versus Adaptation during Referring Expression Generation

Martijn Goudbeek
University of Tilburg
Tilburg, The Netherlands
m.b.goudbeek@uvt.nl

Emiel Krahmer
University of Tilburg
Tilburg, The Netherlands
e.j.krahmer@uvt.nl

Abstract

Current Referring Expression Generation algorithms rely on domain dependent preferences for both content selection and linguistic realization. We present two experiments showing that human speakers may opt for dispreferred properties and dispreferred modifier orderings when these were salient in a preceding interaction (without speakers being consciously aware of this). We discuss the impact of these findings for current generation algorithms.

1 Introduction

The generation of referring expressions is a core ingredient of most Natural Language Generation (NLG) systems (Reiter and Dale, 2000; Mellish et al., 2006). These systems usually approach Referring Expression Generation (REG) as a two-step procedure, where first it is decided which properties to include (content selection), after which the selected properties are turned into a natural language referring expression (linguistic realization). The basic problem in both stages is one of choice; there are many ways in which one could refer to a target object and there are multiple ways in which these could be realized in natural language. Typically, these choice problems are tackled by giving preference to some solutions over others. For example, the Incremental Algorithm (Dale and Reiter, 1995), one of the most widely used REG algorithms, assumes that certain attributes are preferred over others, partly based on evidence provided by Pechmann (1989); a chair would first be described in terms of its color, and only if this does not result in a unique characterization, other, less preferred attributes such as orientation are tried. The Incremental Algorithm is arguably unique in assuming a complete preference order of attributes, but other REG algo-

rithms rely on similar distinctions. The Graph-based algorithm (Krahmer et al., 2003), for example, searches for the cheapest description for a target, and distinguishes cheap attributes (such as color) from more expensive ones (orientation). Realization of referring expressions has received less attention, yet recent studies on the ordering of modifiers (Shaw and Hatzivassiloglou, 1999; Malouf, 2000; Mitchell, 2009) also work from the assumption that some orderings (*large red*) are preferred over others (*red large*).

We argue that such preferences are less stable when referring expressions are generated in interactive settings, as would be required for applications such as spoken dialogue systems or interactive virtual characters. In these cases, we hypothesize that, besides domain preferences, also the referring expressions that were produced earlier in the interaction are important. It has been shown that if one dialogue participant refers to a couch as a *sofa*, the next speaker is more likely to use the word *sofa* as well (Branigan et al., in press). This kind of micro-planning or “lexical entrainment” (Brennan and Clark, 1996) can be seen as a specific form of “alignment” (Pickering and Garrod, 2004) between speaker and addressee. Pickering and Garrod argue that alignment may take place on all levels of interaction, and indeed it has been shown that participants also align their intonation patterns and syntactic structures. However, as far as we know, experimental evidence for alignment on the level of content planning has never been given, and neither have alignment effects in modifier orderings during realization been shown. With a few notable exceptions, such as Buschmeier et al. (2009) who study alignment in micro-planning, and Janarthanam and Lemon (2009) who study alignment in expertise levels, alignment has received little attention in NLG so far.

This paper is organized as follows. Experiment I studies the trade-off between adaptation

and preferences during content selection while Experiment II looks at this trade-off for modifier orderings during realization. Both studies use a novel interactive reference production paradigm, applied to two domains – the Furniture and People domains of the TUNA data-set (Gatt et al., 2007; Koolen et al., 2009) – to see whether adaptation may be domain dependent. Finally, we contrast our findings with the performance of state-of-the-art REG algorithms, discussing how they could be adapted so as to account for the new data, effectively adding plasticity to the generation process.

2 Experiment I

Experiment I studies what speakers do when referring to a target that can be distinguished in a preferred (*the blue fan*) or a dispreferred way (*the left-facing fan*), when in the prior context either the first or the second variant was made salient.

Method

Participants 26 students (2 male, mean age = 20 years, 11 months), all native speakers of Dutch without hearing or speech problems, participated for course credits.

Materials Target pictures were taken from the TUNA corpus (Gatt et al., 2007) that has been extensively used for REG evaluation. This corpus consists of two domains: one containing pictures of people (famous mathematicians), the other containing furniture items in different colors depicted from different orientations. From previous studies (Gatt et al., 2007; Koolen et al., 2009) it is known that participants show a preference for certain attributes: color in the Furniture domain and glasses in the People domain, and disprefer other attributes (orientation of a furniture piece and wearing a tie, respectively).

Procedure Trials consisted of four turns in an interactive reference understanding and production experiment: a prime, two fillers and the experimental description (see Figure 1). First, participants listened to a pre-recorded female voice referring to one of three objects and had to indicate which one was being referenced. In this sub-task, references either used a preferred or a dispreferred attribute; both were distinguishing. Second, participants themselves described a filler picture, after which, third, they had to indicate which filler picture was being described. The two filler turns always concerned stimuli from the alterna-

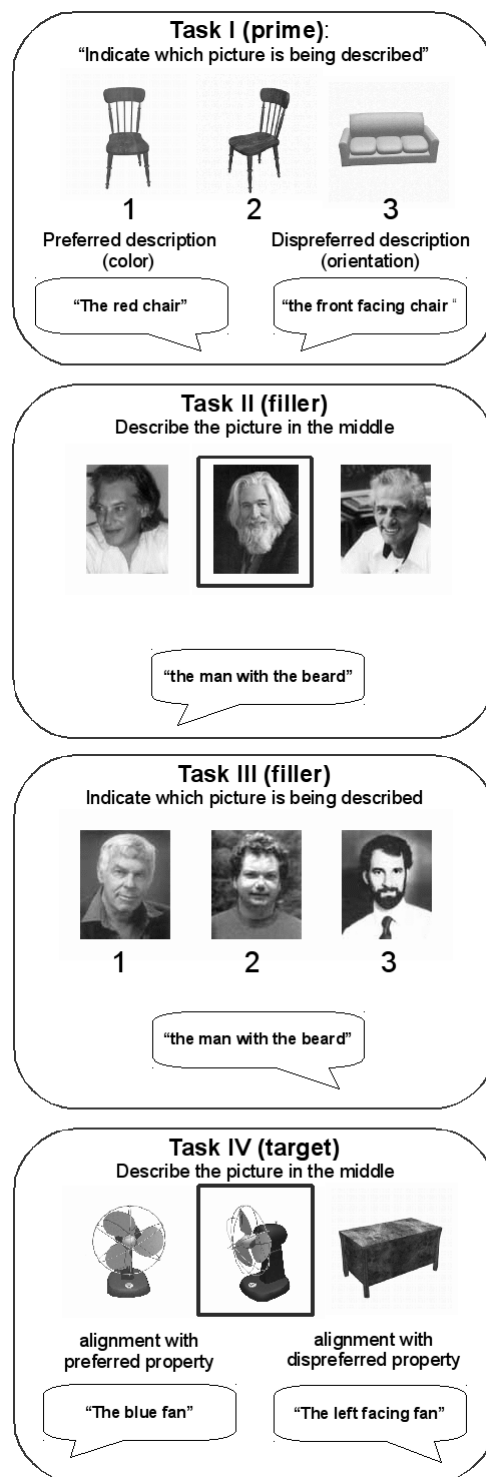


Figure 1: The 4 tasks per trial. A furniture trial is shown; people trials have an identical structure.

tive domain and were intended to prevent a too direct connection between the prime and the target. Fourth, participants described the target object, which could always be distinguished from its distractors in a preferred (*The blue fan*) or a dispreferred (*The left facing fan*) way. Note that at-

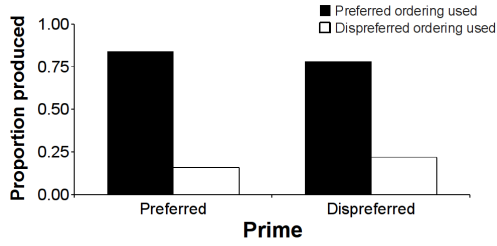


Figure 2: Proportions of preferred and dispreferred attributes in the Furniture domain.

tributes are primed, not values; a participant may have heard *front facing* in the prime turn, while the target has a different value for this attribute (cf. Fig. 1).

For the two domains, there were 20 preferred and 20 dispreferred trials, giving rise to $2 \times (20 + 20) = 80$ critical trials. These were presented in counter-balanced blocks, and within blocks each participant received a different random order. In addition, there were 80 filler trials (each following the same structure as outlined in Figure 1). During debriefing, none of the participants indicated they had been aware of the experiment’s purpose.

Results

We use the proportion of attribute alignment as our dependent measure. Alignment occurs when a participant uses the same attribute in the target as occurred in the prime. This includes overspecified descriptions (Engelhardt et al., 2006; Arnold, 2008), where both the preferred and dispreferred attributes were mentioned by participants. Over-specification occurred in 13% of the critical trials (and these were evenly distributed over the experimental conditions).

The use of the preferred and dispreferred attribute as a function of prime and domain is shown in Figure 2 and Figure 3. In both domains, the preferred attribute is used much more frequently than the dispreferred attribute with the preferred primes, which serves as a manipulation check. As a test of our hypothesis that adaptation processes play an important role in attribute selection for referring expressions, we need to look at participants’ expressions with the *dispreferred* primes (with the preferred primes, effects of adaptation and of preferences cannot be teased apart). Current REG algorithms such as the Incremental Algorithm and the Graph-based algorithm predict that participants will always opt for the preferred

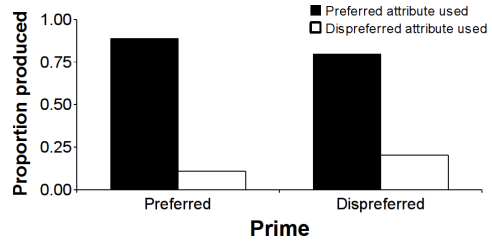


Figure 3: Proportions of preferred and dispreferred attributes in the People domain.

attribute, and hence will not use the dispreferred attribute. This is not what we observe: our participants used the dispreferred attribute at a rate significantly larger than zero when they had been exposed to it three turns earlier ($t_{furniture} [25] = 6.64, p < 0.01$; $t_{people} [25] = 4.78, p < 0.01$). Additionally, they used the dispreferred attribute significantly *more* when they had previously heard the dispreferred attribute rather than the preferred attribute. This difference is especially marked and significant in the Furniture domain ($t_{furniture} [25] = 2.63, p < 0.01, t_{people} [25] = 0.98, p < 0.34$), where participants opt for the dispreferred attribute in 54% of the trials, more frequently than they do for the preferred attribute (Fig. 2).

3 Experiment II

Experiment II uses the same paradigm used for Experiment I to study whether speaker’s preferences for modifier orderings can be changed by exposing them to dispreferred orderings.

Method

Participants 28 Students (ten males, mean age = 23 years and two months) participated for course credits. All were native speakers of Dutch, without hearing and speech problems. None participated in Experiment I.

Materials The materials were identical to those used in Experiment I, except for their arrangement in the critical trials. In these trials, the participants could only identify the target picture using two attributes. In the Furniture domain these were color and size, in the People domain these were having a beard and wearing glasses. In the prime turn (Task I, Fig. 1), these attributes were realized in a preferred way (“size first”: e.g., *the big red sofa*, or “glasses first”: *the bespectacled and bearded man*) or in a dispreferred way (“color first”: *the red big sofa* or “beard first” *the bespectacled and bearded*

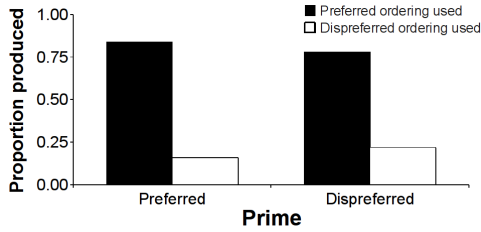


Figure 4: Proportions of preferred and dispreferred modifier orderings in the Furniture domain.

man). Google counts for the original Dutch modifier orderings reveal that the ratio of preferred to dispreferred is in the order of 40:1 in the Furniture domain and 3:1 in the People domain.

Procedure As above.

Results

We use the proportion of modifier ordering alignments as our dependent measure, where alignment occurs when the participant’s ordering coincides with the primed ordering. Figure 4 and 5 show the use of the preferred and dispreferred modifier ordering per prime and domain. It can be seen that in the preferred prime conditions, participants produce the expected orderings, more or less in accordance with the Google counts.

State-of-the-art realizers would always opt for the most frequent ordering of a given pair of modifiers and hence would never predict the dispreferred orderings to occur. Still, the use of the dispreferred modifier ordering occurred significantly more often than one would expect given this prediction, $t_{furniture} [27] = 6.56, p < 0.01$ and $t_{people} [27] = 9.55, p < 0.01$. To test our hypotheses concerning adaptation, we looked at the dispreferred realizations when speakers were exposed to dispreferred primes (compared to preferred primes). In both domains this resulted in an increase of the amount of dispreferred realizations, which was significant in the People domain ($t_{people} [27] = 1.99, p < 0.05$, $t_{furniture} [25] = 2.63, p < 0.01$).

4 Discussion

Current state-of-the-art REG algorithms often rest upon the assumption that some attributes and some realizations are preferred over others. The two experiments described in this paper show that this assumption is incorrect, when references are produced in an interactive setting. In both experiments, speakers were more likely to select a dis-

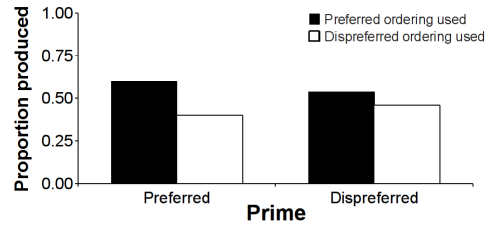


Figure 5: Proportions of preferred and dispreferred modifier orderings in the People domain.

preferred attribute or produce a dispreferred modifier ordering when they had previously been exposed to these attributes or orderings, without being aware of this. These findings fit in well with the adaptation and alignment models proposed by psycholinguists, but ours, as far as we know, is the first experimental evidence of alignment in attribute selection and in modifier ordering. Interestingly, we found that effect sizes differ for the different domains, indicating that the trade-off between preferences and adaptations is a gradual one, also influenced by the *a priori* differences in preference (it is more difficult to make people say something truly dispreferred than something more marginally dispreferred).

To account for these findings, GRE algorithms that function in an interactive setting should be made sensitive to the production of dialogue partners. For the Incremental Algorithm (Dale and Reiter, 1995), this could be achieved by augmenting the list of preferred attributes with a list of “previously mentioned” attributes. The relative weighting of these two lists will be corpus dependent, and can be estimated in a data-driven way. Alternatively, in the Graph-based algorithm (Krahmer et al., 2003), costs of properties could be based on two components: a relatively fixed domain component (preferred is cheaper) and a flexible interactive component (recently used is cheaper). Which approach would work best is an open, empirical question, but either way this would constitute an important step towards interactive REG.

Acknowledgments

The research reported in this paper forms part of the VICI project “Bridging the gap between psycholinguistics and Computational linguistics: the case of referring expressions”, funded by the Netherlands Organization for Scientific Research (NWO grant 277-70-007).

References

- Jennifer Arnold. 2008. Reference production: Production-internal and addressee-oriented processes. *Language and Cognitive Processes*, 23(4):495–527.
- Holly P. Branigan, Martin J. Pickering, Jamie Pearson, and Janet F. McLean. in press. Linguistic alignment between people and computers. *Journal of Pragmatics*, 23:1–2.
- Susan E. Brennan and Herbert H. Clark. 1996. Conceptual pacts and lexical choice in conversation. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 22:1482–1493.
- Hendrik Buschmeier, Kirsten Bergmann, and Stefan Kopp. 2009. An alignment-capable microplanner for Natural Language Generation. In *Proceedings of the 12th European Workshop on Natural Language Generation (ENLG 2009)*, pages 82–89, Athens, Greece, March. Association for Computational Linguistics.
- Robert Dale and Ehud Reiter. 1995. Computational interpretations of the gricean maxims in the generation of referring expressions. *Cognitive Science*, 19(2):233–263.
- Paul E. Engelhardt, Karl G. Bailey, and Fernanda Ferreira. 2006. Do speakers and listeners observe the gricean maxim of quantity? *Journal of Memory and Language*, 54(4):554–573.
- Albert Gatt, Ielka van der Sluis, and Kees van Deemter. 2007. Evaluating algorithms for the generation of referring expressions using a balanced corpus. In *Proceedings of the 11th European Workshop on Natural Language Generation*.
- Srinivasan Janarthanam and Oliver Lemon. 2009. Learning lexical alignment policies for generating referring expressions for spoken dialogue systems. In *Proceedings of the 12th European Workshop on Natural Language Generation (ENLG 2009)*, pages 74–81, Athens, Greece, March. Association for Computational Linguistics.
- Ruud Koolen, Albert Gatt, Martijn Goudbeek, and Emiel Krahmer. 2009. Need I say more? on factors causing referential overspecification. In *Proceedings of the PRE-CogSci 2009 Workshop on the Production of Referring Expressions: Bridging the Gap Between Computational and Empirical Approaches to Reference*.
- Emiel Krahmer, Sebastiaan van Erk, and André Verleg. 2003. Graph-based generation of referring expressions. *Computational Linguistics*, 29(1):53–72.
- Robert Malouf. 2000. The order of pronominal adjectives in natural language generation. In *Proceedings of the 38th Annual Meeting of the Association for Computational Linguistics*, pages 85–92.
- Chris Mellish, Donia Scott, Lynn Cahill, Daniel Paiva, Roger Evans, and Mike Reape. 2006. A reference architecture for natural language generation systems. *Natural Language Engineering*, 12:1–34.
- Margaret Mitchell. 2009. Class-based ordering of pronominal modifiers. In *ENLG '09: Proceedings of the 12th European Workshop on Natural Language Generation*, pages 50–57, Morristown, NJ, USA. Association for Computational Linguistics.
- Thomas Pechmann. 1989. Incremental speech production and referential overspecification. *Linguistics*, 27:89–110.
- Martin Pickering and Simon Garrod. 2004. Towards a mechanistic psychology of dialogue. *Behavioural and Brain Sciences*, 27:169–226.
- Ehud Reiter and Robert Dale. 2000. *Building Natural Language Generation Systems*. Cambridge University Press.
- James Shaw and Vasileios Hatzivassiloglou. 1999. Ordering among premodifiers. In *Proceedings of the 37th annual meeting of the Association for Computational Linguistics on Computational Linguistics*, pages 135–143.

Cognitively Plausible Models of Human Language Processing

Frank Keller

School of Informatics, University of Edinburgh
10 Crichton Street, Edinburgh EH8 9AB, UK
keller@inf.ed.ac.uk

Abstract

We pose the development of cognitively plausible models of human language processing as a challenge for computational linguistics. Existing models can only deal with isolated phenomena (e.g., garden paths) on small, specifically selected data sets. The challenge is to build models that integrate multiple aspects of human language processing at the syntactic, semantic, and discourse level. Like human language processing, these models should be incremental, predictive, broad coverage, and robust to noise. This challenge can only be met if standardized data sets and evaluation measures are developed.

1 Introduction

In many respects, human language processing is the ultimate goldstandard for computational linguistics. Humans understand and generate language with amazing speed and accuracy, they are able to deal with ambiguity and noise effortlessly and can adapt to new speakers, domains, and registers. Most surprisingly, they achieve this competency on the basis of limited training data (Hart and Risley, 1995), using learning algorithms that are largely unsupervised.

Given the impressive performance of humans as language processors, it seems natural to turn to psycholinguistics, the discipline that studies human language processing, as a source of information about the design of efficient language processing systems. Indeed, psycholinguists have uncovered an impressive array of relevant facts (reviewed in Section 2), but computational linguists are often not aware of this literature, and results about human language processing rarely inform the design, implementation, or evaluation of artificial language processing systems.

At the same time, research in psycholinguistics is often oblivious of work in computational

linguistics (CL). To test their theories, psycholinguists construct computational models of human language processing, but these models often fall short of the engineering standards that are generally accepted in the CL community (e.g., broad coverage, robustness, efficiency): typical psycholinguistic models only deal with isolated phenomena and fail to scale to realistic data sets. A particular issue is evaluation, which is typically anecdotal, performed on a small set of hand-crafted examples (see Sections 3).

In this paper, we propose a challenge that requires the combination of research efforts in computational linguistics and psycholinguistics: the development of cognitively plausible models of human language processing. This task can be decomposed into a modeling challenge (building models that instantiate known properties of human language processing) and a data and evaluation challenge (accounting for experimental findings and evaluating against standardized data sets), which we will discuss in turn.

2 Modeling Challenge

2.1 Key Properties

The first part of the challenge is to develop a model that instantiates key properties of human language processing, as established by psycholinguistic experimentation (see Table 1 for an overview and representative references).¹ A striking property of the human language processor is its *efficiency and robustness*. For the vast majority of sentences, it will effortlessly and rapidly deliver the correct analysis, even in the face of noise and ungrammaticalities. There is considerable experimental evi-

¹Here in the following, we will focus on sentence processing, which is often regarded as a central aspect of human language processing. A more comprehensive answer to our modeling challenge should also include phonological and morphological processing, semantic inference, discourse processing, and other non-syntactic aspects of language processing. Furthermore, established results regarding the interface between language processing and non-linguistic cognition (e.g., the sensorimotor system) should ultimately be accounted for in a fully comprehensive model.

Property	Evidence	Model			
		Rank	Surp	Pred	Stack
Efficiency and robustness	Ferreira et al. (2001); Sanford and Sturt (2002)	–	–	–	+
Broad coverage	Crocker and Brants (2000)	+	+	–	+
Incrementality and connectedness	Tanenhaus et al. (1995); Sturt and Lombardo (2005)	+	+	+	+
Prediction	Kamide et al. (2003); Staub and Clifton (2006)	–	±	+	–
Memory cost	Gibson (1998); Vasishth and Lewis (2006)	–	–	+	+

Table 1: Key properties of human language processing and their instantiation in various models of sentence processing (see Section 2 for details)

dence that shallow processing strategies are used to achieve this. The processor also achieves *broad coverage*: it can deal with a wide variety of syntactic constructions, and is not restricted by the domain, register, or modality of the input.

Human language processing is also word-by-word *incremental*. There is strong evidence that a new word is integrated as soon as it is available into the representation of the sentence thus far. Readers and listeners experience differential processing difficulty during this integration process, depending on the properties of the new word and its relationship to the preceding context. There is evidence that the processor instantiates a strict form of incrementality by building only fully connected trees. Furthermore, the processor is able to make *predictions* about upcoming material on the basis of sentence prefixes. For instance, listeners can predict an upcoming post-verbal element based on the semantics of the preceding verb. Or they can make syntactic predictions, e.g., if they encounter the word *either*, they predict an upcoming *or* and the type of complement that follows it.

Another key property of human language processing is the fact that it operates with limited memory, and that structures in memory are subject to decay and interference. In particular, the processor is known to incur a distance-based *memory cost*: combining the head of a phrase with its syntactic dependents is more difficult the more dependents have to be integrated and the further away they are. This integration process is also subject to interference from similar items that have to be held in memory at the same time.

2.2 Current Models

The challenge is to develop a computational model that captures the key properties of human language processing outlined in the previous section. A number of relevant models have been developed, mostly based on probabilistic parsing techniques, but none of them instantiates all the key properties discussed above (Table 1 gives an overview of

model properties).²

The earliest approaches were *ranking-based models* (Rank), which make psycholinguistic predictions based on the ranking of the syntactic analyses produced by a probabilistic parser. Jurafsky (1996) assumes that processing difficulty is triggered if the correct analysis falls below a certain probability threshold (i.e., is pruned by the parser). Similarly, Crocker and Brants (2000) assume that processing difficulty ensures if the highest-ranked analysis changes from one word to the next. Both approaches have been shown to successfully model garden path effects. Being based on probabilistic parsing techniques, ranking-based models generally achieve a broad coverage, but their efficiency and robustness has not been evaluated. Also, they are not designed to capture syntactic prediction or memory effects (other than search with a narrow beam in Brants and Crocker 2000).

The ranking-based approach has been generalized by *surprisal models* (Surp), which predict processing difficulty based on the change in the probability distribution over possible analyses from one word to the next (Hale, 2001; Levy, 2008; Demberg and Keller, 2008a; Ferrara Boston et al., 2008; Roark et al., 2009). These models have been successful in accounting for a range of experimental data, and they achieve broad coverage. They also instantiate a limited form of prediction, viz., they build up expectations about the next word in the input. On the other hand, the efficiency and robustness of these models has largely not been evaluated, and memory costs are not modeled (again except for restrictions in beam size).

The *prediction model* (Pred) explicitly predicts syntactic structure for upcoming words (Demberg and Keller, 2008b, 2009), thus accounting for experimental results on predictive language processing. It also implements a strict form of incre-

²We will not distinguish between model and linking theory, i.e., the set of assumptions that links model quantities to behavioral data (e.g., more probable structures are easier to process). It is conceivable, for instance, that a stack-based model is combined with a linking theory based on surprisal.

Factor	Evidence
Word senses	Roland and Jurafsky (2002)
Selectional restrictions	Garnsey et al. (1997); Pickering and Traxler (1998)
Thematic roles	McRae et al. (1998); Pickering et al. (2000)
Discourse reference	Altmann and Steedman (1988); Grodner and Gibson (2005)
Discourse coherence	Stewart et al. (2000); Kehler et al. (2008)

Table 2: Semantic factors in human language processing

mentality by building fully connected trees. Memory costs are modeled directly as a distance-based penalty that is incurred when a prediction has to be verified later in the sentence. However, the current implementation of the prediction model is neither robust and efficient nor offers broad coverage.

Recently, a *stack-based model* (Stack) has been proposed that imposes explicit, cognitively motivated memory constraints on the parser, in effect limiting the stack size available to the parser (Schuler et al., 2010). This delivers robustness, efficiency, and broad coverage, but does not model syntactic prediction. Unlike the other models discussed here, no psycholinguistic evaluation has been conducted on the stack-based model, so its cognitive plausibility is preliminary.

2.3 Beyond Parsing

There is strong evidence that human language processing is driven by an interaction of syntactic, semantic, and discourse processes (see Table 2 for an overview and references). Considerable experimental work has focused on the semantic properties of the verb of the sentence, and verb sense, selectional restrictions, and thematic roles have all been shown to interact with syntactic ambiguity resolution. Another large body of research has elucidated the interaction of discourse processing and syntactic processing. The most-well known effect is probably that of referential context: syntactic ambiguities can be resolved if a discourse context is provided that makes one of the syntactic alternatives more plausible. For instance, in a context that provides two possible antecedents for a noun phrase, the processor will prefer attaching a PP or a relative clause such that it disambiguates between the two antecedents; garden paths are reduced or disappear. Other results point to the importance of discourse coherence for sentence processing, an example being implicit causality.

The challenge facing researchers in computational and psycholinguistics therefore includes

the development of language processing models that combine syntactic processing with semantic and discourse processing. So far, this challenge is largely unmet: there are some examples of models that integrate semantic processes such as thematic role assignment into a parsing model (Narayanan and Jurafsky, 2002; Padó et al., 2009). However, other semantic factors are not accounted for by these models, and incorporating non-lexical aspects of semantics into models of sentence processing is a challenge for ongoing research. Recently, Dubey (2010) has proposed an approach that combines a probabilistic parser with a model of co-reference and discourse inference based on probabilistic logic. An alternative approach has been taken by Pynte et al. (2008) and Mitchell et al. (2010), who combine a vector-space model of semantics (Landauer and Dumais, 1997) with a syntactic parser and show that this results in predictions of processing difficulty that can be validated against an eye-tracking corpus.

2.4 Acquisition and Crosslinguistics

All models of human language processing discussed so far rely on supervised training data. This raises another aspect of the modeling challenge: the human language processor is the product of an acquisition process that is largely unsupervised and has access to only limited training data: children aged 12–36 months are exposed to between 10 and 35 million words of input (Hart and Risley, 1995). The challenge therefore is to develop a model of language acquisition that works with such small training sets, while also giving rise to a language processor that meets the key criteria in Table 1. The CL community is in a good position to rise to this challenge, given the significant progress in unsupervised parsing in recent years (starting from Klein and Manning 2002). However, none of the existing unsupervised models has been evaluated against psycholinguistic data sets, and they are not designed to meet even basic psycholinguistic criteria such as incrementality.

A related modeling challenge is the development of processing models for languages other than English. There is a growing body of experimental research investigating human language processing in other languages, but virtually all existing psycholinguistic models only work for English (the only exceptions we are aware of are Dubey et al.’s (2008) and Ferrara Boston et al.’s

(2008) parsing models for German). Again, the CL community has made significant progress in crosslinguistic parsing, especially using dependency grammar (Hajič, 2009), and psycholinguistic modeling could benefit from this in order to meet the challenge of developing crosslinguistically valid models of human language processing.

3 Data and Evaluation Challenge

3.1 Test Sets

The second key challenge that needs to be addressed in order to develop cognitively plausible models of human language processing concerns test data and model evaluation. Here, the state of the art in psycholinguistic modeling lags significantly behind standards in the CL community. Most of the models discussed in Section 2 have not been evaluated rigorously. The authors typically describe their performance on a small set of hand-picked examples; no attempts are made to test on a range of items from the experimental literature and determine model fit directly against behavioral measures (e.g., reading times). This makes it very hard to obtain a realistic estimate of how well the models achieve their aim of capturing human language processing behavior.

We therefore suggest the development of standard test sets for psycholinguistic modeling, similar to what is commonplace for tasks in computational linguistics: parsers are evaluated against the Penn Treebank, word sense disambiguation systems against the SemEval data sets, co-reference systems against the Tipster or ACE corpora, etc. Two types of test data are required for psycholinguistic modeling. The first type of test data consists of a collection of representative experimental results. This collection should contain the actual experimental materials (sentences or discourse fragments) used in the experiments, together with the behavioral measurements obtained (reading times, eye-movement records, rating judgments, etc.). The experiments included in this test set would be chosen to cover a wide range of experimental phenomena, e.g., garden paths, syntactic complexity, memory effects, semantic and discourse factors. Such a test set will enable the standardized evaluation of psycholinguistic models by comparing the model predictions (rankings, surprisal values, memory costs, etc.) against behavioral measures on a large set of items. This way both the coverage of a model (how many phenom-

ena can it account for) and its accuracy (how well does it fit the behavioral data) can be assessed.

Experimental test sets should be complemented by test sets based on corpus data. In order to assess the efficiency, robustness, and broad coverage of a model, a corpus of unrestricted, naturally occurring text is required. The use of contextualized language data makes it possible to assess not only syntactic models, but also models that capture discourse effects. These corpora need to be annotated with behavioral measures, e.g., eye-tracking or reading time data. Some relevant corpora have already been constructed, see the overview in Table 3, and various authors have used them for model evaluation (Demberg and Keller, 2008a; Pynte et al., 2008; Frank, 2009; Ferrara Boston et al., 2008; Patil et al., 2009; Roark et al., 2009; Mitchell et al., 2010).

However, the usefulness of the psycholinguistic corpora in Table 3 is restricted by the absence of gold-standard linguistic annotation (though the French part of the Dundee corpus, which is syntactically annotated). This makes it difficult to test the accuracy of the linguistic structures computed by a model, and restricts evaluation to behavioral predictions. The challenge is therefore to collect a standardized test set of naturally occurring text or speech enriched not only with behavioral variables, but also with syntactic and semantic annotation. Such a data set could for example be constructed by eye-tracking section 23 of the Penn Treebank (which is also part of Propbank, and thus has both syntactic and thematic role annotation).

In computational linguistics, the development of new data sets is often stimulated by competitions in which systems are compared on a standardized task, using a data set specifically designed for the competition. Examples include the CoNLL shared task, SemEval, or TREC in computational syntax, semantics, and discourse, respectively. A similar competition could be developed for computational psycholinguistics – maybe along the lines of the model comparison challenges that held at the International Conference on Cognitive Modeling. These challenges provide standardized task descriptions and data sets; participants can enter their cognitive models, which were then compared using a pre-defined evaluation metric.³

³The ICCM 2009 challenge was the Dynamic Stock and Flows Task, for more information see <http://www.hss.cmu.edu/departments/sds/ddmlab/modeldsf/>.

Corpus	Language	Words	Participants	Method	Reference
Dundee Corpus	English, French	50,000	10	Eye-tracking	Kennedy and Pynte (2005)
Potsdam Corpus	German	1,138	222	Eye-tracking	Kliegl et al. (2006)
MIT Corpus	English	3,534	23	Self-paced reading	Bachrach (2008)

Table 3: Test corpora that have been used for psycholinguistic modeling of sentence processing; note that the Potsdam Corpus consists of isolated sentences, rather than of continuous text

3.2 Behavioral and Neural Data

As outlined in the previous section, a number of authors have evaluated psycholinguistic models against eye-tracking or reading time corpora. Part of the data and evaluation challenge is to extend this evaluation to neural data as provided by event-related potential (ERP) or brain imaging studies (e.g., using functional magnetic resonance imaging, fMRI). Neural data sets are considerably more complex than behavioral ones, and modeling them is an important new task that the community is only beginning to address. Some recent work has evaluated models of word semantics against ERP (Murphy et al., 2009) or fMRI data (Mitchell et al., 2008).⁴ This is a very promising direction, and the challenge is to extend this approach to the sentence and discourse level (see Bachrach 2008). Again, it will again be necessary to develop standardized test sets of both experimental data and corpus data.

3.3 Evaluation Measures

We also anticipate that the availability of new test data sets will facilitate the development of new evaluation measures that specifically test the validity of psycholinguistic models. Established CL evaluation measures such as Parseval are of limited use, as they can only test the linguistic, but not the behavioral or neural predictions of a model.

So far, many authors have relied on qualitative evaluation: if a model predicts a difference in (for instance) reading time between two types of sentences where such a difference was also found experimentally, then that counts as a successful test. In most cases, no quantitative evaluation is performed, as this would require modeling the reading times for individual item and individual participants. Suitable procedures for performing such tests do not currently exist; linear mixed effects models (Baayen et al., 2008) provide a way of dealing with item and participant variation, but crucially do not enable direct comparisons between models in terms of goodness of fit.

⁴These data sets were released as part of the NAACL-2010 Workshop on Computational Neurolinguistics.

Further issues arise from the fact that we often want to compare model fit for multiple experiments (ideally without reparametrizing the models), and that various mutually dependent measures are used for evaluation, e.g., processing effort at the sentence, word, and character level. An important open challenge is there to develop evaluation measures and associated statistical procedures that can deal with these problems.

4 Conclusions

In this paper, we discussed the modeling and data/evaluation challenges involved in developing cognitively plausible models of human language processing. Developing computational models is of scientific importance in so far as models are implemented theories: models of language processing allow us to test scientific hypothesis about the cognitive processes that underpin language processing. This type of precise, formalized hypothesis testing is only possible if standardized data sets and uniform evaluation procedures are available, as outlined in the present paper. Ultimately, this approach enables qualitative and quantitative comparisons between theories, and thus enhances our understanding of a key aspect of human cognition, language processing.

There is also an applied side to the proposed challenge. Once computational models of human language processing are available, they can be used to predict the difficulty that humans experience when processing text or speech. This is useful for a number applications: for instance, natural language generation would benefit from being able to assess whether machine-generated text or speech is easy to process. For text simplification (e.g., for children or impaired readers), such a model is even more essential. It could also be used to assess the readability of text, which is of interest in educational applications (e.g., essay scoring). In machine translation, evaluating the fluency of system output is crucial, and a model that predicts processing difficulty could be used for this, or to guide the choice between alternative translations, and maybe even to inform human post-editing.

References

- Altmann, Gerry T. M. and Mark J. Steedman. 1988. Interaction with context during human sentence processing. *Cognition* 30(3):191–238.
- Baayen, R. H., D. J. Davidson, and D. M. Bates. 2008. Mixed-effects modeling with crossed random effects for subjects and items. *Journal of Memory and Language* to appear.
- Bachrach, Asaf. 2008. *Imaging Neural Correlates of Syntactic Complexity in a Naturalistic Context*. Ph.D. thesis, Massachusetts Institute of Technology, Cambridge, MA.
- Brants, Thorsten and Matthew W. Crocker. 2000. Probabilistic parsing and psychological plausibility. In *Proceedings of the 18th International Conference on Computational Linguistics*. Saarbrücken/Luxembourg/Nancy, pages 111–117.
- Crocker, Matthew W. and Thorsten Brants. 2000. Wide-coverage probabilistic sentence processing. *Journal of Psycholinguistic Research* 29(6):647–669.
- Demberg, Vera and Frank Keller. 2008a. Data from eye-tracking corpora as evidence for theories of syntactic processing complexity. *Cognition* 101(2):193–210.
- Demberg, Vera and Frank Keller. 2008b. A psycholinguistically motivated version of TAG. In *Proceedings of the 9th International Workshop on Tree Adjoining Grammars and Related Formalisms*. Tübingen, pages 25–32.
- Demberg, Vera and Frank Keller. 2009. A computational model of prediction in human parsing: Unifying locality and surprisal effects. In Niels Taatgen and Hedderik van Rijn, editors, *Proceedings of the 31st Annual Conference of the Cognitive Science Society*. Cognitive Science Society, Amsterdam, pages 1888–1893.
- Dubey, Amit. 2010. The influence of discourse on syntax: A psycholinguistic model of sentence processing. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*. Uppsala.
- Dubey, Amit, Frank Keller, and Patrick Sturt. 2008. A probabilistic corpus-based model of syntactic parallelism. *Cognition* 109(3):326–344.
- Ferrara Boston, Marisa, John Hale, Reinhold Kliegl, Umesh Patil, and Shravan Vasishth. 2008. Parsing costs as predictors of reading difficulty: An evaluation using the Potsdam Sentence Corpus. *Journal of Eye Movement Research* 2(1):1–12.
- Ferreira, Fernanda, Kiel Christianson, and Andrew Hollingworth. 2001. Misinterpretations of garden-path sentences: Implications for models of sentence processing and reanalysis. *Journal of Psycholinguistic Research* 30(1):3–20.
- Frank, Stefan L. 2009. Surprisal-based comparison between a symbolic and a connectionist model of sentence processing. In Niels Taatgen and Hedderik van Rijn, editors, *Proceedings of the 31st Annual Conference of the Cognitive Science Society*. Cognitive Science Society, Amsterdam, pages 1139–1144.
- Garnsey, Susan M., Neal J. Pearlmutter, Elisabeth M. Myers, and Melanie A. Lotocky. 1997. The contributions of verb bias and plausibility to the comprehension of temporarily ambiguous sentences. *Journal of Memory and Language* 37(1):58–93.
- Gibson, Edward. 1998. Linguistic complexity: locality of syntactic dependencies. *Cognition* 68:1–76.
- Grodner, Dan and Edward Gibson. 2005. Consequences of the serial nature of linguistic input. *Cognitive Science* 29:261–291.
- Hajič, Jan, editor. 2009. *Proceedings of the 13th Conference on Computational Natural Language Learning: Shared Task*. Association for Computational Linguistics, Boulder, CO.
- Hale, John. 2001. A probabilistic Earley parser as a psycholinguistic model. In *Proceedings of the 2nd Conference of the North American Chapter of the Association for Computational Linguistics*. Association for Computational Linguistics, Pittsburgh, PA, volume 2, pages 159–166.
- Hart, Betty and Todd R. Risley. 1995. *Meaningful Differences in the Everyday Experience of Young American Children*. Paul H. Brookes, Baltimore, MD.
- Jurafsky, Daniel. 1996. A probabilistic model of lexical and syntactic access and disambiguation. *Cognitive Science* 20(2):137–194.
- Kamide, Yuki, Gerry T. M. Altmann, and Sarah L. Haywood. 2003. The time-course of prediction in incremental sentence processing: Evidence

- from anticipatory eye movements. *Journal of Memory and Language* 49:133–156.
- Kehler, Andrew, Laura Kertz, Hannah Rohde, and Jeffrey L. Elman. 2008. Coherence and coreference revisited. *Journal of Semantics* 25(1):1–44.
- Kennedy, Alan and Joel Pynte. 2005. Parafoveal-on-foveal effects in normal reading. *Vision Research* 45:153–168.
- Klein, Dan and Christopher Manning. 2002. A generative constituent-context model for improved grammar induction. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*. Philadelphia, pages 128–135.
- Kliegl, Reinhold, Antje Nuthmann, and Ralf Engbert. 2006. Tracking the mind during reading: The influence of past, present, and future words on fixation durations. *Journal of Experimental Psychology: General* 135(1):12–35.
- Landauer, Thomas K. and Susan T. Dumais. 1997. A solution to Plato’s problem: The latent semantic analysis theory of acquisition, induction and representation of knowledge. *Psychological Review* 104(2):211–240.
- Levy, Roger. 2008. Expectation-based syntactic comprehension. *Cognition* 106(3):1126–1177.
- McRae, Ken, Michael J. Spivey-Knowlton, and Michael K. Tanenhaus. 1998. Modeling the influence of thematic fit (and other constraints) in on-line sentence comprehension. *Journal of Memory and Language* 38(3):283–312.
- Mitchell, Jeff, Mirella Lapata, Vera Demberg, and Frank Keller. 2010. Syntactic and semantic factors in processing difficulty: An integrated measure. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*. Uppsala.
- Mitchell, Tom M., Svetlana V. Shinkareva, Andrew Carlson, Kai-Min Chang, Vicente L. Malave, Robert A. Mason, and Marcel Adam Just. 2008. Predicting human brain activity associated with the meanings of nouns. *Science* 320(5880):1191–1195.
- Murphy, Brian, Marco Baroni, and Massimo Poesio. 2009. EEG responds to conceptual stimuli and corpus semantics. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. Singapore, pages 619–627.
- Narayanan, Sridhar and Daniel Jurafsky. 2002. A Bayesian model predicts human parse preference and reading time in sentence processing. In Thomas G. Dietterich, Sue Becker, and Zoubin Ghahramani, editors, *Advances in Neural Information Processing Systems 14*. MIT Press, Cambridge, MA, pages 59–65.
- Padó, Ulrike, Matthew W. Crocker, and Frank Keller. 2009. A probabilistic model of semantic plausibility in sentence processing. *Cognitive Science* 33(5):794–838.
- Patil, Umesh, Shravan Vasishth, and Reinhold Kliegl. 2009. Compound effect of probabilistic disambiguation and memory retrievals on sentence processing: Evidence from an eye-tracking corpus. In A. Howes, D. Peebles, and R. Cooper, editors, *Proceedings of 9th International Conference on Cognitive Modeling*. Manchester.
- Pickering, Martin J. and Martin J. Traxler. 1998. Plausibility and recovery from garden paths: An eye-tracking study. *Journal of Experimental Psychology: Learning Memory and Cognition* 24(4):940–961.
- Pickering, Martin J., Matthew J. Traxler, and Matthew W. Crocker. 2000. Ambiguity resolution in sentence processing: Evidence against frequency-based accounts. *Journal of Memory and Language* 43(3):447–475.
- Pynte, Joel, Boris New, and Alan Kennedy. 2008. On-line contextual influences during reading normal text: A multiple-regression analysis. *Vision Research* 48(21):2172–2183.
- Roark, Brian, Asaf Bachrach, Carlos Cardenas, and Christophe Pallier. 2009. Deriving lexical and syntactic expectation-based measures for psycholinguistic modeling via incremental top-down parsing. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. Singapore, pages 324–333.
- Roland, Douglas and Daniel Jurafsky. 2002. Verb sense and verb subcategorization probabilities. In Paola Merlo and Suzanne Stevenson, editors, *The Lexical Basis of Sentence Processing: Formal, Computational, and Experimental Issues*, John Bejamins, Amsterdam, pages 325–346.
- Sanford, Anthony J. and Patrick Sturt. 2002.

- Depth of processing in language comprehension: Not noticing the evidence. *Trends in Cognitive Sciences* 6:382–386.
- Schuler, William, Samir AbdelRahman, Tim Miller, and Lane Schwartz. 2010. Broad-coverage parsing using human-like memory constraints. *Computational Linguistics* 26(1):1–30.
- Staub, Adrian and Charles Clifton. 2006. Syntactic prediction in language comprehension: Evidence from either . . . or. *Journal of Experimental Psychology: Learning, Memory, and Cognition* 32:425–436.
- Stewart, Andrew J., Martin J. Pickering, and Anthony J. Sanford. 2000. The time course of the influence of implicit causality information: Focusing versus integration accounts. *Journal of Memory and Language* 42(3):423–443.
- Sturt, Patrick and Vincenzo Lombardo. 2005. Processing coordinated structures: Incrementality and connectedness. *Cognitive Science* 29(2):291–305.
- Tanenhaus, Michael K., Michael J. Spivey-Knowlton, Kathleen M. Eberhard, and Julie C. Sedivy. 1995. Integration of visual and linguistic information in spoken language comprehension. *Science* 268:1632–1634.
- Vasishth, Shravan and Richard L. Lewis. 2006. Argument-head distance and processing complexity: Explaining both locality and antilocality effects. *Language* 82(4):767–794.

The Manually Annotated Sub-Corpus: A Community Resource For and By the People

Nancy Ide

Department of Computer Science
Vassar College
Poughkeepsie, NY, USA
ide@cs.vassar.edu

Christiane Fellbaum

Princeton University
Princeton, New Jersey USA
fellbaum@princeton.edu

Collin Baker

International Computer Science Institute
Berkeley, California USA
collinb@icsi.berkeley.edu

Rebecca Passonneau

Columbia University
New York, New York USA
becky@cs.columbia.edu

Abstract

The Manually Annotated Sub-Corpus (MASC) project provides data and annotations to serve as the base for a community-wide annotation effort of a subset of the American National Corpus. The MASC infrastructure enables the incorporation of contributed annotations into a single, usable format that can then be analyzed as it is or ported to any of a variety of other formats. MASC includes data from a much wider variety of genres than existing multiply-annotated corpora of English, and the project is committed to a fully open model of distribution, without restriction, for all data and annotations produced or contributed. As such, MASC is the first large-scale, open, community-based effort to create much needed language resources for NLP. This paper describes the MASC project, its corpus and annotations, and serves as a call for contributions of data and annotations from the language processing community.

1 Introduction

The need for corpora annotated for multiple phenomena across a variety of linguistic layers is keenly recognized in the computational linguistics community. Several multiply-annotated corpora exist, especially for Western European languages and for spoken data, but, interestingly, broad-based English language corpora with robust annotation for diverse linguistic phenomena are relatively rare. The most widely-used corpus of English, the British National Corpus, contains only part-of-speech annotation; and although it contains a wider range of annotation types, the fif-

teen million word Open American National Corpus annotations are largely unvalidated. The most well-known multiply-annotated and validated corpus of English is the one million word *Wall Street Journal* corpus known as the Penn Treebank (Marcus et al., 1993), which over the years has been fully or partially annotated for several phenomena over and above the original part-of-speech tagging and phrase structure annotation. The usability of these annotations is limited, however, by the fact that many of them were produced by independent projects using their own tools and formats, making it difficult to combine them in order to study their inter-relations. More recently, the OntoNotes project (Pradhan et al., 2007) released a one million word English corpus of newswire, broadcast news, and broadcast conversation that is annotated for Penn Treebank syntax, PropBank predicate argument structures, coreference, and named entities. OntoNotes comes closest to providing a corpus with multiple layers of annotation that can be analyzed as a unit via its representation of the annotations in a “normal form”. However, like the *Wall Street Journal* corpus, OntoNotes is limited in the range of genres it includes. It is also limited to only those annotations that may be produced by members of the OntoNotes project. In addition, use of the data and annotations with software other than the OntoNotes database API is not necessarily straightforward.

The sparseness of reliable multiply-annotated corpora can be attributed to several factors. The greatest obstacle is the high cost of manual production and validation of linguistic annotations. Furthermore, the production and annotation of corpora, even when they involve significant scientific research, often do not, *per se*, lead to publishable research results. It is therefore understand-

able that many research groups are unwilling to get involved in such a massive undertaking for relatively little reward.

The Manually Annotated Sub-Corpus (MASC) (Ide et al., 2008) project has been established to address many of these obstacles to the creation of large-scale, robust, multiply-annotated corpora. The project is providing appropriate data and annotations to serve as the base for a community-wide annotation effort, together with an infrastructure that enables the representation of internally-produced and contributed annotations in a single, usable format that can then be analyzed as it is or ported to any of a variety of other formats, thus enabling its immediate use with many common annotation platforms as well as off-the-shelf concordance and analysis software. The MASC project’s aim is to offset some of the high costs of producing high quality linguistic annotations via a distribution of effort, and to solve some of the usability problems for annotations produced at different sites by harmonizing their representation formats.

The MASC project provides a resource that is significantly different from OntoNotes and similar corpora. It provides data from a much wider variety of genres than existing multiply-annotated corpora of English, and all of the data in the corpus are drawn from current American English so as to be most useful for NLP applications. Perhaps most importantly, the MASC project is committed to a fully open model of distribution, without restriction, for all data and annotations. It is also committed to incorporating diverse annotations contributed by the community, regardless of format, into the corpus. As such, MASC is the first large-scale, open, community-based effort to create a much-needed language resource for NLP. This paper describes the MASC project, its corpus and annotations, and serves as a call for contributions of data and annotations from the language processing community.

2 MASC: The Corpus

MASC is a balanced subset of 500K words of written texts and transcribed speech drawn primarily from the Open American National Corpus (OANC)¹. The OANC is a 15 million word (and growing) corpus of American English produced since 1990, all of which is in the public domain

¹<http://www.anc.org>

Genre	No. texts	Total words
Email	2	468
Essay	4	17516
Fiction	4	20413
Gov’t documents	1	6064
Journal	10	25635
Letters	31	10518
Newspaper/newswire	41	17951
Non-fiction	4	17118
Spoken	11	25783
Debate transcript	2	32325
Court transcript	1	20817
Technical	3	15417
Travel guides	4	12463
Total	118	222488

Table 1: MASC Composition (first 220K)

or otherwise free of usage and redistribution restrictions.

Where licensing permits, data for inclusion in MASC is drawn from sources that have already been heavily annotated by others. So far, the first 80K increment of MASC data includes a 40K subset consisting of OANC data that has been previously annotated for PropBank predicate argument structures, Pittsburgh Opinion annotation (opinions, evaluations, sentiments, etc.), TimeML time and events², and several other linguistic phenomena. It also includes a handful of small texts from the so-called Language Understanding (LU) Corpus³ that has been annotated by multiple groups for a wide variety of phenomena, including events and committed belief. All of the first 80K increment is annotated for Penn Treebank syntax. The second 120K increment includes 5.5K words of *Wall Street Journal* texts that have been annotated by several projects, including Penn Treebank, PropBank, Penn Discourse Treebank, TimeML, and the Pittsburgh Opinion project. The composition of the 220K portion of the corpus annotated so far is shown in Table 1. The remaining 280K of the corpus fills out the genres that are under-represented in the first portion and includes a few additional genres such as blogs and tweets.

3 MASC Annotations

Annotations for a variety of linguistic phenomena, either manually produced or corrected from output of automatic annotation systems, are being added

²The TimeML annotations of the data are not yet completed.

³MASC contains about 2K words of the 10K LU corpus, eliminating non-English and translated LU texts as well as texts that are not free of usage and redistribution restrictions.

Annotation type	Method	No. texts	No. words
Token	Validated	118	222472
Sentence	Validated	118	222472
POS/lemma	Validated	118	222472
Noun chunks	Validated	118	222472
Verb chunks	Validated	118	222472
Named entities	Validated	118	222472
FrameNet frames	Manual	21	17829
HSPG	Validated	40*	30106
Discourse	Manual	40*	30106
Penn Treebank	Validated	97	87383
PropBank	Validated	92	50165
Opinion	Manual	97	47583
TimeBank	Validated	34	5434
Committed belief	Manual	13	4614
Event	Manual	13	4614
Coreference	Manual	2	1877

Table 2: Current MASC Annotations (* projected)

to MASC data in increments of roughly 100K words. To date, validated or manually produced annotations for 222K words have been made available.

The MASC project is itself producing annotations for portions of the corpus for WordNet senses and FrameNet frames and frame elements. To derive maximal benefit from the semantic information provided by these resources, the entire corpus is also annotated and manually validated for shallow parses (noun and verb chunks) and named entities (person, location, organization, date and time). Several additional types of annotation have either been contracted by the MASC project or contributed from other sources. The 220K words of MASC I and II include seventeen different types of linguistic annotation⁴, shown in Table 2.

All MASC annotations, whether contributed or produced in-house, are transduced to the Graph Annotation Framework (GrAF) (Ide and Suderman, 2007) defined by ISO TC37 SC4’s Linguistic Annotation Framework (LAF) (Ide and Romary, 2004). GrAF is an XML serialization of the LAF abstract model of annotations, which consists of a directed graph decorated with feature structures providing the annotation content. GrAF’s primary role is to serve as a “pivot” format for transducing among annotations represented in different formats. However, because the underlying data structure is a graph, the GrAF representation itself can serve as the basis for analysis via application of

⁴This includes WordNet sense annotations, which are not listed in Table 2 because they are not applied to full texts; see Section 3.1 for a description of the WordNet sense annotations in MASC.

graph-analytic algorithms such as common subtree detection.

The layering of annotations over MASC texts dictates the use of a stand-off annotation representation format, in which each annotation is contained in a separate document linked to the primary data. Each text in the corpus is provided in UTF-8 character encoding in a separate file, which includes no annotation or markup of any kind. Each file is associated with a set of GrAF standoff files, one for each annotation type, containing the annotations for that text. In addition to the annotation types listed in Table 2, a document containing annotation for logical structure (titles, headings, sections, etc. down to the level of paragraph) is included. Each text is also associated with (1) a header document that provides appropriate metadata together with machine-processable information about associated annotations and interrelations among the annotation layers; and (2) a segmentation of the primary data into minimal regions, which enables the definition of different tokenizations over the text. Contributed annotations are also included in their original format, where available.

3.1 WordNet Sense Annotations

A focus of the MASC project is to provide corpus evidence to support an effort to harmonize sense distinctions in WordNet and FrameNet (Baker and Fellbaum, 2009), (Fellbaum and Baker, to appear). The WordNet and FrameNet teams have selected for this purpose 100 common polysemous words whose senses they will study in detail, and the MASC team is annotating occurrences of these words in the MASC. As a first step, fifty occurrences of each word are annotated using the WordNet 3.0 inventory and analyzed for problems in sense assignment, after which the WordNet team may make modifications to the inventory if needed. The revised inventory (which will be released as part of WordNet 3.1) is then used to annotate 1000 occurrences. Because of its small size, MASC typically contains less than 1000 occurrences of a given word; the remaining occurrences are therefore drawn from the 15 million words of the OANC. Furthermore, the FrameNet team is also annotating one hundred of the 1000 sentences for each word with FrameNet frames and frame elements, providing direct comparisons of WordNet and FrameNet sense assignments in

attested sentences.⁵

For convenience, the annotated sentences are provided as a stand-alone corpus, with the WordNet and FrameNet annotations represented in standoff files. Each sentence in this corpus is linked to its occurrence in the original text, so that the context and other annotations associated with the sentence may be retrieved.

3.2 Validation

Automatically-produced annotations for sentence, token, part of speech, shallow parses (noun and verb chunks), and named entities (person, location, organization, date and time) are hand-validated by a team of students. Each annotation set is first corrected by one student, after which it is checked (and corrected where necessary) by a second student, and finally checked by both automatic extraction of the annotated data and a third pass over the annotations by a graduate student or senior researcher. We have performed inter-annotator agreement studies for shallow parses in order to establish the number of passes required to achieve near-100% accuracy.

Annotations produced by other projects and the FrameNet and Penn Treebank annotations produced specifically for MASC are semi-automatically and/or manually produced by those projects and subjected to their internal quality controls. No additional validation is performed by the ANC project.

The WordNet sense annotations are being used as a base for an extensive inter-annotator agreement study, which is described in detail in (Pasonneau et al., 2009), (Pasonneau et al., 2010). All inter-annotator agreement data and statistics are published along with the sense tags. The release also includes documentation on the words annotated in each round, the sense labels for each word, the sentences for each word, and the annotator or annotators for each sense assignment to each word in context. For the multiply annotated data in rounds 2-4, we include raw tables for each word in the form expected by Ron Artstein's `calculate_alpha.pl` perl script⁶, so that the agreement numbers can be regenerated.

⁵Note that several MASC texts have been fully annotated for FrameNet frames and frame elements, in addition to the WordNet-tagged sentences.

⁶<http://ron.artstein.org/resources/calculate-alpha.perl>

4 MASC Availability and Distribution

Like the OANC, MASC is distributed without license or other restrictions from the American National Corpus website⁷. It is also available from the Linguistic Data Consortium (LDC)⁸ for a nominal processing fee.

In addition to enabling download of the entire MASC, we provide a web application that allows users to select some or all parts of the corpus and choose among the available annotations via a web interface (Ide et al., 2010). Once generated, the corpus and annotation bundle is made available to the user for download. Thus, the MASC user need never deal directly with or see the underlying representation of the stand-off annotations, but gains all the advantages that representation offers. The following output formats are currently available:

1. in-line XML (XCES⁹), suitable for use with the BNCs XAIRA search and access interface and other XML-aware software;
2. token / part of speech, a common input format for general-purpose concordance software such as MonoConc¹⁰, as well as the Natural Language Toolkit (NLTK) (Bird et al., 2009);
3. CONLL IOB format, used in the Conference on Natural Language Learning shared tasks.¹¹

5 Tools

The ANC project provides an API for GrAF annotations that can be used to access and manipulate GrAF annotations directly from Java programs and render GrAF annotations in a format suitable for input to the open source GraphViz¹² graph visualization application.¹³ Beyond this, the ANC project does not provide specific tools for use of the corpus, but rather provides the data in formats suitable for use with a variety of available applications, as described in section 4, together with means to import GrAF annotations into major annotation software platforms. In particular, the ANC project provides plugins for the General

⁷<http://www.anc.org>

⁸<http://www ldc.upenn.edu>

⁹XML Corpus Encoding Standard, <http://www.xces.org>

¹⁰<http://www.athel.com/mono.html>

¹¹<http://ifarm.nl/signll/conll>

¹²<http://www.graphviz.org/>

¹³<http://www.anc.org/graf-api>

Architecture for Text Engineering (GATE) (Cunningham et al., 2002) to input and/or output annotations in GrAF format; a “CAS Consumer” to enable using GrAF annotations in the Unstructured Information Management Architecture (UIMA) (Ferrucci and Lally, 2004); and a corpus reader for importing MASC data and annotations into NLTK¹⁴.

Because the GrAF format is isomorphic to input to many graph-analytic tools, existing graph-analytic software can also be exploited to search and manipulate MASC annotations. Trivial merging of GrAF-based annotations involves simply combining the graphs for each annotation, after which graph minimization algorithms¹⁵ can be applied to collapse nodes with edges to common subgraphs to identify commonly annotated components. Graph-traversal and graph-coloring algorithms can also be applied in order to identify and generate statistics that could reveal interactions among linguistic phenomena that may have previously been difficult to observe. Other graph-analytic algorithms — including common sub-graph analysis, shortest paths, minimum spanning trees, connectedness, identification of articulation vertices, topological sort, graph partitioning, etc. — may also prove to be useful for mining information from a graph of annotations at multiple linguistic levels.

6 Community Contributions

The ANC project solicits contributions of annotations of any kind, applied to any part or all of the MASC data. Annotations may be contributed in any format, either inline or standoff. All contributed annotations are ported to GrAF standoff format so that they may be used with other MASC annotations and rendered in the various formats the ANC tools generate. To accomplish this, the ANC project has developed a suite of internal tools and methods for automatically transducing other annotation formats to GrAF and for rapid adaptation of previously unseen formats.

Contributions may be emailed to anc@cs.vassar.edu or uploaded via the ANC website¹⁶. The validity of annotations and supplemental documentation (if appropriate) are the responsibility of the contributor. MASC

¹⁴Available in September, 2010.

¹⁵Efficient algorithms for graph merging exist; see, e.g., (Habib et al., 2000).

¹⁶<http://www.anc.org/contributions.html>

users may contribute evaluations and error reports for the various annotations on the ANC/MASC wiki¹⁷.

Contributions of unvalidated annotations for MASC and OANC data are also welcomed and are distributed separately. Contributions of unencumbered texts in any genre, including stories, papers, student essays, poetry, blogs, and email, are also solicited via the ANC web site and the ANC Facebook page¹⁸, and may be uploaded at the contribution page cited above.

7 Conclusion

MASC is already the most richly annotated corpus of English available for widespread use. Because the MASC is an open resource that the community can continually enhance with additional annotations and modifications, the project serves as a model for community-wide resource development in the future. Past experience with corpora such as the *Wall Street Journal* shows that the community is eager to annotate available language data, and we anticipate even greater interest in MASC, which includes language data covering a range of genres that no existing resource provides. Therefore, we expect that as MASC evolves, more and more annotations will be contributed, thus creating a massive, inter-linked linguistic infrastructure for the study and processing of current American English in its many genres and varieties. In addition, by virtue of its WordNet and FrameNet annotations, MASC will be linked to parallel WordNets and FrameNets in languages other than English, thus creating a global resource for multi-lingual technologies, including machine translation.

Acknowledgments

The MASC project is supported by National Science Foundation grant CRI-0708952. The WordNet-FrameNet alignment work is supported by NSF grant IIS 0705155.

References

Collin F. Baker and Christiane Fellbaum. 2009. WordNet and FrameNet as complementary resources for annotation. In *Proceedings of the Third Linguistic*

¹⁷<http://www.anc.org/masc-wiki>

¹⁸<http://www.facebook.com/pages/American-National-Corpus/42474226671>

- Annotation Workshop*, pages 125–129, Suntec, Singapore, August. Association for Computational Linguistics.
- Steven Bird, Ewan Klein, and Edward Loper. 2009. *Natural Language Processing with Python*. O’Reilly Media, 1st edition.
- Hamish Cunningham, Diana Maynard, Kalina Bontcheva, and Valentin Tablan. 2002. GATE: A framework and graphical development environment for robust nlp tools and applications. In *Proceedings of ACL’02*.
- Christiane Fellbaum and Collin Baker. to appear. Aligning verbs in WordNet and FrameNet. *Linguistics*.
- David Ferrucci and Adam Lally. 2004. UIMA: An architectural approach to unstructured information processing in the corporate research environment. *Natural Language Engineering*, 10(3-4):327–348.
- Michel Habib, Christophe Paul, and Laurent Viennot. 2000. Partition refinement techniques: an interesting algorithmic tool kit. *International Journal of Foundations of Computer Science*, 175.
- Nancy Ide and Laurent Romary. 2004. International standard for a linguistic annotation framework. *Natural Language Engineering*, 10(3-4):211–225.
- Nancy Ide and Keith Suderman. 2007. GrAF: A graph-based format for linguistic annotations. In *Proceedings of the Linguistic Annotation Workshop*, pages 1–8, Prague, Czech Republic, June. Association for Computational Linguistics.
- Nancy Ide, Collin Baker, Christiane Fellbaum, Charles Fillmore, and Rebecca Passonneau. 2008. MASC: The Manually Annotated Sub-Corpus of American English. In *Proceedings of the Sixth International Conference on Language Resources and Evaluation (LREC)*, Marrakech, Morocco.
- Nancy Ide, Keith Suderman, and Brian Simms. 2010. ANC2Go: A web application for customized corpus creation. In *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC)*, Valletta, Malta, May. European Language Resources Association.
- Mitchell P. Marcus, Mary Ann Marcinkiewicz, and Beatrice Santorini. 1993. Building a large annotated corpus of English: the Penn Treebank. *Computational Linguistics*, 19(2):313–330.
- Rebecca J. Passonneau, Ansaf Salieb-Aouissi, and Nancy Ide. 2009. Making sense of word sense variation. In *SEW ’09: Proceedings of the Workshop on Semantic Evaluations: Recent Achievements and Future Directions*, pages 2–9, Morristown, NJ, USA. Association for Computational Linguistics.
- Rebecca Passonneau, Ansaf Salieb-Aouissi, Vikas Bhardwaj, and Nancy Ide. 2010. Word sense annotation of polysemous words by multiple annotators. In *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC)*, Valletta, Malta.
- Sameer S. Pradhan, Eduard Hovy, Mitch Marcus, Martha Palmer, Lance Ramshaw, and Ralph Weischedel. 2007. OntoNotes: A unified relational semantic representation. In *ICSC ’07: Proceedings of the International Conference on Semantic Computing*, pages 517–526, Washington, DC, USA. IEEE Computer Society.

Correcting Errors in a Treebank Based on Synchronous Tree Substitution Grammar

Yoshihide Kato¹ and Shigeki Matsubara²

¹Information Technology Center, Nagoya University

²Graduate School of Information Science, Nagoya University

Furo-cho, Chikusa-ku, Nagoya, 464-8601 Japan

yoshihide@el.itc.nagoya-u.ac.jp

Abstract

This paper proposes a method of correcting annotation errors in a treebank. By using a synchronous grammar, the method transforms parse trees containing annotation errors into the ones whose errors are corrected. The synchronous grammar is automatically induced from the treebank. We report an experimental result of applying our method to the Penn Treebank. The result demonstrates that our method corrects syntactic annotation errors with high precision.

1 Introduction

Annotated corpora play an important role in the fields such as theoretical linguistic researches or the development of NLP systems. However, they often contain annotation errors which are caused by a manual or semi-manual mark-up process. These errors are problematic for corpus-based researches.

To solve this problem, several error detection and correction methods have been proposed so far (Eskin, 2000; Nakagawa and Matsumoto, 2002; Dickinson and Meurers, 2003a; Dickinson and Meurers, 2003b; Ule and Simov, 2004; Murata et al., 2005; Dickinson and Meurers, 2005; Boyd et al., 2008). These methods detect corpus positions which are marked up incorrectly, and find the correct labels (e.g. pos-tags) for those positions. However, the methods cannot correct errors in structural annotation. This means that they are insufficient to correct annotation errors in a treebank.

This paper proposes a method of correcting errors in structural annotation. Our method is based on a synchronous grammar formalism, called *synchronous tree substitution grammar* (STSG) (Eisner, 2003), which defines a tree-to-tree transfor-

mation. By using an STSG, our method transforms parse trees containing errors into the ones whose errors are corrected. The grammar is automatically induced from the treebank. To select STSG rules which are useful for error correction, we define a score function based on the occurrence frequencies of the rules. An experimental result shows that the selected rules archive high precision.

This paper is organized as follows: Section 2 gives an overview of previous work. Section 3 explains our method of correcting errors in a treebank. Section 4 reports an experimental result using the Penn Treebank.

2 Previous Work

This section summarizes previous methods for correcting errors in corpus annotation and discusses their problem.

Some research addresses the detection of errors in pos-annotation (Nakagawa and Matsumoto, 2002; Dickinson and Meurers, 2003a), syntactic annotation (Dickinson and Meurers, 2003b; Ule and Simov, 2004; Dickinson and Meurers, 2005), and dependency annotation (Boyd et al., 2008). These methods only detect corpus positions where errors occur. It is unclear how we can correct the errors.

Several methods can correct annotation errors (Eskin, 2000; Murata et al., 2005). These methods are to correct tag-annotation errors, that is, they simply suggest a candidate tag for each position where an error is detected. The methods cannot correct syntactic annotation errors, because syntactic annotation is structural. There is no approach to correct structural annotation errors.

To clarify the problem, let us consider an example. Figure 1 depicts two parse trees annotated according to the Penn Treebank annotation ¹. The

¹0 and *T* are null elements.

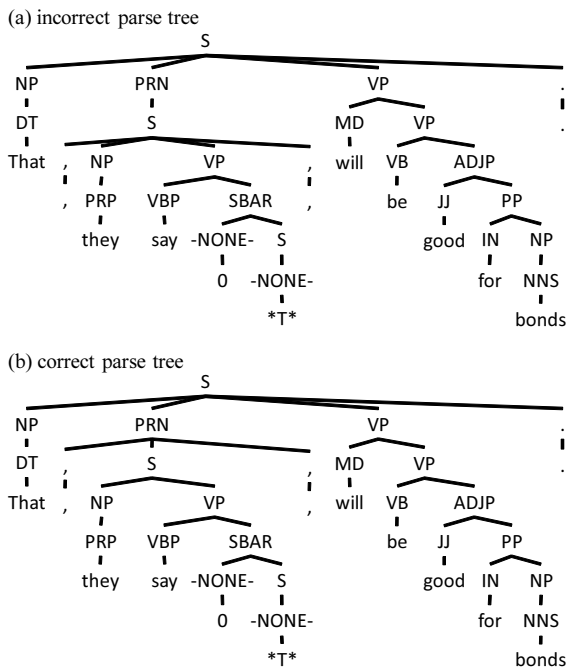


Figure 1: An example of a treebank error

parse tree (a) contains errors and the parse tree (b) is the corrected version. In the parse tree (a), the positions of the two subtrees (,) are erroneous. To correct the errors, we need to move the subtrees to the positions which are directly dominated by the node PRN. This example demonstrates that we need a framework of transforming tree structures to correct structural annotation errors.

3 Correcting Errors by Using Synchronous Grammar

To solve the problem described in Section 2, this section proposes a method of correcting structural annotation errors by using a synchronous tree substitution grammar (STSG) (Eisner, 2003). An STSG defines a tree-to-tree transformation. Our method induces an STSG which transforms parse trees containing errors into the ones whose errors are corrected.

3.1 Synchronous Tree Substitution Grammar

First of all, we describe the STSG formalism. An STSG defines a set of tree pairs. An STSG can be treated as a tree transducer which takes a tree as input and produces a tree as output. Each grammar rule consists of the following elements:

- a pair of trees called *elementary trees*

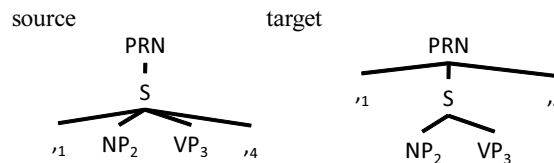


Figure 2: An example of an STSG rule

- a one-to-one alignment between nodes in the elementary trees

For a tree pair $\langle t, t' \rangle$, the tree t and t' are called *source* and *target*, respectively. The non-terminal leaves of elementary trees are called *frontier nodes*. There exists a one-to-one alignment between the frontier nodes in t and t' . The rule means that the structure which matches the source elementary tree is transformed into the structure which is represented by the target elementary tree. Figure 2 shows an example of an STSG rule. The subscripts indicate the alignment. This rule can correct the errors in the parse tree (a) depicted in Figure 1.

An STSG derives tree pairs. Any derivation process starts with the pair of nodes labeled with special symbols called *start symbols*. A derivation proceeds in the following steps:

1. Choose a pair of frontier nodes $\langle \eta, \eta' \rangle$ for which there exists an alignment.
2. Choose a rule $\langle t, t' \rangle$ s.t. $label(\eta) = root(t)$ and $label(\eta') = root(t')$ where $label(\eta)$ is the label of η and $root(t)$ is the root label of t .
3. Substitute t and t' into η and η' , respectively.

Figure 3 shows a derivation process in an STSG.

In the rest of the paper, we focus on the rules in which the source elementary tree is not identical to its target, since such identical rules cannot contribute to error correction.

3.2 Inducing an STSG for Error Correction

This section describes a method of inducing an STSG for error correction. The basic idea of our method is similar to the method presented by Dickinson and Meurers (2003b). Their method detects errors by seeking word sequences satisfying the following conditions:

- The word sequence occurs more than once in the corpus.

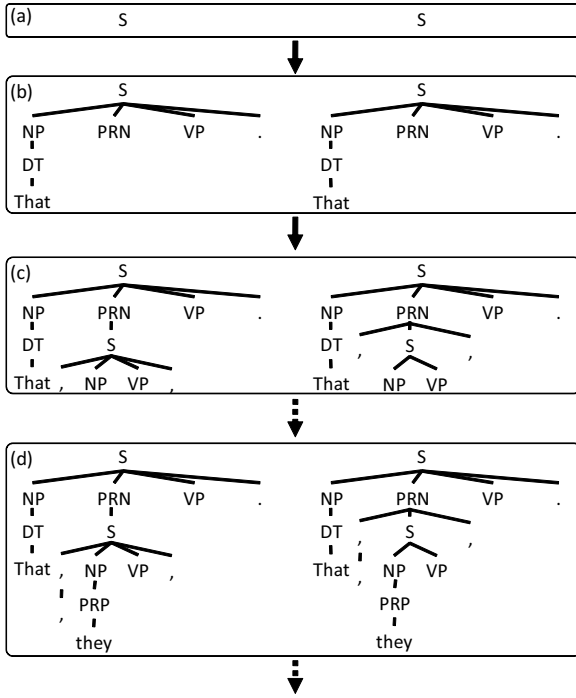


Figure 3: A derivation process of tree pairs in an STSG

- Different syntactic labels are assigned to the occurrences of the word sequence.

Unlike their method, our method seeks word sequences whose occurrences have different partial parse trees. We call a collection of these word sequences with partial parse trees *pseudo parallel corpus*. Moreover, our method extracts STSG rules which transform the one partial tree into the other.

3.2.1 Constructing a Pseudo Parallel Corpus

Our method firstly constructs a pseudo parallel corpus which represents a correspondence between parse trees containing errors and the ones whose errors are corrected. The procedure is as follows: Let T be the set of the parse trees occurring in the corpus. We write $Sub(\sigma)$ for the set which consists of the partial parse trees included in the parse tree σ . A pseudo parallel corpus $Para(T)$ is constructed as follows:

$$\begin{aligned}
 Para(T) = \{ \langle \tau, \tau' \rangle \mid & \tau, \tau' \in \bigcup_{\sigma \in T} Sub(\sigma) \\
 & \wedge \tau \neq \tau' \\
 & \wedge yield(\tau) = yield(\tau') \\
 & \wedge root(\tau) = root(\tau') \}
 \end{aligned}$$

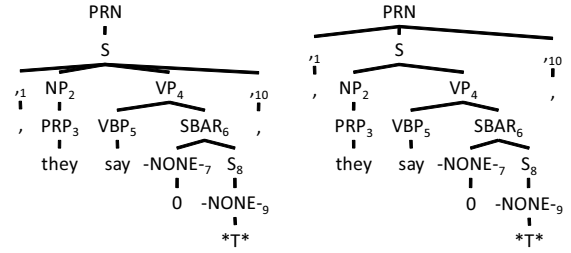


Figure 4: An example of a partial parse tree pair in a pseudo parallel corpus

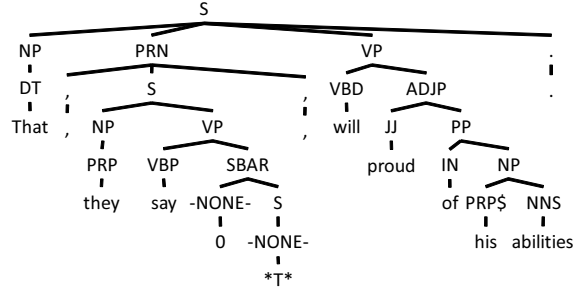


Figure 5: Another example of a parse tree containing a word sequence ‘, they say ,’

where $yield(\tau)$ is the word sequence dominated by τ .

Let us consider an example. If the parse trees depicted in Figure 1 exist in the treebank T , the pair of partial parse trees depicted in Figure 4 is an element of $Para(T)$. We also obtain this pair in the case where there exists not the parse tree (b) depicted in Figure 1 but the parse tree depicted in Figure 5, which contains the word sequence ‘, they say ,’.

3.2.2 Inducing a Grammar from a Pseudo Parallel Corpus

Our method induces an STSG from the pseudo parallel corpus according to the method proposed by Cohn and Lapata (2009). Cohn and Lapata’s method can induce an STSG which represents a correspondence in a parallel corpus. Their method firstly determine an alignment of nodes between pairs of trees in the parallel corpus and extracts STSG rules according to the alignments.

For partial parse trees τ and τ' , we define a node alignment $C(\tau, \tau')$ as follows:

$$\begin{aligned}
 C(\tau, \tau') = \{ \langle \eta, \eta' \rangle \mid & \eta \in Node(\tau) \\
 & \wedge \eta' \in Node(\tau') \\
 & \wedge \eta \text{ is not the root of } \tau
 \end{aligned}$$

$$\begin{aligned} \wedge \eta' \text{ is not the root of } \tau' \\ \wedge \text{label}(\eta) = \text{label}(\eta') \\ \wedge \text{yield}(\eta) = \text{yield}(\eta') \end{aligned}$$

where $\text{Node}(\tau)$ is the set of the nodes in τ , and $\text{yield}(\eta)$ is the word sequence dominated by η . Figure 4 shows an example of a node alignment. The subscripts indicate the alignment.

An STSG rule is extracted by deleting nodes in a partial parse tree pair $\langle \tau, \tau' \rangle \in \text{Para}(T)$. The procedure is as follows:

- For each $\langle \eta, \eta' \rangle \in C(\tau, \tau')$, delete the descendants of η and η' .

For example, the rule shown in Figure 2 is extracted from the pair shown in Figure 4.

3.3 Rule Selection

Some rules extracted by the procedure in Section 3.2 are not useful for error correction, since the pseudo parallel corpus contains tree pairs whose source tree is correct or whose target tree is incorrect. The rules which are extracted from such pairs can be harmful. To select rules which are useful for error correction, we define a score function which is based on the occurrence frequencies of elementary trees in the treebank. The score function is defined as follows:

$$\text{Score}(\langle t, t' \rangle) = \frac{f(t')}{f(t) + f(t')}$$

where $f(\cdot)$ is the occurrence frequency in the treebank. The score function ranges from 0 to 1. We assume that the occurrence frequency of an elementary tree matching incorrect parse trees is very low. According to this assumption, the score function $\text{Score}(\langle t, t' \rangle)$ is high when the source elementary tree t matches incorrect parse trees and the target elementary tree t' matches correct parse trees. Therefore, STSG rules with high scores are regarded to be useful for error correction.

4 An Experiment

To evaluate the effectiveness of our method, we conducted an experiment using the Penn Treebank (Marcus et al., 1993).

We used 49208 sentences in Wall Street Journal sections. We induced STSG rules by applying our method to the corpus. We obtained 8776 rules. We

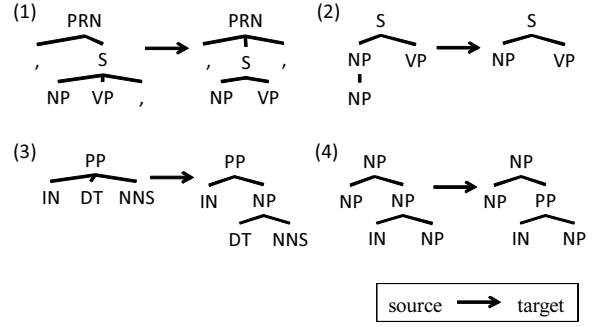


Figure 6: Examples of error correction rules induced from the Penn Treebank

measured the precision of the rules. The precision is defined as follows:

$$\text{precision} = \frac{\# \text{ of the positions where an error is corrected}}{\# \text{ of the positions to which some rule is applied}}$$

We manually checked whether each rule application corrected an error, because the corrected treebank does not exist². Furthermore, we only evaluated the first 100 rules which are ordered by the score function described in Section 3.3, since it is time-consuming and expensive to evaluate all of the rules. These 100 rules were applied at 331 positions. The precision of the rules is 71.9%. For each rule, we measured the precision of it. 70 rules achieved 100% precision. These results demonstrate that our method can correct syntactic annotation errors with high precision. Moreover, 30 rules of the 70 rules transformed bracketed structures. This fact shows that the treebank contains structural errors which cannot be dealt with by the previous methods.

Figure 6 depicts examples of error correction rules which achieved 100% precision. Rule (1), (2) and (3) are rules which transform bracketed structures. Rule (4) simply replaces a node label. Rule (1) corrects an erroneous position of a comma (see Figure 7 (a)). Rule (2) deletes a useless node NP in a subject position (see Figure 7 (b)). Rule (3) inserts a node NP (see Figure 7 (c)). Rule (4) replaces a node label NP with the correct label PP (see Figure 7 (d)). These examples demonstrate that our method can correct syntactic annotation errors.

Figure 8 depicts an example where our method detected an annotation error but could not correct it. To correct the error, we need to attach the node

²This also means that we cannot measure the recall of the rules.

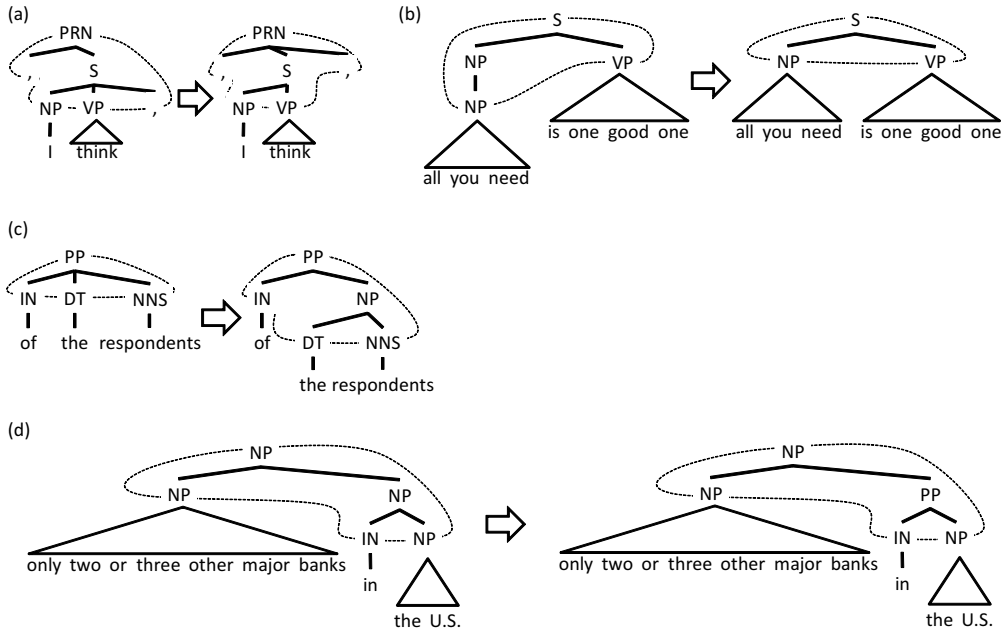


Figure 7: Examples of correcting syntactic annotation errors

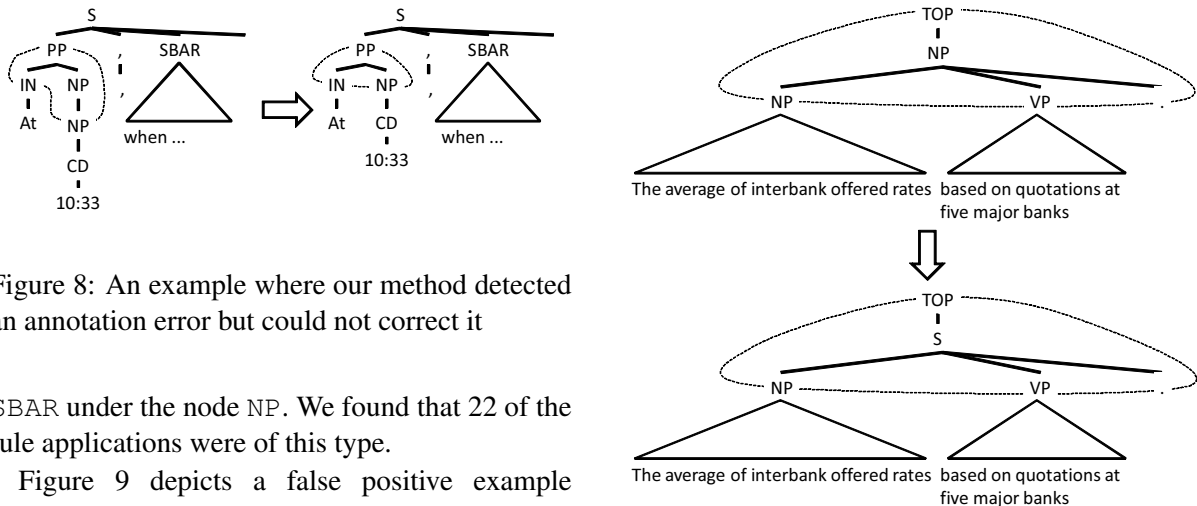


Figure 8: An example where our method detected an annotation error but could not correct it

SBAR under the node NP. We found that 22 of the rule applications were of this type.

Figure 9 depicts a false positive example where our method mistakenly transformed a correct syntactic structure. The score of the rule is very high, since the source elementary tree (TOP (NP NP VP .)) is less frequent. This example shows that our method has a risk of changing correct annotations of less frequent syntactic structures.

5 Conclusion

This paper proposes a method of correcting errors in a treebank by using a synchronous tree substitution grammar. Our method constructs a pseudo parallel corpus from the treebank and extracts STSG rules from the parallel corpus. The experimental result demonstrates that we can obtain error correction rules with high precision.

Figure 9: A false positive example where a correct syntactic structure was mistakenly transformed

In future work, we will explore a method of increasing the recall of error correction by constructing a wide-coverage STSG.

Acknowledgements

This research is partially supported by the Grant-in-Aid for Scientific Research (B) (No. 22300051) of JSPS and by the Kayamori Foundation of Informational Science Advancement.

References

- Adriane Boyd, Markus Dickinson, and Detmar Meurers. 2008. On detecting errors in dependency treebanks. *Research on Language and Computation*, 6(2):113–137.
- Trevor Cohn and Mirella Lapata. 2009. Sentence compression as tree transduction. *Journal of Artificial Intelligence Research*, 34(1):637–674.
- Markus Dickinson and Detmar Meurers. 2003a. Detecting errors in part-of-speech annotation. In *Proceedings of the 10th Conference of the European Chapter of the Association for Computational Linguistics*, pages 107–114.
- Markus Dickinson and Detmar Meurers. 2003b. Detecting inconsistencies in treebanks. In *Proceedings of the Second Workshop on Treebanks and Linguistic Theories*.
- Markus Dickinson and W. Detmar Meurers. 2005. Prune diseased branches to get healthy trees! how to find erroneous local trees in a treebank and why it matters. In *Proceedings of the 4th Workshop on Treebanks and Linguistic Theories*.
- Jason Eisner. 2003. Learning non-isomorphic tree mappings for machine translation. In *Proceedings of the 41st Annual Meeting of the Association for Computational Linguistics, Companion Volume*, pages 205–208.
- Eleazar Eskin. 2000. Detecting errors within a corpus using anomaly detection. In *Proceedings of the 1st North American chapter of the Association for Computational Linguistics Conference*, pages 148–153.
- Mitchell P. Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. 1993. Building a large annotated corpus of English: the Penn Treebank. *Computational Linguistics*, 19(2):310–330.
- Masaki Murata, Masao Utiyama, Kiyotaka Uchimoto, Hitoshi Isahara, and Qing Ma. 2005. Correction of errors in a verb modality corpus for machine translation with a machine-learning method. *ACM Transactions on Asian Language Information Processing*, 4(1):18–37.
- Tetsuji Nakagawa and Yuji Matsumoto. 2002. Detecting errors in corpora using support vector machines. In *Proceedings of the 19th International Conference on Computational Linguistics*, pages 709–715.
- Tylman Ule and Kiril Simov. 2004. Unexpected productions may well be errors. In *Proceedings of 4th International Conference on Language Resources and Evaluation*, pages 1795–1798.

Evaluating Machine Translations using mNCD

Marcus Dobrinkat and Tero Tapiovaara and Jaakko Väyrynen

Adaptive Informatics Research Centre

Aalto University School of Science and Technology

P.O. Box 15400, FI-00076 Aalto, Finland

{marcus.dobrinkat, jaakko.j.vayrynen, tero.tapiovaara}@tkk.fi

Kimmo Kettunen

Kymenlaakso University of Applied Sciences

P.O. Box 9, FI-48401 Kotka, Finland

kimmo.kettunen@kyamk.fi

Abstract

This paper introduces mNCD, a method for automatic evaluation of machine translations. The measure is based on normalized compression distance (NCD), a general information theoretic measure of string similarity, and flexible word matching provided by stemming and synonyms. The mNCD measure outperforms NCD in system-level correlation to human judgments in English.

1 Introduction

Automatic evaluation of machine translation (MT) systems requires automated procedures to ensure consistency and efficient handling of large amounts of data. In statistical MT systems, automatic evaluation of translations is essential for parameter optimization and system development. Human evaluation is too labor intensive, time consuming and expensive for daily evaluations. However, manual evaluation is important in the comparison of different MT systems and for the validation and development of automatic MT evaluation measures, which try to model human assessments of translations as closely as possible. Furthermore, the ideal evaluation method would be language independent, fast to compute and simple.

Recently, normalized compression distance (NCD) has been applied to the evaluation of machine translations. NCD is a general information theoretic measure of string similarity, whereas most MT evaluation measures, e.g., BLEU and METEOR, are specifically constructed for the task. Parker (2008) introduced BADGER, an MT evaluation measure that uses NCD and a language independent word normalization

method. BADGER scores were directly compared against the scores of METEOR and word error rate (WER). The correlation between BADGER and METEOR were low and correlations between BADGER and WER high. Kettunen (2009) uses the NCD directly as an MT evaluation measure. He showed with a small corpus of three language pairs that NCD and METEOR 0.6 correlated for translations of 10–12 MT systems. NCD was not compared to human assessments of translations, but correlations of NCD and METEOR scores were very high for all the three language pairs.

Väyrynen et al. (2010) have extended the work by including NCD in the ACL WMT08 evaluation framework and showing that NCD is correlated to human judgments. The NCD measure did not match the performance of the state-of-the-art MT evaluation measures in English, but it presented a viable alternative to de facto standard BLEU (Papineni et al., 2001), which is simple and effective but has been shown to have a number of drawbacks (Callison-Burch et al., 2006).

Some recent advances in automatic MT evaluation have included non-binary matching between compared items (Banerjee and Lavie, 2005; Agarwal and Lavie, 2008; Chan and Ng, 2009), which is implicitly present in the string-based NCD measure. Our motivation is to investigate whether including additional language dependent resources would improve the NCD measure. We experiment with relaxed word matching using stemming and a lexical database to allow lexical changes. These additional modules attempt to make the reference sentences more similar to the evaluated translations on the string level. We report an experiment showing that document-level NCD and aggregated NCD scores for individual sentences produce very similar correlations to human judgments.

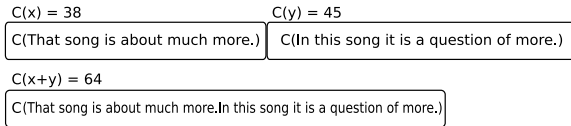


Figure 1: An example showing the compressed sizes of two strings separately and concatenated.

2 Normalized Compression Distance

Normalized compression distance (NCD) is a similarity measure based on the idea that a string x is similar to another string y when both share substrings. The description of y can reference shared substrings in the known x without repetition, indicating shared information. Figure 1 shows an example in which the compression of the concatenation of x and y results in a shorter output than individual compressions of x and y .

The normalized compression distance, as defined by Cilibrasi and Vitanyi (2005), is given in Equation 1, with $C(x)$ as length of the compression of x and $C(x, y)$ as the length of the compression of the concatenation of x and y .

$$NCD(x, y) = \frac{C(x, y) - \min \{C(x), C(y)\}}{\max \{C(x), C(y)\}} \quad (1)$$

NCD computes the distance as a score closer to one for very different strings and closer to zero for more similar strings.

NCD is an approximation of the uncomputable normalized information distance (NID), a general measure for the similarity of two objects. NID is based on the notion of Kolmogorov complexity $K(x)$, a theoretical measure for the information content of a string x , defined as the shortest universal Turing machine that prints x and stops (Solomonoff, 1964). NCD approximates NID by the use of a compressor $C(x)$ that is an upper bound of the Kolmogorov complexity $K(x)$.

3 mNCD

Normalized compression distance was not conceived with MT evaluation in mind, but rather it is a general measure of string similarity. Implicit non-binary matching with NCD is indicated by preliminary experiments which show that NCD is less sensitive to random changes on the character level than, for instance, BLEU, which only counts the exact matches between word n-grams. Thus comparison of sentences at the character level could account better for morphological changes.

Variation in language leads to several acceptable translations for each source sentence, which is why multiple reference translations are preferred in evaluation. Unfortunately, it is typical to have only one reference translation. Paraphrasing techniques can produce additional translation variants (Russo-Lassner et al., 2005; Kauchak and Barzilay, 2006). These can be seen as new reference translations, similar to pseudo references (Ma et al., 2007).

The proposed method, mNCD, works analogously to M-BLEU and M-TER, which use the flexible word matching modules from METEOR to find relaxed word-to-word alignments (Agarwal and Lavie, 2008). The modules are able to align words even if they do not share the same surface form, but instead have a common stem or are synonyms of each other. A similarized translation reference is generated by replacing words in the reference with their aligned counterparts from the translation hypothesis. The NCD score is computed between the translations and the similarized references to get the mNCD score.

Table 1 shows some hand-picked German-English candidate translations along with a) the reference translations including the 1-NCD score to easily compare with METEOR and b) the similarized references including the mNCD score. For comparison, the corresponding METEOR scores without implicit relaxed matching are shown.

4 Experiments

The proposed mNCD and the basic NCD measure were evaluated by computing correlation to human judgments of translations. A high correlation value between an MT evaluation measure and human judgments indicates that the measure is able to evaluate translations in a more similar way to humans.

Relaxed alignments with the METEOR modules `exact`, `stem` and `synonym` were created for English for the computation of the mNCD score. The `synonym` module was not available with other target languages.

4.1 Evaluation Data

The 2008 ACL Workshop on Statistical Machine Translation (Callison-Burch et al., 2008) shared task data includes translations from a total of 30 MT systems between English and five European languages, as well as automatic and human trans-

	Candidate C/ Reference R/ Similarized Reference S	1-NCD	METEOR
C	There is no effective means to stop a Tratsch, which was already included in the world.		
R	There is no good way to halt gossip that has already begun to spread.	.41	.31
S	There is no effective means to stop gossip that has already begun to spread.	.56	.55
C	Crisis, not only in America		
R	A Crisis Not Only in the U.S.	.51	.44
S	A Crisis not only in the America	.72	.56
C	Influence on the whole economy should not have this crisis.		
R	Nevertheless, the crisis should not have influenced the entire economy.	.60	.37
S	Nevertheless, the crisis should not have Influence the entire economy.	.62	.44
C	Or the lost tight meeting will be discovered at the hands of a gentlemen?		
R	Perhaps you see the pen you thought you lost lying on your colleague’s desk.	.42	.09
S	Perhaps you meeting the pen you thought you lost lying on your colleague’s desk.	.40	.13

Table 1: Example German–English translations showing the effect of relaxed matching in the 1-mNCD score (for rows S) compared with METEOR using the `exact` module only, since the modules `stem` and `synonym` are already used in the similarized reference. Replaced words are emphasized.

lation evaluations for the translations. There are several tasks, defined by the language pair and the domain of translated text.

The human judgments include three different categories. The RANK category has human quality rankings of five translations for one sentence from different MT systems. The CONST category contains rankings for short phrases (constituents), and the YES/NO category contains binary answers if a short phrase is an acceptable translation or not.

For the translation tasks into English, the relaxed alignment using a `stem` module and the `synonym` module affected 7.5% of all words, whereas only 5.1% of the words were changed in the tasks from English into the other languages.

The data was preprocessed in two different ways. For NCD we kept the data as is, which we called real casing (rc). Since the used METEOR align module lowercases all text, we restored the case information in mNCD by copying the correct case from the reference translation to the similarized reference, based on METEOR’s alignment. The other way was to lowercase all data (lc).

4.2 System-level correlation

We follow the same evaluation methodology as in Callison-Burch et al. (2008), which allows us to measure how well MT evaluation measures correlate with human judgments on the system level.

Spearman’s rank correlation coefficient ρ was calculated between each MT evaluation measure and human judgment category using the simplified equation

$$\rho = 1 - \frac{6 \sum_i d_i}{n(n^2 - 1)} \quad (2)$$

where for each system i , d_i is the difference between the rank derived from annotators’ input and the rank obtained from the measure. From the annotators’ input, the n systems were ranked based on the number of times each system’s output was selected as the best translation divided by the number of times each system was part of a judgment.

We computed system-level correlations for tasks with English, French, Spanish and German as the target language¹.

5 Results

We compare mNCD against NCD and relate their performance to other MT evaluation measures.

5.1 Block size effect on NCD scores

Väyrynen et al. (2010) computed NCD between a set of candidate translations and references at the same time regardless of the sentence alignments, analogously to document comparison. We experimented with segmentation of the candidate translations into smaller blocks, which were individually evaluated with NCD and aggregated into a single value with arithmetic mean. The resulting system-level correlations between NCD and human judgments are shown in Figure 2 as a function of the block size. The correlations are very similar with all block sizes, except for Spanish, where smaller block size produces higher correlation. An experiment with geometric mean produced similar results. The reported results with mNCD use maximum block size, similar to Väyrynen et al. (2010).

¹The English-Spanish news task was left out as most measures had negative correlation with human judgments.

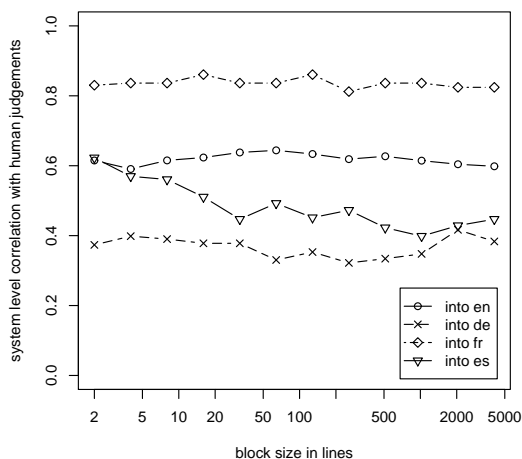


Figure 2: The block size has very little effect on the correlation between NCD and human judgements. The right side corresponds to document comparison and the left side to aggregated NCD scores for sentences.

5.2 mNCD against NCD

Table 2 shows the average system level correlation of different NCD and mNCD variants for translations into English. The two compressors that worked best in our experiments were PPMZ and bz2. PPMZ is slower to compute but performs slightly better compared to bz2, except for the

Method	Parameters	RANK	CONST	YES/NO	Mean
mNCD	PPMZ rc	.69	.74	.80	.74
NCD	PPMZ rc	.60	.66	.71	.66
mNCD	bz2 rc	.64	.73	.73	.70
NCD	bz2 rc	.57	.64	.69	.64
mNCD	PPMZ lc	.66	.80	.79	.75
NCD	PPMZ lc	.56	.79	.75	.70
mNCD	bz2 lc	.59	.85	.74	.73
NCD	bz2 lc	.54	.82	.71	.69

Table 2: Mean system level correlations over all translation tasks into English for variants of mNCD and NCD. Higher values are emphasized. Parameters are the compressor PPMZ or bz2 and the preprocessing choice lowercasing (lc) or real casing (rc).

Method	Parameters		Target Lang Corr			
			EN	DE	FR	Es
mNCD	PPMZ	rc	.69	.37	.82	.38
NCD	PPMZ	rc	.60	.37	.84	.39
mNCD	bz2	rc	.64	.32	.75	.25
NCD	bz2	rc	.57	.34	.85	.42
mNCD	PPMZ	lc	.66	.33	.79	.23
NCD	PPMZ	lc	.56	.37	.77	.21
mNCD	bz2	lc	.59	.25	.78	.16
NCD	bz2	lc	.54	.26	.77	.15

Table 3: mNCD versus NCD system correlation RANK results with different parameters (the same as in Table 2) for each target language. Higher values are emphasized. Target languages DE, FR and ES use only the stem module.

lowercased CONST category.

Table 2 shows that real casing improves RANK correlation slightly throughout NCD and mNCD variants, whereas it reduces correlation in the categories CONST, YES/NO as well as the mean.

The best mNCD (PPMZ rc) improves the best NCD (PPMZ rc) method by 15% in the RANK category. In the CONST category the best mNCD (bz2 lc) improves the best NCD (bz2 lc) by 3.7%. For the total average, the best mNCD (PPMZ rc) improves the the best NCD (bz2 lc) by 7.2%.

Table 3 shows the correlation results for the RANK category by target language. As shown already in Table 2, mNCD clearly outperforms NCD for English. Correlations for other languages show mixed results and on average, mNCD gives lower correlations than NCD.

5.3 mNCD versus other methods

Table 4 presents the results for the selected mNCD (PPMZ rc) and NCD (bz2 rc) variants along with the correlations for other MT evaluation methods from the WMT'08 data, based on the results in Callison-Burch et al. (2008). The results are averages over language pairs into English, sorted by RANK, which we consider the most significant category. Although mNCD correlation with human evaluations improved over NCD, the ranking among other measures was not affected. Language and task specific results not shown here, reveal very low mNCD and NCD correlations in the Spanish-English news task, which significantly

Method	RANK	CONST	YES/NO	Mean
DP	.81	.66	.74	.73
ULCh	.80	.68	.78	.75
DR	.79	.53	.65	.66
meteor-ranking	.78	.55	.63	.65
ULC	.77	.72	.81	.76
posbleu	.75	.69	.78	.74
SR	.75	.66	.76	.72
posF4gram-gm	.74	.60	.71	.68
meteor-baseline	.74	.60	.63	.66
posF4gram-am	.74	.58	.69	.67
mNCD (PPMZ rc)	.69	.74	.80	.74
NCD (PPMZ rc)	.60	.66	.71	.66
mbleu	.50	.76	.70	.65
bleu	.50	.72	.74	.65
mter	.38	.74	.68	.60
svm-rank	.37	.10	.23	.23
Mean	.67	.62	.69	.66

Table 4: Average system-level correlations over translation tasks into English for NCD, mNCD and other MT evaluations measures

degrades the averages. Considering the mean of the categories instead, mNCD’s correlation of .74 is third best together with ‘posbleu’.

Table 5 shows the results from English. The table is shorter since many of the better MT measures use language specific linguistic resources that are not easily available for languages other than English. mNCD performs competitively only for French, otherwise it falls behind NCD and other methods as already shown earlier.

6 Discussion

We have introduced a new MT evaluation measure, mNCD, which is based on normalized compression distance and METEOR’s relaxed alignment modules. The mNCD measure outperforms NCD in English with all tested parameter combinations, whereas results with other target languages are unclear. The improved correlations with mNCD did not change the position in the RANK category of the MT evaluation measures in the 2008 ACL WMT shared task.

The improvement in English was expected on the grounds of the synonym module, and indicated also by the larger number of affected words in the

Method	Target Lang Corr			
	DE	FR	ES	Mean
posbleu	.75	.80	.75	.75
posF4gram-am	.74	.82	.79	.74
posF4gram-gm	.74	.82	.79	.74
bleu	.47	.83	.80	.68
NCD (bz2 rc)	.34	.85	.42	.66
svm-rank	.44	.80	.80	.66
mbleu	.39	.77	.83	.63
mNCD (PPMZ rc)	.37	.82	.38	.63
meteor-baseline	.43	.61	.84	.58
meteor-ranking	.26	.70	.83	.55
mter	.26	.69	.73	.52
Mean	.47	.77	.72	.65

Table 5: Average system-level correlations for the RANK category from English for NCD, mNCD and other MT evaluation measures.

similarized references. We believe there is potential for improvement in other languages as well if synonym lexicons are available.

We have also extended the basic NCD measure to scale between a document comparison measure and aggregated sentence-level measure. The rather surprising result is that NCD produces quite similar scores with all block sizes. The different result with Spanish may be caused by differences in the data or problems in the calculations.

After using the same evaluation methodology as in Callison-Burch et al. (2008), we have doubts whether it presents the most effective method exploiting all the given human evaluations in the best way. The system-level correlation measure only awards the winner of the ranking of five different systems. If a system always scored second, it would never be awarded and therefore be overly penalized. In addition, the human knowledge that gave the lower rankings is not exploited.

In future work with mNCD as an MT evaluation measure, we are planning to evaluate synonym dictionaries for other languages than English. The synonym module for English does not distinguish between different senses of words. Therefore, synonym lexicons found with statistical methods might provide a viable alternative for manually constructed lexicons (Kauchak and Barzilay, 2006).

References

- Abhaya Agarwal and Alon Lavie. 2008. METEOR, M-BLEU and M-TER: evaluation metrics for high-correlation with human rankings of machine translation output. In *StatMT '08: Proceedings of the Third Workshop on Statistical Machine Translation*, pages 115–118, Morristown, NJ, USA. Association for Computational Linguistics.
- Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan, June. Association for Computational Linguistics.
- Chris Callison-Burch, Miles Osborne, and Philipp Koehn. 2006. Re-evaluating the role of BLEU in machine translation research. In *Proceedings of EACL-2006*, pages 249–256.
- Chris Callison-Burch, Cameron Fordyce, Philipp Koehn, Christoph Monz, and Josh Schroeder. 2008. Further meta-evaluation of machine translation. *ACL Workshop on Statistical Machine Translation*.
- Yee Seng Chan and Hwee Tou Ng. 2009. MaxSim: performance and effects of translation fluency. *Machine Translation*, 23(2-3):157–168.
- Rudi Cilibrasi and Paul Vitanyi. 2005. Clustering by compression. *IEEE Transactions on Information Theory*, 51:1523–1545.
- David Kauchak and Regina Barzilay. 2006. Paraphrasing for automatic evaluation. In *Proceedings of the main conference on Human Language Technology Conference of the North American Chapter of the Association of Computational Linguistics*, pages 455–462, Morristown, NJ, USA. Association for Computational Linguistics.
- Kimmo Kettunen. 2009. Packing it all up in search for a language independent MT quality measure tool. In *In Proceedings of LTC-09, 4th Language and Technology Conference*, pages 280–284, Poznan.
- Yanjun Ma, Nicolas Stroppa, and Andy Way. 2007. Bootstrapping word alignment via word packing. In *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, pages 304–311, Prague, Czech Republic, June. Association for Computational Linguistics.
- K. Papineni, S. Roukos, T. Ward, and W. Zhu. 2001. BLEU: a method for automatic evaluation of machine translation. Technical Report RC22176 (W0109-022), IBM Research Division, Thomas J. Watson Research Center.
- Steven Parker. 2008. BADGER: A new machine translation metric. In *Metrics for Machine Translation Challenge 2008*, Waikiki, Hawai'i, October. AMTA.
- Grazia Russo-Lassner, Jimmy Lin, and Philip Resnik. 2005. A paraphrase-based approach to machine translation evaluation. Technical Report LAMP-TR-125/CS-TR-4754/UMIACS-TR-2005-57, University of Maryland, College Park.
- Ray Solomonoff. 1964. Formal theory of inductive inference. Part I. *Information and Control*, 7(1):1–22.
- Jaakko J. Väyrynen, Tero Tapiovaara, Kimmo Kettunen, and Marcus Dobrinkat. 2010. Normalized compression distance as an automatic MT evaluation metric. In *Proceedings of MT 25 years on*. To appear.

Tackling Sparse Data Issue in Machine Translation Evaluation *

Ondřej Bojar, Kamil Kos, and David Mareček

Charles University in Prague, Institute of Formal and Applied Linguistics
{bojar, marecek}@ufal.mff.cuni.cz, kamilkos@email.cz

Abstract

We illustrate and explain problems of n -grams-based machine translation (MT) metrics (e.g. BLEU) when applied to morphologically rich languages such as Czech. A novel metric SemPOS based on the deep-syntactic representation of the sentence tackles the issue and retains the performance for translation to English as well.

1 Introduction

Automatic metrics of machine translation (MT) quality are vital for research progress at a fast pace. Many automatic metrics of MT quality have been proposed and evaluated in terms of correlation with human judgments while various techniques of manual judging are being examined as well, see e.g. MetricsMATR08 (Przybocki et al., 2008)¹, WMT08 and WMT09 (Callison-Burch et al., 2008; Callison-Burch et al., 2009)².

The contribution of this paper is twofold. Section 2 illustrates and explains severe problems of a widely used BLEU metric (Papineni et al., 2002) when applied to Czech as a representative of languages with rich morphology. We see this as an instance of the sparse data problem well known for MT itself: too much detail in the formal representation leading to low coverage of e.g. a translation dictionary. In MT evaluation, too much detail leads to the lack of comparable parts of the hypothesis and the reference.

* This work has been supported by the grants EuroMatrixPlus (FP7-ICT-2007-3-231720 of the EU and 7E09003 of the Czech Republic), FP7-ICT-2009-4-247762 (Faust), GA201/09/H057, GAUK 1163/2010, and MSM 0021620838. We are grateful to the anonymous reviewers for further research suggestions.

¹<http://nist.gov/speech/tests/metricmatr/2008/results/>

²<http://www.statmt.org/wmt08> and [wmt09](http://www.statmt.org/wmt09)

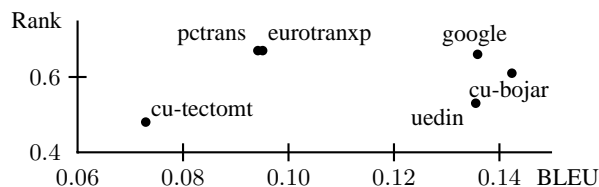


Figure 1: BLEU and human ranks of systems participating in the English-to-Czech WMT09 shared task.

Section 3 introduces and evaluates some new variations of SemPOS (Kos and Bojar, 2009), a metric based on the deep syntactic representation of the sentence performing very well for Czech as the target language. Aside from including dependency and n -gram relations in the scoring, we also apply and evaluate SemPOS for English.

2 Problems of BLEU

BLEU (Papineni et al., 2002) is an established language-independent MT metric. Its correlation to human judgments was originally deemed high (for English) but better correlating metrics (esp. for other languages) were found later, usually employing language-specific tools, see e.g. Przybocki et al. (2008) or Callison-Burch et al. (2009). The unbeaten advantage of BLEU is its simplicity.

Figure 1 illustrates a very low correlation to human judgments when translating to Czech. We plot the official BLEU score against the rank established as the percentage of sentences where a system ranked no worse than all its competitors (Callison-Burch et al., 2009). The systems developed at Charles University (cu-) are described in Bojar et al. (2009), uedin is a vanilla configuration of Moses (Koehn et al., 2007) and the remaining ones are commercial MT systems.

In a manual analysis, we identified the reasons for the low correlation: BLEU is overly sensitive to *sequences* and *forms* in the hypothesis matching

Con- firmed	Error Flags	1-grams	2-grams	3-grams	4-grams
Yes	Yes	6.34%	1.58%	0.55%	0.29%
Yes	No	36.93%	13.68%	5.87%	2.69%
No	Yes	22.33%	41.83%	54.64%	63.88%
No	No	34.40%	42.91%	38.94%	33.14%
Total n -grams		35,531	33,891	32,251	30,611

Table 1: n -grams confirmed by the reference and containing error flags.

the reference translation. This focus goes directly against the properties of Czech: relatively free word order allows many permutations of words and rich morphology renders many valid word forms not confirmed by the reference.³ These problems are to some extent mitigated if several reference translations are available, but this is often not the case.

Figure 2 illustrates the problem of “sparse data” in the reference. Due to the lexical and morphological variance of Czech, only a single word in each hypothesis matches a word in the reference. In the case of pctrans, the match is even a false positive, “do” (to) is a preposition that should be used for the “minus” phrase and not for the “end of the day” phrase. In terms of BLEU, both hypotheses are equally poor but 90% of their tokens were not evaluated.

Table 1 estimates the overall magnitude of this issue: For 1-grams to 4-grams in 1640 instances (different MT outputs and different annotators) of 200 sentences with manually flagged errors⁴, we count how often the n -gram is confirmed by the reference and how often it contains an error flag. The suspicious cases are n -grams confirmed by the reference but still containing a flag (false positives) and n -grams not confirmed despite containing no error flag (false negatives).

Fortunately, there are relatively few false positives in n -gram based metrics: 6.3% of unigrams and far fewer higher n -grams.

The issue of false negatives is more serious and confirms the problem of sparse data if only one reference is available. 30 to 40% of n -grams do not contain any error and yet they are not con-

³Condon et al. (2009) identify similar issues when evaluating translation to Arabic and employ rule-based normalization of MT output to improve the correlation. It is beyond the scope of this paper to describe the rather different nature of morphological richness in Czech, Arabic and also other languages, e.g. German or Finnish.

⁴The dataset with manually flagged errors is available at <http://ufal.mff.cuni.cz/euromatrixplus/>

firmed by the reference. This amounts to 34% of running unigrams, giving enough space to differ in human judgments and still remain unscored.

Figure 3 documents the issue across languages: the lower the BLEU score itself (i.e. fewer confirmed n -grams), the lower the correlation to human judgments regardless of the target language (WMT09 shared task, 2025 sentences per language).

Figure 4 illustrates the overestimation of scores caused by too much attention to sequences of tokens. A phrase-based system like Moses (cubojar) can sometimes produce a long sequence of tokens exactly as required by the reference, leading to a high BLEU score. The framed words in the illustration are not confirmed by the reference, but the actual error in these words is very severe for comprehension: nouns were used twice instead of finite verbs, and a misleading translation of a preposition was chosen. The output by pctrans preserves the meaning much better despite not scoring in either of the finite verbs and producing far shorter confirmed sequences.

3 Extensions of SemPOS

SemPOS (Kos and Bojar, 2009) is inspired by metrics based on overlapping of linguistic features in the reference and in the translation (Giménez and Márquez, 2007). It operates on so-called “tectogrammatical” (deep syntactic) representation of the sentence (Sgall et al., 1986; Hajič et al., 2006), formally a dependency tree that includes only autosemantic (content-bearing) words.⁵ SemPOS as defined in Kos and Bojar (2009) disregards the syntactic structure and uses the semantic part of speech of the words (noun, verb, etc.). There are 19 fine-grained parts of speech. For each semantic part of speech t , the overlapping $O(t)$ is set to zero if the part of speech does not occur in the reference or the candidate set and otherwise it is computed as given in Equation 1 below.

⁵We use TectoMT (Žabokrtský and Bojar, 2008), <http://ufal.mff.cuni.cz/tectomt/>, for the linguistic pre-processing. While both our implementation of SemPOS as well as TectoMT are in principle freely available, a stable public version has yet to be released. Our plans include experiments with approximating the deep syntactic analysis with a simple tagger, which would also decrease the installation burden and computation costs, at the expense of accuracy.

SRC	Prague Stock Market falls to minus by the end of the trading day
REF	pražská burza se ke konci obchodování propadla do minusu
cu-bojar	praha stock market klesne k minus na konci obchodního dne
pctrans	praha trh cenných papírů padá minus do konce obchodního dne

Figure 2: Sparse data in BLEU evaluation: Large chunks of hypotheses are not compared at all. Only a single unigram in each hypothesis is confirmed in the reference.

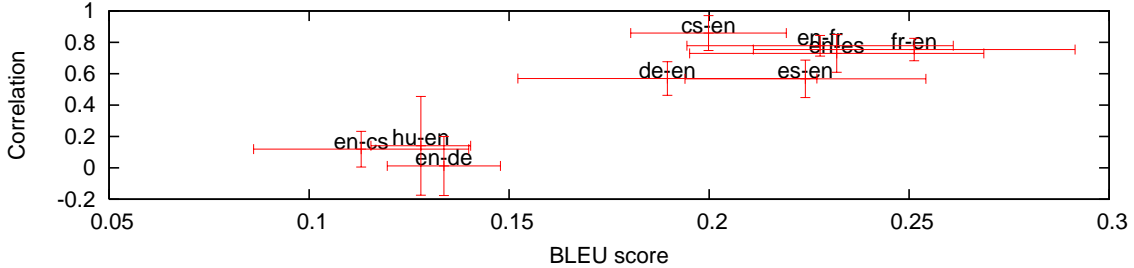


Figure 3: BLEU correlates with its correlation to human judgments. BLEU scores around 0.1 predict little about translation quality.

$$O(t) = \frac{\sum_{i \in I} \sum_{w \in r_i \cap c_i} \min(\text{cnt}(w, t, r_i), \text{cnt}(w, t, c_i))}{\sum_{i \in I} \sum_{w \in r_i \cup c_i} \max(\text{cnt}(w, t, r_i), \text{cnt}(w, t, c_i))} \quad (1)$$

The semantic part of speech is denoted t ; c_i and r_i are the candidate and reference translations of sentence i , and $\text{cnt}(w, t, rc)$ is the number of words w with type t in rc (the reference or the candidate). The matching is performed on the level of lemmas, i.e. no morphological information is preserved in ws . See Figure 5 for an example; the sentence is the same as in Figure 4.

The final SemPOS score is obtained by macro-averaging over all parts of speech:

$$\text{SemPOS} = \frac{1}{|T|} \sum_{t \in T} O(t) \quad (2)$$

where T is the set of all possible semantic parts of speech types. (The degenerate case of blank candidate and reference has SemPOS zero.)

3.1 Variations of SemPOS

This section describes our modifications of SemPOS. All methods are evaluated in Section 3.2.

Different Classification of Autosemantic Words. SemPOS uses semantic parts of speech to classify autosemantic words. The tectogrammatical layer offers also a feature called *Functor* describing the relation of a word to its governor

similarly as semantic roles do. There are 67 functor types in total.

Using *Functor* instead of *SemPOS* increases the number of word classes that independently require a high overlap. For a contrast we also completely remove the classification and use only one global class (*Void*).

Deep Syntactic Relations in SemPOS. In SemPOS, an autosemantic word of a class is confirmed if its lemma matches the reference. We utilize the dependency relations at the tectogrammatical layer to validate valence by refining the overlap and requiring also the lemma of 1) the parent (denoted “par”), or 2) all the children regardless of their order (denoted “sons”) to match.

Combining BLEU and SemPOS. One of the major drawbacks of *SemPOS* is that it completely ignores word order. This is too coarse even for languages with relatively free word order like Czech. Another issue is that it operates on lemmas and it completely disregards correct word forms. Thus, a weighted linear combination of SemPOS and BLEU (computed on the surface representation of the sentence) should compensate for this. For the purposes of the combination, we compute BLEU *only* on unigrams up to fourgrams (denoted $\text{BLEU}_1, \dots, \text{BLEU}_4$) but including the brevity penalty as usual. Here we try only a few weight settings in the linear combination but given a held-out dataset, one could optimize the weights for the best performance.

SRC	Congress yields: US government can pump 700 billion dollars into banks					
REF	kongres ustoupil : vláda usa může do bank napumpovat 700 miliard dolarů					
cu-bojar	kongres	výnosy	: vláda usa může	čerpadlo	700 miliard dolarů	v bankách
pctrans	<u>kongres</u>	<u>vynáší</u>	<u>: us vláda může</u>	<u>čerpat</u>	<u>700 miliardu dolarů</u>	<u>do bank</u>

Figure 4: Too much focus on sequences in BLEU: pctrans’ output is better but does not score well. BLEU gave credit to cu-bojar for 1, 3, 5 and 8 fourgrams, trigrams, bigrams and unigrams, resp., but only for 0, 0, 1 and 8 n -grams produced by pctrans. Confirmed sequences of tokens are underlined and important errors (not considered by BLEU) are framed.

REF	<u>kongres/n</u> <u>ustoupit/v</u> :/n <u>vláda/n</u> <u>usa/n</u> <u>banka/n</u> <u>napumpovat/v</u> <u>700/n</u> <u>miliarda/n</u> <u>dolar/n</u>
cu-bojar	<u>kongres/n</u> <u>výnos/n</u> :/n <u>vláda/n</u> <u>usa/n</u> <u>moci/v</u> <u>čerpadlo/n</u> <u>700/n</u> <u>miliarda/n</u> <u>dolar/n</u> <u>banka/n</u>
pctrans	<u>kongres/n</u> <u>vynášet/v</u> :/n <u>us/n</u> <u>vláda/n</u> <u>čerpat/v</u> <u>700/n</u> <u>miliarda/n</u> <u>dolar/n</u> <u>banka/n</u>

Figure 5: SemPOS evaluates the overlap of lemmas of autosemantic words given their semantic part of speech (n, v, ...). Underlined words are confirmed by the reference.

SemPOS for English. The tectogrammatical layer is being adapted for English (Cinková et al., 2004; Hajič et al., 2009) and we are able to use the available tools to obtain all SemPOS features for English sentences as well.

3.2 Evaluation of SemPOS and Friends

We measured the metric performance on data used in MetricsMATR08, WMT09 and WMT08. For the evaluation of metric correlation with human judgments at the system level, we used the Pearson correlation coefficient ρ applied to ranks. In case of a tie, the systems were assigned the average position. For example if three systems achieved the same highest score (thus occupying the positions 1, 2 and 3 when sorted by score), each of them would obtain the average rank of $2 = \frac{1+2+3}{3}$. When correlating ranks (instead of exact scores) and with this handling of ties, the Pearson coefficient is equivalent to Spearman’s rank correlation coefficient.

The MetricsMATR08 human judgments include preferences for pairs of MT systems saying which one of the two systems is better, while the WMT08 and WMT09 data contain system scores (for up to 5 systems) on the scale 1 to 5 for a given sentence. We assigned a human ranking to the systems based on the percent of time that their translations were judged to be better than or equal to the translations of any other system in the manual evaluation. We converted automatic metric scores to ranks.

Metrics’ performance for translation to English and Czech was measured on the following testsets (the number of human judgments for a given source language in brackets):

To English: MetricsMATR08 (cn+ar: 1652), WMT08 News Articles (de: 199, fr: 251), WMT08 Europarl (es: 190, fr: 183), WMT09 (cz: 320, de: 749, es: 484, fr: 786, hu: 287)

To Czech: WMT08 News Articles (en: 267), WMT08 Commentary (en: 243), WMT09 (en: 1425)

The MetricsMATR08 testset contained 4 reference translations for each sentence whereas the remaining testsets only one reference.

Correlation coefficients for English are shown in Table 2. The best metric is Void_{par} closely followed by Void_{sons}. The explanation is that Void compared to SemPOS or Functor does not lose points by an erroneous assignment of the POS or the functor, and that Void_{par} profits from checking the dependency relations between autosemantic words. The combination of BLEU and SemPOS⁶ outperforms both individual metrics, but in case of SemPOS only by a minimal difference. Additionally, we confirm that 4-grams alone have little discriminative power both when used as a metric of their own (BLEU₄) as well as in a linear combination with SemPOS.

The best metric for Czech (see Table 3) is a linear combination of SemPOS and 4-gram BLEU closely followed by other SemPOS and BLEU_n combinations. We assume this is because BLEU₄ can capture correctly translated fixed phrases, which is positively reflected in human judgments. Including BLEU₁ in the combination favors translations with word forms as expected by the refer-

⁶For each $n \in \{1, 2, 3, 4\}$, we show only the best weight setting for SemPOS and BLEU_n.

Metric	Avg	Best	Worst
Void _{par}	0.75	0.89	0.60
Void _{sons}	0.75	0.90	0.54
Void	0.72	0.91	0.59
Functor _{sons}	0.72	1.00	0.43
GTM	0.71	0.90	0.54
4·SemPOS+1·BLEU ₂	0.70	0.93	0.43
SemPOS _{par}	0.70	0.93	0.30
1·SemPOS+4·BLEU ₃	0.70	0.91	0.26
4·SemPOS+1·BLEU ₁	0.69	0.93	0.43
NIST	0.69	0.90	0.53
SemPOS _{sons}	0.69	0.94	0.40
SemPOS	0.69	0.95	0.30
2·SemPOS+1·BLEU ₄	0.68	0.91	0.09
BLEU ₁	0.68	0.87	0.43
BLEU ₂	0.68	0.90	0.26
BLEU ₃	0.66	0.90	0.14
BLEU	0.66	0.91	0.20
TER	0.63	0.87	0.29
PER	0.63	0.88	0.32
BLEU ₄	0.61	0.90	-0.31
Functor _{par}	0.57	0.83	-0.03
Functor	0.55	0.82	-0.09

Table 2: Average, best and worst system-level correlation coefficients for translation to English from various source languages evaluated on 10 different testsets.

ence, thus allowing to spot bad word forms. In all cases, the linear combination puts more weight on SemPOS. Given the negligible difference between SemPOS alone and the linear combinations, we see that word forms are not the major issue for humans interpreting the translation—most likely because the systems so far often make more important errors. This is also confirmed by the observation that using BLEU alone is rather unreliable for Czech and BLEU-1 (which judges unigrams only) is even worse. Surprisingly BLEU-2 performed better than any other n -grams for reasons that have yet to be examined. The error metrics PER and TER showed the lowest correlation with human judgments for translation to Czech.

4 Conclusion

This paper documented problems of single-reference BLEU when applied to morphologically rich languages such as Czech. BLEU suffers from a sparse data problem, unable to judge the quality of tokens not confirmed by the reference. This is confirmed for other languages as well: the lower the BLEU score the lower the correlation to human judgments.

We introduced a refinement of SemPOS, an automatic metric of MT quality based on deep-syntactic representation of the sentence tackling

Metric	Avg	Best	Worst
3·SemPOS+1·BLEU ₄	0.55	0.83	0.14
2·SemPOS+1·BLEU ₂	0.55	0.83	0.14
2·SemPOS+1·BLEU ₁	0.53	0.83	0.09
4·SemPOS+1·BLEU ₃	0.53	0.83	0.09
SemPOS	0.53	0.83	0.09
BLEU ₂	0.43	0.83	0.09
SemPOS _{par}	0.37	0.53	0.14
Functor _{sons}	0.36	0.53	0.14
GTM	0.35	0.53	0.14
BLEU ₄	0.33	0.53	0.09
Void	0.33	0.53	0.09
NIST	0.33	0.53	0.09
Void _{sons}	0.33	0.53	0.09
BLEU	0.33	0.53	0.09
BLEU ₃	0.33	0.53	0.09
BLEU ₁	0.29	0.53	-0.03
SemPOS _{sons}	0.28	0.42	0.03
Functor _{par}	0.23	0.40	0.14
Functor	0.21	0.40	0.09
Void _{par}	0.16	0.53	-0.08
PER	0.12	0.53	-0.09
TER	0.07	0.53	-0.23

Table 3: System-level correlation coefficients for English-to-Czech translation evaluated on 3 different testsets.

the sparse data issue. SemPOS was evaluated on translation to Czech and to English, scoring better than or comparable to many established metrics.

References

- Ondřej Bojar, David Mareček, Václav Novák, Martin Popel, Jan Ptáček, Jan Rouš, and Zdeněk Žabokrtský. 2009. English-Czech MT in 2008. In *Proceedings of the Fourth Workshop on Statistical Machine Translation*, Athens, Greece, March. Association for Computational Linguistics.
- Chris Callison-Burch, Cameron Fordyce, Philipp Koehn, Christof Monz, and Josh Schroeder. 2008. Further meta-evaluation of machine translation. In *Proceedings of the Third Workshop on Statistical Machine Translation*, pages 70–106, Columbus, Ohio, June. Association for Computational Linguistics.
- Chris Callison-Burch, Philipp Koehn, Christof Monz, and Josh Schroeder. 2009. Findings of the 2009 workshop on statistical machine translation. In *Proceedings of the Fourth Workshop on Statistical Machine Translation*, Athens, Greece. Association for Computational Linguistics.
- Silvie Cinková, Jan Hajič, Marie Mikulová, Lucie Mladová, Anja Nedolužko, Petr Pajas, Jarmila Panevová, Jiří Semecký, Jana Šindlerová, Josef Toman, Zdeňka Urešová, and Zdeněk Žabokrtský. 2004. Annotation of English on the tectogrammatical level. Technical Report TR-2006-35, ÚFAL/CKL, Prague, Czech Republic, December.

- Sherri Condon, Gregory A. Sanders, Dan Parvaz, Alan Rubenstein, Christy Doran, John Aberdeen, and Beatrice Oshika. 2009. Normalization for Automated Metrics: English and Arabic Speech Translation. In *MT Summit XII*.
- Jesús Giménez and Lluís Márquez. 2007. Linguistic Features for Automatic Evaluation of Heterogeneous MT Systems. In *Proceedings of the Second Workshop on Statistical Machine Translation*, pages 256–264, Prague, June. Association for Computational Linguistics.
- Jan Hajič, Silvie Cinková, Kristýna Čermáková, Lucie Mladová, Anja Nedolužko, Petr Pajas, Jiří Sebecký, Jana Šindlerová, Josef Toman, Kristýna Tomšů, Matěj Korvas, Magdaléna Rysová, Kateřina Veselovská, and Zdeněk Žabokrtský. 2009. Prague English Dependency Treebank 1.0. Institute of Formal and Applied Linguistics, Charles University in Prague, ISBN 978-80-904175-0-2, January.
- Jan Hajič, Jarmila Panevová, Eva Hajičová, Petr Sgall, Petr Pajas, Jan Štěpánek, Jiří Havelka, Marie Mikulová, Zdeněk Žabokrtský, and Magda Ševčíková Razímová. 2006. Prague Dependency Treebank 2.0. LDC2006T01, ISBN: 1-58563-370-4.
- Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondřej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open Source Toolkit for Statistical Machine Translation. In *ACL 2007, Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions*, pages 177–180, Prague, Czech Republic, June. Association for Computational Linguistics.
- Kamil Kos and Ondřej Bojar. 2009. Evaluation of Machine Translation Metrics for Czech as the Target Language. *Prague Bulletin of Mathematical Linguistics*, 92.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a Method for Automatic Evaluation of Machine Translation. In *ACL 2002, Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania.
- M. Przybocki, K. Peterson, and S. Bronsart. 2008. Official results of the NIST 2008 "Metrics for Machine TRanslation" Challenge (MetricsMATR08).
- Petr Sgall, Eva Hajičová, and Jarmila Panevová. 1986. *The Meaning of the Sentence and Its Semantic and Pragmatic Aspects*. Academia/Reidel Publishing Company, Prague, Czech Republic/Dordrecht, Netherlands.
- Zdeněk Žabokrtský and Ondřej Bojar. 2008. TectoMT, Developer's Guide. Technical Report TR-2008-39, Institute of Formal and Applied Linguistics, Faculty of Mathematics and Physics, Charles University in Prague, December.

Exemplar-Based Models for Word Meaning In Context

Katrin Erk

Department of Linguistics
University of Texas at Austin
katrin.erk@mail.utexas.edu

Sebastian Padó

Institut für maschinelle Sprachverarbeitung
Stuttgart University
pado@ims.uni-stuttgart.de

Abstract

This paper describes ongoing work on distributional models for word meaning in context. We abandon the usual one-vector-per-word paradigm in favor of an exemplar model that activates only relevant occurrences. On a paraphrasing task, we find that a simple exemplar model outperforms more complex state-of-the-art models.

1 Introduction

Distributional models are a popular framework for representing word meaning. They describe a lemma through a high-dimensional vector that records co-occurrence with context features over a large corpus. Distributional models have been used in many NLP analysis tasks (Salton et al., 1975; McCarthy and Carroll, 2003; Salton et al., 1975), as well as for cognitive modeling (Baroni and Lenci, 2009; Landauer and Dumais, 1997; McDonald and Ramscar, 2001). Among their attractive properties are their simplicity and versatility, as well as the fact that they can be acquired from corpora in an unsupervised manner.

Distributional models are also attractive as a model of word meaning in context, since they do not have to rely on fixed sets of dictionary sense with their well-known problems (Kilgarriff, 1997; McCarthy and Navigli, 2009). Also, they can be used directly for testing paraphrase applicability (Szpektor et al., 2008), a task that has recently become prominent in the context of textual entailment (Bar-Haim et al., 2007). However, polysemy is a fundamental problem for distributional models. Typically, distributional models compute a single “type” vector for a target word, which contains co-occurrence counts for all the occurrences of the target in a large corpus. If the target is polysemous, this vector mixes contextual features for all the senses of the target. For example, among the

top 20 features for *coach*, we get *match* and *team* (for the “trainer” sense) as well as *driver* and *car* (for the “bus” sense). This problem has typically been approached by modifying the type vector for a target to better match a given context (Mitchell and Lapata, 2008; Erk and Padó, 2008; Thater et al., 2009).

In the terms of research on human concept representation, which often employs feature vector representations, the use of type vectors can be understood as a *prototype*-based approach, which uses a single vector per category. From this angle, computing prototypes throws away much interesting distributional information. A rival class of models is that of *exemplar* models, which memorize each seen instance of a category and perform categorization by comparing a new stimulus to each remembered exemplar vector.

We can address the polysemy issue through an exemplar model by simply removing all exemplars that are “not relevant” for the present context, or conversely *activating* only the relevant ones. For the *coach* example, in the context of a text about motorways, presumably an instance like “The coach drove a steady 45 mph” would be activated, while “The team lost all games since the new coach arrived” would not.

In this paper, we present an exemplar-based distributional model for modeling word meaning in context, applying the model to the task of deciding paraphrase applicability. With a very simple vector representation and just using activation, we outperform the state-of-the-art prototype models. We perform an in-depth error analysis to identify stable parameters for this class of models.

2 Related Work

Among distributional models of word, there are some approaches that address polysemy, either by inducing a fixed clustering of contexts into senses (Schütze, 1998) or by dynamically modi-

fying a word’s type vector according to each given sentence context (Landauer and Dumais, 1997; Mitchell and Lapata, 2008; Erk and Padó, 2008; Thater et al., 2009). Polysemy-aware approaches also differ in their notion of *context*. Some use a bag-of-words representation of words in the current sentence (Schütze, 1998; Landauer and Dumais, 1997), some make use of syntactic context (Mitchell and Lapata, 2008; Erk and Padó, 2008; Thater et al., 2009). The approach that we present in the current paper computes a representation dynamically for each sentence context, using a simple bag-of-words representation of context.

In cognitive science, prototype models predict degree of category membership through similarity to a single prototype, while exemplar theory represents a concept as a collection of all previously seen exemplars (Murphy, 2002). Griffiths et al. (2007) found that the benefit of exemplars over prototypes grows with the number of available exemplars. The problem of representing meaning in context, which we consider in this paper, is closely related to the problem of *concept combination* in cognitive science, i.e., the derivation of representations for complex concepts (such as “metal spoon”) given the representations of base concepts (“metal” and “spoon”). While most approaches to concept combination are based on prototype models, Voorspoels et al. (2009) show superior results for an exemplar model based on exemplar *activation*.

In NLP, exemplar-based (memory-based) models have been applied to many problems (Daelemans et al., 1999). In the current paper, we use an exemplar model for computing distributional representations for word meaning in context, using the context to *activate* relevant exemplars. Comparing representations of context, bag-of-words (BOW) representations are more informative and noisier, while syntax-based representations deliver sparser and less noisy information. Following the hypothesis that richer, topical information is more suitable for exemplar activation, we use BOW representations of sentential context in the current paper.

3 Exemplar Activation Models

We now present an exemplar-based model for meaning in context. It assumes that each target lemma is represented by a set of exemplars, where an exemplar is a sentence in which the target occurs, represented as a vector. We use lowercase letters for individual exemplars (vectors), and uppercase

Sentential context	Paraphrase
After a fire extinguisher is used, it must always be returned for recharging and its use recorded.	bring back (3), take back (2), send back (1), give back (1)
We return to the young woman who is reading the Wrigley’s wrapping paper.	come back (3), revert (1), revisit (1), go (1)

Table 1: The Lexical Substitution (LexSub) dataset.

letters for sets of exemplars.

We model polysemy by *activating* relevant exemplars of a lemma E in a given sentence context s . (Note that we use E to refer to both a lemma and its exemplar set, and that s can be viewed as just another exemplar vector.) In general, we define *activation* of a set E by exemplar s as

$$act(E, s) = \{e \in E \mid sim(e, s) > \theta(E, s)\}$$

where E is an exemplar set, s is the “point of comparison”, sim is some similarity measure such as Cosine or Jaccard, and $\theta(E, s)$ is a threshold. Exemplars belong to the activated set if their similarity to s exceeds $\theta(E, s)$.¹ We explore two variants of activation. In **k NN activation**, the k most similar exemplars to s are activated by setting θ to the similarity of the k -th most similar exemplar. In **q -percentage activation**, we activate the top $q\%$ of E by setting θ to the $(100-q)$ -th percentile of the $sim(e, s)$ distribution. Note that, while in the kNN activation scheme the number of activated exemplars is the same for every lemma, this is not the case for percentage activation: There, a more frequent lemma (i.e., a lemma with more exemplars) will have more exemplars activated.

Exemplar activation for paraphrasing. A paraphrase is typically only applicable to a particular sense of a target word. Table 1 illustrates this on two examples from the Lexical Substitution (LexSub) dataset (McCarthy and Navigli, 2009), both featuring the target *return*. The right column lists appropriate paraphrases of *return* in each context (given by human annotators).² We apply the exemplar activation model to the task of predicting paraphrase felicity: Given a target lemma T in a particular sentential context s , and given a list of

¹In principle, activation could be treated not just as binary inclusion/exclusion, but also as a graded weighting scheme. However, weighting schemes introduce a large number of parameters, which we wanted to avoid.

²Each annotator was allowed to give up to three paraphrases per target in context. As a consequence, the number of gold paraphrases per target sentence varies.

potential paraphrases of T , the task is to predict which of the paraphrases are applicable in s .

Previous approaches (Mitchell and Lapata, 2008; Erk and Padó, 2008; Erk and Padó, 2009; Thater et al., 2009) have performed this task by modifying the type vector for T to the context s and then comparing the resulting vector T' to the type vector of a paraphrase candidate P . In our exemplar setting, we select a contextually adequate subset of contexts in which T has been observed, using $T' = act(T, s)$ as a generalized representation of meaning of target T in the context of s .

Previous approaches used all of P as a representation for a paraphrase candidate P . However, P includes also irrelevant exemplars, while for a paraphrase to be judged as good, it is sufficient that one plausible reading exists. Therefore, we use $P' = act(P, s)$ to represent the paraphrase.

4 Experimental Evaluation

Data. We evaluate our model on predicting paraphrases from the Lexical Substitution (LexSub) dataset (McCarthy and Navigli, 2009). This dataset consists of 2000 instances of 200 target words in sentential contexts, with paraphrases for each target word instance generated by up to 6 participants. Paraphrases are ranked by the number of annotators that chose them (cf. Table 1). Following Erk and Padó (2008), we take the list of paraphrase candidates for a target as given (computed by pooling all paraphrases that LexSub annotators proposed for the target) and use the models to rank them for any given sentence context.

As exemplars, we create bag-of-words co-occurrence vectors from the BNC. These vectors represent instances of a target word by the other words in the same sentence, lemmatized and POS-tagged, minus stop words. E.g., if the lemma *gnurge* occurs twice in the BNC, once in the sentence “The dog will gnurge the other dog”, and once in “The old windows gnurged”, the exemplar set for *gnurge* contains the vectors [*dog-n: 2, other-a:1*] and [*old-a: 1, window-n: 1*]. For exemplar similarity, we use the standard Cosine similarity, and for the similarity of two exemplar sets, the Cosine of their centroids.

Evaluation. The model’s prediction for an item is a list of paraphrases ranked by their predicted goodness of fit. To evaluate them against a weighted list of gold paraphrases, we follow Thater et al. (2009) in using Generalized Average Preci-

parameter	<i>actT</i>		<i>actP</i>	
	kNN	perc.	kNN	perc.
10	36.1	35.5	36.5	38.6
20	36.2	35.2	36.2	37.9
30	36.1	35.3	35.8	37.8
40	36.0	35.3	35.8	37.7
50	35.9	35.1	35.9	37.5
60	36.0	35.0	36.1	37.5
70	35.9	34.8	36.1	37.5
80	36.0	34.7	36.0	37.4
90	35.9	34.5	35.9	37.3
no act.	34.6		35.7	
random BL	28.5			

Table 2: Activation of T or P individually on the full LexSub dataset (GAP evaluation)

sion (GAP), which interpolates the precision values of top- n prediction lists for increasing n . Let $G = \langle q_1, \dots, q_m \rangle$ be the list of gold paraphrases with gold weights $\langle y_1, \dots, y_m \rangle$. Let $P = \langle p_1, \dots, p_n \rangle$ be the list of model predictions as ranked by the model, and let $\langle x_1, \dots, x_n \rangle$ be the *gold* weights associated with them (assume $x_i = 0$ if $p_i \notin G$), where $G \subseteq P$. Let $I(x_i) = 1$ if $p_i \in G$, and zero otherwise. We write $\bar{x}_i = \frac{1}{i} \sum_{k=1}^i x_k$ for the average gold weight of the first i model predictions, and analogously \bar{y}_i . Then

$$GAP(P, G) = \frac{1}{\sum_{j=1}^m I(y_j) \bar{y}_j} \sum_{i=1}^n I(x_i) \bar{x}_i$$

Since the model may rank multiple paraphrases the same, we average over 10 random permutations of equally ranked paraphrases. We report mean GAP over all items in the dataset.

Results and Discussion. We first computed two models that activate either the paraphrase or the target, but not both. Model 1, *actT*, activates only the target, using the complete P as paraphrase, and ranking paraphrases by $sim(P, act(T, s))$. Model 2, *actP*, activates only the paraphrase, using s as the target word, ranking by $sim(act(P, s), s)$.

The results for these models are shown in Table 2, with both kNN and percentage activation: kNN activation with a parameter of 10 means that the 10 closest neighbors were activated, while percentage with a parameter of 10 means that the closest 10% of the exemplars were used. Note first that we computed a random baseline (last row) with a GAP of 28.5. The second-to-last row (“no activation”) shows two more informed baselines.

The *actT* “no act” result (34.6) corresponds to a prototype-based model that ranks paraphrase candidates by the distance between their type vectors and the target’s type vector. Virtually all exemplar models outperform this prototype model. Note also that both *actT* and *actP* show the best results for small values of the activation parameter. This indicates paraphrases can be judged on the basis of a rather small number of exemplars. Nevertheless, *actT* and *actP* differ with regard to the details of their optimal activation. For *actT*, a small absolute number of activated exemplars (here, 20) works best, while *actP* yields the best results for a small percentage of paraphrase exemplars. This can be explained by the different functions played by *actT* and *actP* (cf. Section 3): Activation of the paraphrase must allow a guess about whether there is reasonable interpretation of *P* in the context *s*. This appears to require a reasonably-sized sample from *P*. In contrast, target activation merely has to counteract the sparsity of *s*, and activation of too many exemplars from *T* leads to oversmoothing.

We obtained significances by computing 95% and 99% confidence intervals with bootstrap resampling. As a rule of thumb, we find that 0.4% difference in GAP corresponds to a significant difference at the 95% level, and 0.7% difference in GAP to significance at the 99% level. The four activation methods (i.e., columns in Table 2) are significantly different from each other, with the exception of the pair *actT*/kNN and *actP*/kNN (n.s.), so that we get the following order:

$$actP/perc > actP/kNN \approx actT/kNN > actT/perc$$

where $>$ means “significantly outperforms”. In particular, the best method (*actT*/kNN) outperforms all other methods at $p < 0.01$. Here, the best parameter setting (10% activation) is also significantly better than the next-one one (20% activation). With the exception of *actT*/perc, all activation methods significantly outperform the best baseline (*actP*, no activation).

Based on these observations, we computed a third model, *actTP*, that activates both *T* (by kNN) and *P* (by percentage), ranking paraphrases by $sim(act(P, s), act(T, s))$. Table 3 shows the results. We find the overall best model at a similar location in parameter space as for *actT* and *actP* (cf. Table 2), namely by setting the activation parameters to small values. The sensitivity of the parameters changes considerably, though. When

<i>P</i> activation (%) \Rightarrow	10	20	30
<i>T</i> activation (kNN) \Downarrow			
5	38.2	38.1	38.1
10	37.6	37.8	37.7
20	37.3	37.4	37.3
40	37.2	37.2	36.1

Table 3: Joint activation of *P* and *T* on the full LexSub dataset (GAP evaluation)

we fix the *actP* activation level, we find comparatively large performance differences between the *T* activation settings $k=5$ and $k=10$ (highly significant for 10% *actP*, and significant for 20% and 30% *actP*). On the other hand, when we fix the *actT* activation level, changes in *actP* activation generally have an insignificant impact.

Somewhat disappointingly, we are not able to surpass the best result for *actP* alone. This indicates that – at least in the current vector space – the sparsity of *s* is less of a problem than the “dilution” of *s* that we face when we representing the target word by exemplars of *T* close to *s*. Note, however, that the numerically worse performance of the best *actTP* model is still not significantly different from the best *actP* model.

Influence of POS and frequency. An analysis of the results by target part-of-speech showed that the globally optimal parameters also yield the best results for individual POS, even though there are substantial differences among POS. For *actT*, the best results emerge for all POS with kNN activation with k between 10 and 30. For $k=20$, we obtain a GAP of 35.3 (verbs), 38.2 (nouns), and 35.1 (adjectives). For *actP*, the best parameter for all POS was activation of 10%, with GAPs of 36.9 (verbs), 41.4 (nouns), and 37.5 (adjectives). Interestingly, the results for *actTP* (verbs: 38.4, nouns: 40.6, adjectives: 36.9) are better than *actP* for verbs, but worse for nouns and adjectives, which indicates that the sparsity problem might be more prominent than for the other POS. In all three models, we found a clear effect of target and paraphrase frequency, with deteriorating performance for the highest-frequency targets as well as for the lemmas with the highest average paraphrase frequency.

Comparison to other models. Many of the other models are syntax-based and are therefore only applicable to a subset of the LexSub data. We have re-evaluated our exemplar models on the subsets we used in Erk and Padó (2008, EP08, 367

	Models		
	EP08	EP09	TDP09
EP08 dataset	27.4	NA	NA
EP09 dataset	NA	32.2	36.5
	<i>actT</i>	<i>actP</i>	<i>actTP</i>
EP08 dataset	36.5	38.0	39.9
EP09 dataset	39.1	39.9	39.6

Table 4: Comparison to other models on two subsets of LexSub (GAP evaluation)

datapoints) and Erk and Padó (2009, EP09, 100 datapoints). The second set was also used by Thater et al. (2009, TDP09). The results in Table 4 compare these models against our best previous exemplar models and show that our models outperform these models across the board.³ Due to the small sizes of these datasets, statistical significance is more difficult to attain. On EP09, the differences among our models are not significant, but the difference between them and the original EP09 model is.⁴ On EP08, all differences are significant except for *actP* vs. *actTP*.

We note that both the EP08 and the EP09 datasets appear to be simpler to model than the complete Lexical Substitution dataset, at least by our exemplar-based models. This underscores an old insight: namely, that direct syntactic neighbors, such as arguments and modifiers, provide strong clues as to word sense.

5 Conclusions and Outlook

This paper reports on work in progress on an exemplar activation model as an alternative to one-vector-per-word approaches to word meaning in context. Exemplar activation is very effective in handling polysemy, even with a very simple (and sparse) bag-of-words vector representation. On both the EP08 and EP09 datasets, our models surpass more complex prototype-based approaches (Tab. 4). It is also noteworthy that the exemplar activation models work best when few exemplars are used, which bodes well for their efficiency.

We found that the best target representations re-

³Since our models had the advantage of being tuned on the dataset, we also report the range of results across the parameters we tested. On the EP08 dataset, we obtained 33.1–36.5 for *actT*; 33.3–38.0 for *actP*; 37.7–39.9 for *actTP*. On the EP09 dataset, the numbers were 35.8–39.1 for *actT*; 38.1–39.9 for *actP*; 37.2–39.8 for *actTP*.

⁴We did not have access to the TDP09 predictions to do significance testing.

sult from activating a low absolute number of exemplars. Paraphrase representations are best activated with a percentage-based threshold. Overall, we found that paraphrase activation had a much larger impact on performance than target activation, and that drawing on target exemplars other than *s* to represent the target meaning in context improved over using *s* itself only for verbs (Tab. 3). This suggests the possibility of considering *T*’s activated paraphrase candidates as the representation of *T* in the context *s*, rather than some vector of *T* itself, in the spirit of Kintsch (2001).

While it is encouraging that the best parameter settings involved the activation of only few exemplars, computation with exemplar models still requires the management of large numbers of vectors. The computational overhead can be reduced by using data structures that cut down on the number of vector comparisons, or by decreasing vector dimensionality (Gorman and Curran, 2006). We will experiment with those methods to determine the tradeoff of runtime and accuracy for this task.

Another area of future work is to move beyond bag-of-words context: It is known from WSD that syntactic and bag-of-words contexts provide complementary information (Florian et al., 2002; Szpektor et al., 2008), and we hope that they can be integrated in a more sophisticated exemplar model.

Finally, we will to explore task-based evaluations. Relation extraction and textual entailment in particular are tasks where similar models have been used before (Szpektor et al., 2008).

Acknowledgements. This work was supported in part by National Science Foundation grant IIS-0845925, and by a Morris Memorial Grant from the New York Community Trust.

References

- R. Bar-Haim, I. Dagan, I. Greental, and E. Shnarch. 2007. Semantic inference at the lexical-syntactic level. In *Proceedings of AAAI*, pages 871–876, Vancouver, BC.
- M. Baroni and A. Lenci. 2009. One distributional memory, many semantic spaces. In *Proceedings of the EACL Workshop on Geometrical Models of Natural Language Semantics*, Athens, Greece.
- W. Daelemans, A. van den Bosch, and J. Zavrel. 1999. Forgetting exceptions is harmful in language learning. *Machine Learning*, 34(1/3):11–43. Special Issue on Natural Language Learning.
- K. Erk and S. Padó. 2008. A structured vector space

- model for word meaning in context. In *Proceedings of EMNLP*, pages 897–906, Honolulu, HI.
- K. Erk and S. Padó. 2009. Paraphrase assessment in structured vector space: Exploring parameters and datasets. In *Proceedings of the EACL Workshop on Geometrical Models of Natural Language Semantics*, Athens, Greece.
- R. Florian, S. Cucerzan, C. Schafer, and D. Yarowsky. 2002. Combining classifiers for word sense disambiguation. *Journal of Natural Language Engineering*, 8(4):327–341.
- J. Gorman and J. R. Curran. 2006. Scaling distributional similarity to large corpora. In *Proceedings of ACL*, pages 361–368, Sydney.
- T. Griffiths, K. Canini, A. Sanborn, and D. J. Navarro. 2007. Unifying rational models of categorization via the hierarchical Dirichlet process. In *Proceedings of CogSci*, pages 323–328, Nashville, TN.
- A. Kilgarriff. 1997. I don’t believe in word senses. *Computers and the Humanities*, 31(2):91–113.
- W. Kintsch. 2001. Predication. *Cognitive Science*, 25:173–202.
- T. Landauer and S. Dumais. 1997. A solution to Platos problem: the latent semantic analysis theory of acquisition, induction, and representation of knowledge. *Psychological Review*, 104(2):211–240.
- D. McCarthy and J. Carroll. 2003. Disambiguating nouns, verbs, and adjectives using automatically acquired selectional preferences. *Computational Linguistics*, 29(4):639–654.
- D. McCarthy and R. Navigli. 2009. The English lexical substitution task. *Language Resources and Evaluation*, 43(2):139–159. Special Issue on Computational Semantic Analysis of Language: SemEval-2007 and Beyond.
- S. McDonald and M. Ramscar. 2001. Testing the distributional hypothesis: The influence of context on judgements of semantic similarity. In *Proceedings of CogSci*, pages 611–616.
- J. Mitchell and M. Lapata. 2008. Vector-based models of semantic composition. In *Proceedings of ACL*, pages 236–244, Columbus, OH.
- G. L. Murphy. 2002. *The Big Book of Concepts*. MIT Press.
- G. Salton, A. Wang, and C. Yang. 1975. A vector-space model for information retrieval. *Journal of the American Society for Information Science*, 18:613–620.
- H. Schütze. 1998. Automatic word sense discrimination. *Computational Linguistics*, 24(1):97–124.
- I. Szpektor, I. Dagan, R. Bar-Haim, and J. Goldberger. 2008. Contextual preferences. In *Proceedings of ACL*, pages 683–691, Columbus, OH.
- S. Thater, G. Dinu, and M. Pinkal. 2009. Ranking paraphrases in context. In *Proceedings of the ACL Workshop on Applied Textual Inference*, pages 44–47, Singapore.
- W. Voorspoels, W. Vanpaemel, and G. Storms. 2009. The role of extensional information in conceptual combination. In *Proceedings of CogSci*.

A Structured Model for Joint Learning of Argument Roles and Predicate Senses

Yotaro Watanabe

Graduate School of Information Sciences
Tohoku University
6-6-05, Aramaki Aza Aoba, Aoba-ku,
Sendai 980-8579, Japan
yotaro-w@ecei.tohoku.ac.jp

Masayuki Asahara Yuji Matsumoto

Graduate School of Information Science
Nara Institute of Science and Technology
8916-5 Takayama, Ikoma,
Nara, 630-0192, Japan
{masayu-a, matsu}@is.naist.jp

Abstract

In predicate-argument structure analysis, it is important to capture non-local dependencies among arguments and inter-dependencies between the sense of a predicate and the semantic roles of its arguments. However, no existing approach explicitly handles both non-local dependencies and semantic dependencies between predicates and arguments. In this paper we propose a structured model that overcomes the limitation of existing approaches; the model captures both types of dependencies simultaneously by introducing four types of factors including a global factor type capturing non-local dependencies among arguments and a pairwise factor type capturing local dependencies between a predicate and an argument. In experiments the proposed model achieved competitive results compared to the state-of-the-art systems without applying any feature selection procedure.

1 Introduction

Predicate-argument structure analysis is a process of assigning *who* does *what* to *whom*, *where*, *when*, etc. for each predicate. Arguments of a predicate are assigned particular *semantic roles*, such as *Agent*, *Theme*, *Patient*, etc. Lately, predicate-argument structure analysis has been regarded as a task of assigning semantic roles of arguments as well as word senses of a predicate (Surdeanu et al., 2008; Hajič et al., 2009).

Several researchers have paid much attention to predicate-argument structure analysis, and the following two important factors have been shown. Toutanova et al. (2008), Johansson and Nugues (2008), and Björkelund et al. (2009) presented importance of capturing non-local dependencies

of core arguments in predicate-argument structure analysis. They used argument sequences tied with a predicate sense (e.g. AGENT-buy.01/Active-PATIENT) as a feature for the re-ranker of the system where predicate sense and argument role candidates are generated by their pipelined architecture. They reported that incorporating this type of features provides substantial gain of the system performance.

The other factor is inter-dependencies between a predicate sense and argument roles, which relate to selectional preference, and motivated us to jointly identify a predicate sense and its argument roles. This type of dependencies has been explored by Riedel and Meza-Ruiz (2008; 2009b; 2009a), all of which use Markov Logic Networks (MLN). The work uses the global formulae that have atoms in terms of both a predicate sense and each of its argument roles, and the system identifies predicate senses and argument roles simultaneously.

Ideally, we want to capture both types of dependencies simultaneously. The former approaches can not explicitly include features that capture inter-dependencies between a predicate sense and its argument roles. Though these are implicitly incorporated by re-ranking where the most plausible assignment is selected from a small subset of predicate and argument candidates, which are generated independently. On the other hand, it is difficult to deal with core argument features in MLN. Because the number of core arguments varies with the role assignments, this type of features cannot be expressed by a single formula.

Thompson et al. (2010) proposed a generative model that captures both predicate senses and its argument roles. However, the first-order markov assumption of the model eliminates ability to capture non-local dependencies among arguments. Also, generative models are in general inferior to discriminatively trained linear or log-

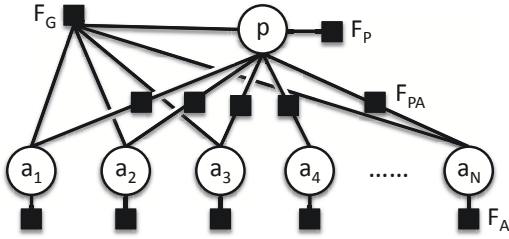


Figure 1: Undirected graphical model representation of the structured model

linear models.

In this paper we propose a structured model that overcomes limitations of the previous approaches. For the model, we introduce several types of features including those that capture both non-local dependencies of core arguments, and inter-dependencies between a predicate sense and its argument roles. By doing this, both tasks are mutually influenced, and the model determines the most plausible set of assignments of a predicate sense and its argument roles simultaneously. We present an exact inference algorithm for the model, and a large-margin learning algorithm that can handle both local and global features.

2 Model

Figure 1 shows the graphical representation of our proposed model. The node p corresponds to a predicate, and the nodes a_1, \dots, a_N to arguments of the predicate. Each node is assigned a particular predicate sense or an argument role label. The black squares are *factors* which provide scores of label assignments. In the model, the nodes for arguments depend on the predicate sense, and by influencing labels of a predicate sense and its argument roles, the most plausible label assignment of the nodes is determined considering all factors.

In this work, we use linear models. Let \mathbf{x} be words in a sentence, p be a sense of a predicate in \mathbf{x} , and $\mathcal{A} = \{a_n\}_1^N$ be a set of possible role label assignments for \mathbf{x} . A predicate-argument structure is represented by a pair of p and \mathcal{A} . We define the score function for predicate-argument structures as $s(p, \mathcal{A}) = \sum_{F_k \in \mathcal{F}} F_k(\mathbf{x}, p, \mathcal{A})$. \mathcal{F} is a set of all the factors, $F_k(\mathbf{x}, p, \mathcal{A})$ corresponds to a particular factor in Figure 1, and gives a score to a predicate or argument label assignments. Since we use linear models, $F_k(\mathbf{x}, p, \mathcal{A}) = \mathbf{w} \cdot \Phi_k(\mathbf{x}, p, \mathcal{A})$.

2.1 Factors of the Model

We define four types of factors for the model.

Predicate Factor F_P scores a sense of p , and does not depend on any arguments. The score function is defined by $F_P(\mathbf{x}, p, \mathcal{A}) = \mathbf{w} \cdot \Phi_P(\mathbf{x}, p)$.

Argument Factor F_A scores a label assignment of a particular argument $a \in \mathcal{A}$. The score is determined independently from a predicate sense, and is given by $F_A(\mathbf{x}, p, a) = \mathbf{w} \cdot \Phi_A(\mathbf{x}, a)$.

Predicate-Argument Pairwise Factor

F_{PA} captures inter-dependencies between a predicate sense and one of its argument roles. The score function is defined as $F_{PA}(\mathbf{x}, p, a) = \mathbf{w} \cdot \Phi_{PA}(\mathbf{x}, p, a)$. The difference from F_A is that F_{PA} influences both the predicate sense and the argument role. By introducing this factor, the role label can be influenced by the predicate sense, and vice versa.

Global Factor F_G is introduced to capture plausibility of the whole predicate-argument structure. Like the other factors, the score function is defined as $F_G(\mathbf{x}, p, \mathcal{A}) = \mathbf{w} \cdot \Phi_G(\mathbf{x}, p, \mathcal{A})$. A possible feature that can be considered by this factor is the mutual dependencies among core arguments. For instance, if a predicate-argument structure has an agent (A0) followed by the predicate and a patient (A1), we encode the structure as a string *A0-PRED-A1* and use it as a feature. This type of features provide *plausibility* of predicate-argument structures. Even if the highest scoring predicate-argument structure with the other factors misses some core arguments, the global feature demands the model to fill the missing arguments.

The numbers of factors for each factor type are: F_P and F_G are 1, F_A and F_{PA} are $|\mathcal{A}|$. By integrating the all factors, the score function becomes $s(p, \mathcal{A}) = \mathbf{w} \cdot \Phi_P(\mathbf{x}, p) + \mathbf{w} \cdot \Phi_G(\mathbf{x}, p, \mathcal{A}) + \mathbf{w} \cdot \sum_{a \in \mathcal{A}} \{\Phi_A(\mathbf{x}, a) + \Phi_{PA}(\mathbf{x}, p, a)\}$.

2.2 Inference

The crucial point of the model is how to deal with the global factor F_G , because enumerating possible assignments is too costly. A number of methods have been proposed for the use of global features for linear models such as (Daumé III and Marcu, 2005; Kazama and Torisawa, 2007). In this work, we use the approach proposed in (Kazama and Torisawa, 2007). Although the approach is proposed for sequence labeling tasks, it

can be easily extended to our structured model. That is, for each possible predicate sense p of the predicate, we provide N-best argument role assignments using three local factors F_P , F_A and F_{PA} , and then add scores of the global factor F_G , finally select the argmax from them. In this case, the argmax is selected from $|\mathcal{P}_l|N$ candidates.

2.3 Learning the Model

For learning of the model, we borrow a fundamental idea of Kazama and Torisawa’s perceptron learning algorithm. However, we use a more sophisticated online-learning algorithm based on the Passive-Aggressive Algorithm (PA) (Crammer et al., 2006).

For the sake of simplicity, we introduce some notations. We denote a predicate-argument structure $\mathbf{y} = \langle p, \mathcal{A} \rangle$, a local feature vector as $\Phi_L(\mathbf{x}, \mathbf{y}) = \Phi_P(\mathbf{x}, p) + \sum_{a \in \mathcal{A}} \{ \Phi_A(\mathbf{x}, a) + \Phi_{PA}(\mathbf{x}, p, a) \}$, a feature vector coupling both local and global features as $\Phi_{L+G}(\mathbf{x}, \mathbf{y}) = \Phi_L(\mathbf{x}, \mathbf{y}) + \Phi_G(\mathbf{x}, p, \mathcal{A})$, the argmax using Φ_{L+G} as $\hat{\mathbf{y}}^{L+G}$, the argmax using Φ_L as $\hat{\mathbf{y}}^L$. Also, we use a loss function $\rho(\mathbf{y}, \mathbf{y}')$, which is a cost function associated with \mathbf{y} and \mathbf{y}' .

The margin perceptron learning proposed by Kazama and Torisawa can be seen as an optimization with the following two constrains.

$$(A) \quad \mathbf{w} \cdot \Phi_{L+G}(\mathbf{x}, \mathbf{y}) - \mathbf{w} \cdot \Phi_{L+G}(\mathbf{x}, \hat{\mathbf{y}}^{L+G}) \geq \rho(\mathbf{y}, \hat{\mathbf{y}}^{L+G})$$

$$(B) \quad \mathbf{w} \cdot \Phi_L(\mathbf{x}, \mathbf{y}) - \mathbf{w} \cdot \Phi_L(\mathbf{x}, \hat{\mathbf{y}}^L) \geq \rho(\mathbf{y}, \hat{\mathbf{y}}^L)$$

(A) is the constraint that ensures a sufficient margin $\rho(\mathbf{y}, \hat{\mathbf{y}}^{L+G})$ between \mathbf{y} and $\hat{\mathbf{y}}^{L+G}$. (B) is the constraint that ensures a sufficient margin $\rho(\mathbf{y}, \hat{\mathbf{y}}^L)$ between \mathbf{y} and $\hat{\mathbf{y}}^L$. The necessity of this constraint is that if we apply only (A), the algorithm does not guarantee a sufficient margin in terms of local features, and it leads to poor quality in the N-best assignments. The Kazama and Torisawa’s perceptron algorithm uses constant values for the cost function $\rho(\mathbf{y}, \hat{\mathbf{y}}^{L+G})$ and $\rho(\mathbf{y}, \hat{\mathbf{y}}^L)$.

The proposed model is trained using the following optimization problem.

$$\begin{aligned} \mathbf{w}_{new} = \arg \min_{\mathbf{w}' \in \mathbb{R}^n} & \frac{1}{2} \|\mathbf{w}' - \mathbf{w}\|^2 + C\xi \\ \left\{ \begin{array}{ll} \text{s.t. } l_{L+G} \leq \xi, \xi \geq 0 & \text{if } \hat{\mathbf{y}}^{L+G} \neq \mathbf{y} \\ \text{s.t. } l_L \leq \xi, \xi \geq 0 & \text{if } \hat{\mathbf{y}}^{L+G} = \mathbf{y} \neq \hat{\mathbf{y}}^L \end{array} \right. & (1) \end{aligned}$$

$$\begin{aligned} l_{L+G} = \mathbf{w} \cdot \Phi_{L+G}(\mathbf{x}, \hat{\mathbf{y}}^{L+G}) \\ - \mathbf{w} \cdot \Phi_{L+G}(\mathbf{x}, \mathbf{y}) + \rho(\mathbf{y}, \hat{\mathbf{y}}^{L+G}) \end{aligned} \quad (2)$$

$$l_L = \mathbf{w} \cdot \Phi_L(\mathbf{x}, \hat{\mathbf{y}}^L) - \mathbf{w} \cdot \Phi_L(\mathbf{x}, \mathbf{y}) + \rho(\mathbf{y}, \hat{\mathbf{y}}^L) \quad (3)$$

l_{L+G} is the loss function for the case of using both local and global features, corresponding to the constraint (A), and l_L is the loss function for the case of using only local features, corresponding to the constraints (B) provided that (A) is satisfied.

2.4 The Role-less Argument Bias Problem

The fact that an argument candidate is not assigned any role (namely it is assigned the label “NONE”) is unlikely to contribute predicate sense disambiguation. However, it remains possible that “NONE” arguments is biased toward a particular predicate sense by F_{PA} (i.e. $\mathbf{w} \cdot \Phi_{PA}(\mathbf{x}, \text{sense}_i, a_k = \text{“NONE”}) > \mathbf{w} \cdot \Phi_{PA}(\mathbf{x}, \text{sense}_j, a_k = \text{“NONE”})$).

In order to avoid this bias, we define a special sense label, sense_{any} , that is used to calculate the score for a predicate and a roll-less argument, regardless of the predicate’s sense. We use the feature vector $\Phi_{PA}(\mathbf{x}, \text{sense}_{any}, a_k)$ if $a_k = \text{“NONE”}$ and $\Phi_{PA}(\mathbf{x}, \text{sense}_i, a_k)$ otherwise.

3 Experiment

3.1 Experimental Settings

We use the CoNLL-2009 Shared Task dataset (Hajič et al., 2009) for experiments. It is a dataset for multi-lingual syntactic and semantic dependency parsing¹. In the SRL-only challenge of the task, participants are required to identify predicate-argument structures of only the specified predicates. Therefore the problems to be solved are predicate sense disambiguation and argument role labeling. We use Semantic Labeled F1 for evaluation.

For generating N-bests, we used the beam-search algorithm, and the number of N-bests was set to $N = 64$. For learning of the joint model, the loss function $\rho(\mathbf{y}_t, \mathbf{y}')$ of the Passive-Aggressive Algorithm was set to the number of incorrect assignments of a predicate sense and its argument roles. Also, the number of iterations of the model used for testing was selected based on the performance on the development data.

Table 1 shows the features used for the structured model. The global features used for F_G are based on those used in (Toutanova et al., 2008; Johansson and Nugues, 2008), and the features

¹The dataset consists of seven languages: Catalan, Chinese, Czech, English, German, Japanese and Spanish.

F_P	Plemma of the predicate and predicate’s head, and ppos of the predicate Dependency label between the predicate and predicate’s head The concatenation of the dependency labels of the predicate’s dependents
F_A	Plemma and ppos of the predicate, the predicate’s head, the argument candidate, and the argument’s head Plemma and ppos of the leftmost/rightmost dependent and leftmost/rightmost sibling The dependency label of predicate, argument candidate and argument candidate’s dependent The position of the argument candidate with respect to the predicate position in the dep. tree (e.g. CHILD) The position of the head of the dependency relation with respect to the predicate position in the sentence The left-to-right chain of the deplabels of the predicate’s dependents Plemma, ppos and dependency label paths between the predicate and the argument candidates The number of dependency edges between the predicate and the argument candidate
F_{PA}	Plemma and plemma&ppos of the argument candidate Dependency label path between the predicate and the argument candidates
F_G	The sequence of the predicate and the argument labels in the predicate-argument structure (e.g. A0-PRED-A1) Whether the semantic roles defined in frames exist in the structure, (e.g. <i>CONTAINS:A1</i>) The conjunction of the predicate sense and the frame information (e.g. <i>wear:01&CONTAINS:A1</i>)

Table 1: Features for the Structured Model

	Avg.	Ca	Ch	Cz	En	Ge	Jp	Sp
F_P+F_A	79.17	78.00	76.02	85.24	83.09	76.76	77.27	77.83
$F_P+F_A+F_{PA}$	79.58	78.38	76.23	85.14	83.36	78.31	77.72	77.92
$F_P+F_A+F_G$	80.42	79.50	76.96	85.88	84.49	78.64	78.32	79.21
ALL	80.75	79.55	77.20	85.94	84.97	79.62	78.69	79.29
Björkelund	80.80	80.01	78.60	85.41	85.63	79.71	76.30	79.91
Zhao	80.47	80.32	77.72	85.19	85.44	75.99	78.15	80.46
Meza-Ruiz	77.46	78.00	77.73	75.75	83.34	73.52	76.00	77.91

Table 2: Results on the CoNLL-2009 Shared Task dataset (Semantic Labeled F1).

	SENSE	ARG
F_P+F_A	89.65	72.20
$F_P+F_A+F_{PA}$	89.78	72.74
$F_P+F_A+F_G$	89.83	74.11
ALL	90.15	74.46

Table 3: Predicate sense disambiguation and argument role labeling results (average).

used for F_{PA} are inspired by formulae used in the MLN-based SRL systems, such as (Meza-Ruiz and Riedel, 2009b). We used the same feature templates for all languages.

3.2 Results

Table 2 shows the results of the experiments, and also shows the results of the top 3 systems in the CoNLL-2009 Shared Task participants of the *SRL-only* system.

By incorporating F_{PA} , we achieved performance improvement for all languages. This results suggest that it is effective to capture local interdependencies between a predicate sense and one of its argument roles. Comparing the results with F_P+F_A and $F_P+F_A+F_G$, incorporating F_G also contributed performance improvements for all languages, especially the substantial F1 improvement of +1.88 is obtained in German.

Next, we compare our system with top 3 systems in the CoNLL-2009 Shared Task. By incorporating both F_{PA} and F_G , our joint model achieved competitive results compared to the top 2 systems (Björkelund and Zhao), and achieved the better results than the Meza-Ruiz’s system². The systems by Björkelund and Zhao applied feature selection algorithms in order to select the best set of feature templates for each language, requiring about 1 to 2 months to obtain the best feature set. On the other hand, our system achieved the competitive results with the top two systems, despite the fact that we used the same feature templates for all languages without applying any feature engineering procedure.

Table 3 shows the performances of predicate sense disambiguation and argument role labeling separately. In terms of sense disambiguation results, incorporating F_{PA} and F_G worked well. Although incorporating either of F_{PA} and F_G provided improvements of +0.13 and +0.18 on average, adding both factors provided improvements of +0.50. We compared the predicate sense dis-

²The result of Meza-Ruiz for Czech is substantially worse than the other systems because of inappropriate preprocessing for predicate sense disambiguation. Excepting Czech, the average F1 value of the Meza-Ruiz is 77.75, where as our system is 79.89.

ambiguity results of $F_P + F_A$ and ALL with the McNemar test, and the difference was statistically significant ($p < 0.01$). This result suggests that combination of these factors is effective for sense disambiguation.

As for argument role labeling results, incorporating F_{PA} and F_G contributed positively for all languages. Especially, we obtained a substantial gain (+4.18) in German. By incorporating F_{PA} , the system achieved the F1 improvements of +0.54 on average. This result shows that capturing inter-dependencies between a predicate and its arguments contributes to argument role labeling. By incorporating F_G , the system achieved the substantial improvement of F1 (+1.91).

Since both tasks improved by using all factors, we can say that the proposed joint model succeeded in *joint learning* of predicate senses and its argument roles.

4 Conclusion

In this paper, we proposed a structured model that captures both non-local dependencies between arguments, and inter-dependencies between a predicate sense and its argument roles. We designed a linear model-based structured model, and defined four types of factors: predicate factor, argument factor, predicate-argument pairwise factor and global factor for the model. In the experiments, the proposed model achieved competitive results compared to the state-of-the-art systems without any feature engineering.

A further research direction we are investigating is exploitation of unlabeled texts. Semi-supervised semantic role labeling methods have been explored by (Collobert and Weston, 2008; Deschacht and Moens, 2009; Fürstenu and Lapata, 2009), and they have achieved successful outcomes. However, we believe that there is still room for further improvement.

References

Anders Björkelund, Love Hafdel, and Pierre Nugues. 2009. Multilingual semantic role labeling. In *CoNLL-2009*.

Ronan Collobert and Jason Weston. 2008. A unified architecture for natural language processing: Deep neural networks with multitask learning. In *ICML 2008*.

Koby Crammer, Ofer Dekel, Joseph Keshet, Shai

Shalev-Shwartz, and Yoram Singer. 2006. Online passive-aggressive algorithms. *JMLR*, 7:551–585.

Hal Daumé III and Daniel Marcu. 2005. Learning as search optimization: Approximate large margin methods for structured prediction. In *ICML-2005*.

Koen Deschacht and Marie-Francine Moens. 2009. Semi-supervised semantic role labeling using the latent words language model. In *EMNLP-2009*.

Hagen Fürstenu and Mirella Lapata. 2009. Graph alignment for semi-supervised semantic role labeling. In *EMNLP-2009*.

Jan Hajič, Massimiliano Ciaramita, Richard Johansson, Daisuke Kawahara, Maria Antònia Martí, Lluís Màrquez, Adam Meyers, Joakim Nivre, Sebastian Padó, Jan Štěpánek, Pavel Straňák, Mihai Surdeanu, Nianwen Xue, and Yi Zhang. 2009. The CoNLL-2009 shared task: Syntactic and semantic dependencies in multiple languages. In *CoNLL-2009*, Boulder, Colorado, USA.

Richard Johansson and Pierre Nugues. 2008. Dependency-based syntactic-semantic analysis with propbank and nombank. In *CoNLL-2008*.

Jun’Ichi Kazama and Kentaro Torisawa. 2007. A new perceptron algorithm for sequence labeling with non-local features. In *EMNLP-CoNLL 2007*.

Ivan Meza-Ruiz and Sebastian Riedel. 2009a. Jointly identifying predicates, arguments and senses using markov logic. In *HLT/NAACL-2009*.

Ivan Meza-Ruiz and Sebastian Riedel. 2009b. Multilingual semantic role labelling with markov logic. In *CoNLL-2009*.

Sebastian Riedel and Ivan Meza-Ruiz. 2008. Collective semantic role labelling with markov logic. In *CoNLL-2008*.

Mihai Surdeanu, Richard Johansson, Adam Meyers, Lluís Màrquez, and Joakim Nivre. 2008. The CoNLL-2008 shared task on joint parsing of syntactic and semantic dependencies. In *CoNLL-2008*.

Synthia A. Thompson, Roger Levy, and Christopher D. Manning. 2010. A generative model for semantic role labeling. In *Proceedings of the 48th Annual Meeting of the Association of Computational Linguistics (to appear)*.

Kristina Toutanova, Aria Haghighi, and Christopher D. Manning. 2008. A global joint model for semantic role labeling. *Computational Linguistics*, 34(2).

Semantics-Driven Shallow Parsing for Chinese Semantic Role Labeling

Weiwei Sun

Department of Computational Linguistics, Saarland University
German Research Center for Artificial Intelligence (DFKI)
D-66123, Saarbrücken, Germany
wsun@coli.uni-saarland.de

Abstract

One deficiency of current shallow parsing based Semantic Role Labeling (SRL) methods is that syntactic chunks are too small to effectively group words. To partially resolve this problem, we propose semantics-driven shallow parsing, which takes into account both syntactic structures and predicate-argument structures. We also introduce several new “path” features to improve shallow parsing based SRL method. Experiments indicate that our new method obtains a significant improvement over the best reported Chinese SRL result.

1 Introduction

In the last few years, there has been an increasing interest in Semantic Role Labeling (SRL) on several languages, which consists of recognizing arguments involved by predicates of a given sentence and labeling their semantic types. Both full parsing based and shallow parsing based SRL methods have been discussed for English and Chinese. In Chinese SRL, shallow parsing based methods that cast SRL as the classification of syntactic chunks into semantic labels has gained promising results. The performance reported in (Sun et al., 2009) outperforms the best published performance of full parsing based SRL systems.

Previously proposed shallow parsing takes into account only syntactic information and basic chunks are usually too small to group words into argument candidates. This causes one main deficiency of Chinese SRL. To partially resolve this problem, we propose a new shallow parsing. The new chunk definition takes into account both syntactic structure and predicate-argument structures

of a given sentence. Because of the semantic information it contains, we call it semantics-driven shallow parsing. The key idea is to make basic chunks as large as possible but not overlap with arguments. Additionally, we introduce several new “path” features to express more structural information, which is important for SRL.

We present encouraging SRL results on Chinese PropBank (CPB) data. With semantics-driven shallow parsing, our SRL system achieves 76.10 F-measure, with gold segmentation and POS tagging. The performance further achieves 76.46 with the help of new “path” features. These results obtain significant improvements over the best reported SRL performance (74.12) in the literature (Sun et al., 2009).

2 Related Work

CPB is a project to add predicate-argument relations to the syntactic trees of the Chinese TreeBank (CTB). Similar to English PropBank, the arguments of a predicate are labeled with a contiguous sequence of integers, in the form of AN (N is a natural number); the adjuncts are annotated as such with the label AM followed by a secondary tag that represents the semantic classification of the adjunct. The assignment of argument labels is illustrated in Figure 1, where the predicate is the verb “提供/provide”. For example, the noun phrase “保险公司/the insurance company” is labeled as A0, meaning that it is the *proto-Agent* of “提供”.

Sun et al. (2009) explore the Chinese SRL problem on the basis of shallow syntactic information at the level of phrase chunks. They present a semantic chunking method to resolve SRL on basis of shallow parsing. Their method casts SRL as the classification of syntactic chunks with IOB2 representation for semantic roles (i.e. semantic

WORD:	保险	公司	已	为	三峡	工程	提供	保险	服务
	insurance	company	already	for	Sanxia	Project	provide	insurance	service
POS:	[NN	NN]	[AD]	[P]	[NR]	[NN]	[VV]	[NN	NN]
SYN CH:	[NP]		[ADVP]	[PP	NP	NP]	[VP]		[NP]
SEM CH:	B-A0		B-AM-ADV	B-A2	I-A2	I-A2	B-V		B-A1
	The insurance company has provided insurance services for the Sanxia Project.								

Figure 1: An example from Chinese PropBank.

chunks). Two labeling strategies are presented: 1) directly tagging semantic chunks in one-stage, and 2) identifying argument boundaries as a chunking task and labeling their semantic types as a classification task. On the basis of syntactic chunks, they define semantic chunks which do not overlap nor embed using IOB2 representation. Syntactic chunks outside a chunk receive the tag O (Outside). For syntactic chunks forming a chunk of type A^* , the first chunk receives the $B-A^*$ tag (Begin), and the remaining ones receive the tag $I-A^*$ (Inside). Then a SRL system can work directly by using sequence tagging technique. Shallow chunk definition presented in (Chen et al., 2006) is used in their experiments. The definition of syntactic and semantic chunks is illustrated Figure 1. For example, “*保险公司/the insurance company*”, consisting of two nouns, is a noun phrase; in the syntactic chunking stage, its two components “*保险*” and “*公司*” should be labeled as $B-NP$ and $I-NP$. Because this phrase is the *Agent* of the predicate “*提供/provide*”, it takes a semantic chunk label $B-A0$. In the semantic chunking stage, this phrase should be labeled as $B-A0$.

Their experiments on CPB indicate that according to current state-of-the-art of Chinese parsing, SRL systems on basis of full parsing do not perform better than systems based on shallow parsing. They report the best SRL performance with gold segmentation and POS tagging as inputs. This is very different from English SRL. In English SRL, previous work shows that full parsing, both constituency parsing and dependency parsing, is necessary.

Ding and Chang (2009) discuss semantic chunking methods without any parsing information. Different from (Sun et al., 2009), their method formulates SRL as the classification of words with semantic chunks. Comparison of experimental results in their work shows that parsing is necessary for Chinese SRL, and the semantic chunking methods on the basis of shallow parsing outperform the ones without any parsing.

Joint learning of syntactic and semantic structures is another hot topic in dependency parsing research. Some models have been well evaluated in CoNLL 2008 and 2009 shared tasks (Surdeanu et al., 2008; Hajič et al., 2009). The CoNLL 2008/2009 shared tasks propose a unified dependency-based formalism to model both syntactic dependencies and semantic roles for multiple languages. Several joint parsing models are presented in the shared tasks. Our focus is different from the shared tasks. In this paper, we hope to find better syntactic representation for semantic role labeling.

3 Semantics-Driven Shallow Parsing

3.1 Motivation

There are two main jobs of semantic chunking: 1) grouping words as argument candidate and 2) classifying semantic types of possible arguments. Previously proposed shallow parsing only considers syntactic information and basic chunks are usually too small to effectively group words. This causes one main deficiency of semantic chunking. E.g. the argument “*为三峡工程/for the Sanxia Project*” consists of three chunks, each of which only consists of one word. To rightly recognize this $A2$, Semantic chunker should rightly predict three chunk labels. Small chunks also make the important “path” feature sparse, since there are more chunks between a target chunk and the predicate in focus. In this section, we introduce a new chunk definition to improve shallow parsing based SRL, which takes both syntactic and predicate-argument structures into account. The key idea is to make syntactic chunks as large as possible for semantic chunking. The formal definition is as follows.

3.2 Chunk Bracketing

Given a sentence $s = w_1, \dots, w_n$, let $c[i : j]$ denote a constituent that is made up of words between w_i and w_j (including w_i and w_j); let $p_v = \{c[i : j] | c[i : j] \text{ is an argument of } v\}$

	WORD	POS	TARGET	PROPOSITION				CHUNK 1	CHUNK 2
China	中国	NR	-	(A0*	*	*	*	B-NP	B-NP^S
tax	税务	NN	-	*	*	*	*	I-NP	I-NP^S
department	部门	NN	-	*)	*	*	*	I-NP	I-NP^S
stipulate	规定	VV	规定	(V*	*	*	*	O	O
:	:	PU	-	*	*	*	*	O	O
owing	欠缴	VV	欠缴	(A1*	(V*	*	(A0*	O	O
tax payment	税款	NN	-	*	(A1*	*	*	B-NP	B-NP^VP
company	企业	NN	-	*	(A0*	*	*	B-NP	B-NP^NP
Function Word	的	DEG	-	*	*	*	*	O	O
leaders	领导人	NN	-	*	*	*	*)	B-NP	B-NP^NP
not	不	AD	-	*	*	*	(AM-ADV*)	B-ADVP	B-ADVP^VP
can	得	VV	得	*	*	(V*	*	O	O
leave the country	出境	VV	出境	*)	*	*	(V*)	B-VP	B-VP^VP

Figure 2: An example for definition of semantics-driven chunks with IOB2 representation.

denote one predicate-argument structure where v is the predicate in focus. Given a syntactic tree $\mathcal{T}_s = \{c[i : j] | c[i : j] \text{ is a constituent of } s\}$, and its all argument structures $\mathcal{P}_s = \{p_v | v \text{ is a verbal predicate in } s\}$, there is one and only one chunk set $\mathcal{C} = \{c[i : j]\}$ s.t.

1. $\forall c[i : j] \in \mathcal{C}, c[i : j] \in \mathcal{T}_s$;
2. $\forall c[i : j] \in \mathcal{C}, \forall c[i^v : j^v] \in \cup \mathcal{P}_s, j < i^v$ or $i > j^v$ or $i^v \leq i \leq j \leq j^v$;
3. $\forall c[i : j] \in \mathcal{C}$, the parent of $c[i : j]$ does not satisfy the condition 2.
4. $\forall \mathcal{C}'$ satisfies above conditions, $\mathcal{C}' \subset \mathcal{C}$.

The first condition guarantees that every chunk is a constituent. The second condition means that chunks do not overlap with arguments, and further guarantees that semantic chunking can recover all arguments with the last condition. The third condition makes new chunks as big as possible. The last one makes sure that \mathcal{C} contains all sub-components of all arguments. Figure 2 is an example to illustrate our new chunk definition. For example, “中国/Chinese 税务/tax 部分/department” is a constituent of current sentence, and is also an argument of “规定/stipulate”. If we take it as a chunk, it does not conflict with any other arguments, so it is a reasonable syntactic chunk. For the phrase “欠缴/owing 税款/tax payment”, though it does not overlap with the first, third and fourth propositions, it is bigger than the argument “税款” (conflicting with condition 2) while labeling the predicate “欠缴”, so it has to be separated into two chunks. Note that the third condition also guarantees the constituents in \mathcal{C} does not overlap with each other since each one is as large as possible.

So we can still formulate our new shallow parsing as an “IOB” sequence labeling problem.

3.3 Chunk Type

We introduce two types of chunks. The first is simply the phrase type, such as NP , PP , of current chunk. The column *CHUNK 1* illustrates this kind of chunk type definition. The second is more complicated. Inspired by (Klein and Manning, 2003), we split one phrase type into several subsymbols, which contain category information of current constituent’s parent. For example, an NP immediately dominated by a S , will be substituted by NP^S . This strategy severely increases the number of chunk types and make it hard to train chunking models. To shrink this number, we linguistically use a cluster of CTB phrasal types, which was introduced in (Sun and Sui, 2009). The column *CHUNK 2* illustrates this definition. E.g., NP^S implicitly represents *Subject* while NP^VP represents *Object*.

3.4 New Path Features

The *Path* feature is defined as a chain of base phrases between the token and the predicate. At both ends, the chain is terminated with the POS tags of the predicate and the headword of the token. For example, the path feature of “保险公司” in Figure 1 is “公司-ADVP-PP-NP-NP-VV”. Among all features, the “path” feature contains more structural information, which is very important for SRL. To better capture structural information, we introduce several new “path” features. They include:

- NP|PP|VP path: only syntactic chunks that takes tag NP , PP or VP are kept.

When labeling the predicate “出境/leave the country” in Figure 2, this feature of “中国税务部门/Chinese tax departments” is $NP+NP+NP+NP+VP$.

- $V|$ 的 path: a sequential container of POS tags of verbal words and “的”; This feature of “中国税务部门” is $NP+VV+VV+$ 的 $+VV+VP$.
- O2POS path: if a word occupies a chunk label O , use its POS in the path feature. This feature of “中国税务部门” is $NP+VV+PU+VV+NP+NP+DEG+ADVP+VV+VP$.

4 Experiments and Analysis

4.1 Experimental Setting

Experiments in previous work are mainly based on CPB 1.0 and CTB 5.0. We use CoNLL-2005 shared task software to process CPB and CTB. To facilitate comparison with previous work, we use the same data setting with (Xue, 2008). Nearly all previous research on Chinese SRL evaluation use this setting, also including (Ding and Chang, 2008, 2009; Sun et al., 2009; Sun, 2010). The data is divided into three parts: files from chtb_081 to chtb_899 are used as training set; files from chtb_041 to chtb_080 as development set; files from chtb_001 to chtb_040, and chtb_900 to chtb_931 as test set. Both syntactic chunkers and semantic chunkers are trained and evaluated by using the same data set. By using CPB and CTB, we can extract gold standard semantics-driven shallow chunks according to our definition. We use this kind of gold chunks automatically generated from training data to train syntactic chunkers.

For both syntactic and semantic chunking, we used conditional random field model. `Crfsgd`¹, is used for experiments. `Crfsgd` provides a feature template that defines a set of strong word and POS features to do syntactic chunking. We use this feature template to resolve shallow parsing. For semantic chunking, we implement a similar *one-stage* shallow parsing based SRL system described in (Sun et al., 2009). There are two differences between our system and Sun et al.’s system. First, our system uses Start/End method to represent semantic chunks (Kudo and Matsumoto, 2001). Second, word formation features are not used.

Test	P(%)	R(%)	$F_{\beta=1}$
(Chen et al., 2006)	93.51	92.81	93.16
Overall (C1)	91.66	89.13	90.38
Bracketing (C1)	92.31	89.72	91.00
Overall (C2)	88.77	86.71	87.73
Bracketing (C2)	92.71	90.55	91.62

Table 1: Shallow parsing performance.

4.2 Syntactic Chunking Performance

Table 1 shows the performance of shallow syntactic parsing. Line *Chen et al., 2006* is the chunking performance evaluated on syntactic chunk definition proposed in (Chen et al., 2006). The second and third blocks present the chunking performance with new semantics-driven shallow parsing. The second block shows the overall performance when the first kind of chunks type is used, while the last block shows the performance when the more complex chunk type definition is used. For the semantic-driven parsing experiments, we add the *path* from current word to the first verb before or after as two new features. Line *Bracketing* evaluates the word grouping ability of these two kinds of chunks. In other words, detailed phrase types are not considered. Because the two new chunk definitions use the same chunk boundaries, the fourth and sixth lines are comparable. There is a clear decrease between the traditional shallow parsing (Chen et al., 2006) and ours. We think one main reason is that syntactic chunks in our new definition are larger than the traditional ones. An interesting phenomenon is that though the second kind of chunk type definition increases the complexity of the parsing job, it achieves better bracketing performance.

4.3 SRL Performance

Table 2 summarizes the SRL performance. Line *Sun et al., 2009* is the SRL performance reported in (Sun et al., 2009). To the author’s knowledge, this is the best published SRL result in the literature. Line *SRL (Chen et al., 2006)* is the SRL performance of our system. These two systems are both evaluated by using syntactic chunking defined in (Chen et al., 2006). From the first block we can see that our semantic chunking system reaches the state-of-the-art. The second and third blocks in Table 2 present the performance with

¹<http://leon.bottou.org/projects/sgd>

new shallow parsing. Line *SRL (C1)* and *SRL (C2)* show the overall performances with the first and second chunk definition. The lines following are the SRL performance when new “path” features are added. We can see that new “path” features are useful for semantic chunking.

Test	P(%)	R(%)	$F_{\beta=1}$
(Sun et al., 2009)	79.25	69.61	74.12
SRL [(Chen et al., 2006)]	80.87	68.74	74.31
SRL [C1]	80.23	71.00	75.33
+ NP PP VP path	80.25	71.19	75.45
+ V 的 path	80.78	71.67	75.96
+ O2POS path	80.44	71.59	75.76
+ All new path	80.73	72.08	76.16
SRL [C2]	80.87	71.86	76.10
+ All new path	81.03	72.38	76.46

Table 2: SRL performance on the test data. Items in the first column *SRL [(Chen et al., 2006)]*, *SRL [C1]* and *SRL [C2]* respectively denote the SRL systems based on shallow parsing defined in (Chen et al., 2006) and Section 3.

5 Conclusion

In this paper we propose a new syntactic shallow parsing for Chinese SRL. The new chunk definition contains both syntactic structure and predicate-argument structure information. To improve SRL, we also introduce several new “path” features. Experimental results show that our new chunk definition is more suitable for Chinese SRL. It is still an open question what kinds of syntactic information is most important for Chinese SRL. We suggest that our attempt at semantics-driven shallow parsing is a possible way to better exploit this problem.

Acknowledgments

The author is funded both by German Academic Exchange Service (DAAD) and German Research Center for Artificial Intelligence (DFKI).

The author would like to thank the anonymous reviewers for their helpful comments.

References

Wenliang Chen, Yujie Zhang, and Hitoshi Isahara. 2006. An empirical study of Chinese chunking. In *Proceedings of the COLING/ACL 2006 Main Conference Poster Sessions*, pages 97–104. As-

sociation for Computational Linguistics, Sydney, Australia.

Weiwei Ding and Baobao Chang. 2008. Improving Chinese semantic role classification with hierarchical feature selection strategy. In *Proceedings of the EMNLP 2008*, pages 324–333. Association for Computational Linguistics, Honolulu, Hawaii.

Weiwei Ding and Baobao Chang. 2009. Fast semantic role labeling for Chinese based on semantic chunking. In *ICCPOL '09: Proceedings of the 22nd International Conference on Computer Processing of Oriental Languages. Language Technology for the Knowledge-based Economy*, pages 79–90. Springer-Verlag, Berlin, Heidelberg.

Jan Hajič, Massimiliano Ciaramita, Richard Johansson, Daisuke Kawahara, Maria Antònia Martí, Lluís Màrquez, Adam Meyers, Joakim Nivre, Sebastian Padó, Jan Štěpánek, Pavel Straňák, Mihai Surdeanu, Nianwen Xue, and Yi Zhang. 2009. The CoNLL-2009 shared task: Syntactic and semantic dependencies in multiple languages. In *Proceedings of the 13th Conference on Computational Natural Language Learning (CoNLL-2009)*, June 4-5. Boulder, Colorado, USA.

Dan Klein and Christopher D. Manning. 2003. Accurate unlexicalized parsing. In *Proceedings of the 41st Annual Meeting of the Association for Computational Linguistics*, pages 423–430. Association for Computational Linguistics, Sapporo, Japan.

Taku Kudo and Yuji Matsumoto. 2001. Chunking with support vector machines. In *NAACL '01: Second meeting of the North American Chapter of the Association for Computational Linguistics on Language technologies 2001*, pages 1–8. Association for Computational Linguistics, Morristown, NJ, USA.

Weiwei Sun. 2010. Improving Chinese semantic role labeling with rich features. In *Proceedings of the ACL 2010*.

Weiwei Sun and Zhifang Sui. 2009. Chinese function tag labeling. In *Proceedings of the 23rd Pacific Asia Conference on Language, Information and Computation*. Hong Kong.

Weiwei Sun, Zhifang Sui, Meng Wang, and Xin Wang. 2009. Chinese semantic role labeling

with shallow parsing. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing*, pages 1475–1483. Association for Computational Linguistics, Singapore.

Mihai Surdeanu, Richard Johansson, Adam Meyers, Lluís Màrquez, and Joakim Nivre. 2008. The conll 2008 shared task on joint parsing of syntactic and semantic dependencies. In *CoNLL 2008: Proceedings of the Twelfth Conference on Computational Natural Language Learning*, pages 159–177. Coling 2008 Organizing Committee, Manchester, England.

Nianwen Xue. 2008. Labeling Chinese predicates with semantic roles. *Comput. Linguist.*, 34(2):225–255.

Collocation Extraction beyond the Independence Assumption

Gerlof Bouma

Universität Potsdam, Department Linguistik
Campus Golm, Haus 24/35
Karl-Liebknecht-Straße 24–25
14476 Potsdam, Germany
gerlof.bouma@uni-potsdam.de

Abstract

In this paper we start to explore two-part collocation extraction association measures that do not estimate expected probabilities on the basis of the independence assumption. We propose two new measures based upon the well-known measures of mutual information and pointwise mutual information. Expected probabilities are derived from automatically trained Aggregate Markov Models. On three collocation gold standards, we find the new association measures vary in their effectiveness.

1 Introduction

Collocation extraction typically proceeds by scoring collocation candidates with an association measure, where high scores are taken to indicate likely collocationhood. Two well-known such measures are pointwise mutual information (PMI) and mutual information (MI). In terms of observing a combination of words w_1, w_2 , these are:

$$i(w_1, w_2) = \log \frac{p(w_1, w_2)}{p(w_1)p(w_2)}, \quad (1)$$

$$I(w_1, w_2) = \sum_{\substack{x \in \{w_1, \neg w_1\} \\ y \in \{w_2, \neg w_2\}}} p(x, y) i(x, y). \quad (2)$$

PMI (1) is the logged ratio of the observed bigramme probability and the expected bigramme probability under independence of the two words in the combination. MI (2) is the expected outcome of PMI, and measures how much information of the distribution of one word is contained in the distribution of the other. PMI was introduced into the collocation extraction field by Church and Hanks (1990). Dunning (1993) proposed the use of the likelihood-ratio test statistic, which is equivalent to MI up to a constant factor.

Two aspects of (P)MI are worth highlighting. First, the observed occurrence probability p_{obs} is compared to the expected occurrence probability p_{exp} . Secondly, the independence assumption underlies the estimation of p_{exp} .

The first aspect is motivated by the observation that interesting combinations are often those that are *unexpectedly* frequent. For instance, the bigramme *of the* is uninteresting from a collocation extraction perspective, although it probably is amongst the most frequent bigrammes for any English corpus. However, we can expect to frequently observe the combination by mere chance, simply because its parts are so frequent. Looking at p_{obs} and p_{exp} together allows us to recognize these cases (Manning and Schütze (1999) and Evert (2007) for more discussion).

The second aspect, the independence assumption in the estimation of p_{exp} , is more problematic, however, even in the context of collocation extraction. As Evert (2007, p42) notes, the assumption of “independence is extremely unrealistic,” because it ignores “a variety of syntactic, semantic and lexical restrictions.” Consider an estimate for $p_{\text{exp}}(\textit{the the})$. Under independence, this estimate will be high, as *the* itself is very frequent. However, with our knowledge of English syntax, we would say $p_{\text{exp}}(\textit{the the})$ is low. The independence assumption leads to overestimated expectation and *the the* will need to be very frequent for it to show up as a likely collocation. A less contrived example of how the independence assumption might mislead collocation extraction is when bigramme distribution is influenced by compositional, non-collocational, semantic dependencies. Investigating adjective-noun combinations in a corpus, we might find that *beige cloth* gets a high PMI, whereas *beige thought* does not. This does not make the former a collocation or multiword unit. Rather, what we would measure is the tendency to use colours with visible things and not with abstract objects. Syntactic and semantic

associations between words are real dependencies, but they need not be *collocational* in nature. Because of the independence assumption, PMI and MI measure these syntactic and semantic associations just as much as they measure collocational association. In this paper, we therefore experimentally investigate the use of a more informed p_{exp} in the context of collocation extraction.

2 Aggregate Markov Models

To replace p_{exp} under independence, one might consider models with explicit linguistic information, such as a POS-tag bigramme model. This would for instance give us a more realistic p_{exp} (*the the*). However, lexical semantic information is harder to incorporate. We might not know exactly what factors are needed to estimate p_{exp} and even if we do, we might lack the resources to train the resulting models. The only thing we know about estimating p_{exp} is that we need more information than a unigramme model but less than a bigramme model (as this would make $p_{\text{obs}}/p_{\text{exp}}$ uninformative). Therefore, we propose to use Aggregate Markov Models (Saul and Pereira, 1997; Hofmann and Puzicha, 1998; Rooth et al., 1999; Blitzer et al., 2005)¹ for the task of estimating p_{exp} . In an AMM, bigramme probability is not directly modeled, but mediated by a hidden class variable c :

$$p_{\text{amm}}(w_2|w_1) = \sum_c p(c|w_1)p(w_2|c). \quad (3)$$

The number of classes in an AMM determines the amount of dependency that can be captured. In the case of just one class, AMM is equivalent to a unigramme model. AMMs become equivalent to the full bigramme model when the number of classes equals the size of the smallest of the vocabularies of the parts of the combination. Between these two extremes, AMMs can capture syntactic, lexical, semantic and even pragmatic dependencies.

AMMs can be trained with EM, using no more information than one would need for ML bigramme probability estimates. Specifications of the E- and M-steps can be found in any of the four papers cited above – here we follow Saul and Pereira (1997). At each iteration, the model components are updated

¹These authors use very similar models, but with differing terminology and with different goals. The term AMM is used in the first and fourth paper. In the second paper, the models are referred to as Separable Mixture Models. Their use in collocation extraction is to our knowledge novel.

according to:

$$p(c|w_1) \leftarrow \frac{\sum_w n(w_1, w)p(c|w_1, w)}{\sum_{w, c'} n(w_1, w)p(c'|w_1, w)}, \quad (4)$$

$$p(w_2|c) \leftarrow \frac{\sum_w n(w, w_2)p(c|w, w_2)}{\sum_{w, w'} n(w, w')p(c|w, w')}, \quad (5)$$

where $n(w_1, w_2)$ are bigramme counts and the posterior probability of a hidden category c is estimated by:

$$p(c|w_1, w_2) = \frac{p(c|w_1)p(w_2|c)}{\sum_{c'} p(c'|w_1)p(w_2|c')}. \quad (6)$$

Successive updates converge to a local maximum of the AMM’s log-likelihood.

The definition of the counterparts to (P)MI without the independence assumption, the AMM-ratio and AMM-divergence, is now straightforward:

$$r_{\text{amm}}(w_1, w_2) = \log \frac{p(w_1, w_2)}{p(w_1)p_{\text{amm}}(w_2|w_1)}, \quad (7)$$

$$d_{\text{amm}}(w_1, w_2) = \sum_{\substack{x \in \{w_1, \neg w_1\} \\ y \in \{w_2, \neg w_2\}}} p(x, y) r_{\text{amm}}(x, y). \quad (8)$$

The free parameter in these association measures is the number of hidden classes in the AMM, that is, the amount of dependency between the bigramme parts used to estimate p_{exp} . Note that AMM-ratio and AMM-divergence with one hidden class are equivalent to PMI and MI, respectively. It can be expected that in different corpora and for different types of collocation, different settings of this parameter are suitable.

3 Evaluation

3.1 Data and procedure

We apply AMM-ratio and AMM-divergence to three collocation gold standards. The effectiveness of association measures in collocation extraction is measured by ranking collocation candidates after the scores defined by the measures, and calculating average precision of these lists against the gold standard annotation. We consider the newly proposed AMM-based measures for a varying number of hidden categories. The new measures are compared against two baselines: ranking by frequency (p_{obs}) and random ordering. Because AMM-ratio and -divergence with one hidden class boil down to PMI and MI (and thus log-likelihood ratio), the evaluation contains an implicit comparison with

these canonical measures, too. However, the results will not be state-of-the-art: for the datasets investigated below, there are more effective extraction methods based on supervised machine learning (Pecina, 2008).

The first gold standard used is the German adjective-noun dataset (Evert, 2008). It contains 1212 A-N pairs taken from a German newspaper corpus. We consider three subtasks, depending on how strict we define true positives. We used the bigramme frequency data included in the resource. We assigned all types with a token count ≤ 5 to one type, resulting in AMM training data of 10k As, 20k Ns and 446k A-N pair types.

The second gold standard consists of 5102 German PP-verb combinations, also sampled from newspaper texts (Krenn, 2008). The data contains annotation for support verb constructions (FVGs) and figurative expressions. This resource also comes with its own frequency data. After frequency thresholding, AMMs are trained on 46k PPs, 7.6k Vs, and 890k PP-V pair types.

Third and last is the English verb-particle construction (VPC) gold standard (Baldwin, 2008), consisting of 3078 verb-particle pairs and annotation for transitive and intransitive idiomatic VPCs. We extract frequency data from the BNC, following the methods described in Baldwin (2005). This results in two slightly different datasets for the two types of VPC. For the intransitive VPCs, we train AMMs on 4.5k Vs, 35 particles, and 43k pair types. For the transitive VPCs, we have 5k Vs, 35 particles and 54k pair types.

All our EM runs start with randomly initialized model vectors. In Section 3.3 we discuss the impact of model variation due to this random factor.

3.2 Results

German A-N collocations The top slice in Table 1 shows results for the three subtasks of the A-N dataset. We see that using AMM-based p_{exp} initially improves average precision, for each task and for both the ratio and the divergence measure. At their maxima, the informed measures outperform both baselines as well as PMI and MI/log-likelihood ratio (# classes=1). The AMM-ratio performs best for 16-class AMMs, the optimum for AMM-divergence varies slightly.

It is likely that the drop in performance for the larger AMM-based measures is due to the AMMs learning the collocations themselves. That is, the

AMMs become rich enough to not only capture the broadly applicative distributional influences of syntax and semantics, but also provide accurate p_{exp} s for individual, distributionally deviant combinations – like collocations. An accurate p_{exp} results in a low association score.

One way of inspecting what kind of dependencies the AMMs pick up is to cluster the data with them. Following Blitzer et al. (2005), we take the 200 most frequent adjectives and assign them to the category that maximizes $p(c|w_1)$; likewise for nouns and $p(w_2|c)$. Four selected clusters (out of 16) are given in Table 2.² The esoteric class 1 contains ordinal numbers and nouns that one typically uses those with, including references to temporal concepts. Class 2 and 3 appear more semantically motivated, roughly containing human and collective denoting nouns, respectively. Class 4 shows a group of adjectives denoting colours and/or political affiliations and a less coherent set of nouns, although the noun cluster can be understood if we consider individual adjectives that are associated with this class. Our informal impression from looking at clusters is that this is a common situation: as a whole, a cluster cannot be easily characterized, although for subsets or individual pairs, one can get an intuition for why they are in the same class. Unfortunately, we also see that some actual collocations are clustered in class 4, such as *gelbe Karte* ‘warning’ (lit.: ‘yellow card’) and *dickes Auto* ‘big (lit.: fat) car’.

German PP-Verb collocations The second slice in Table 1 shows that, for both subtypes of PP-V collocation, better p_{exp} -estimates lead to *decreased* average precision. The most effective AMM-ratio and -distance measures are those equivalent to (P)MI. Apparently, the better p_{exp} s are unfortunate for the extraction of the type of collocations in this dataset.

The poor performance of PMI on these data – clearly below frequency – has been noticed before by Krenn and Evert (2001). A possible explanation for the lack of improvement in the AMMs lies in the relatively high performing frequency baselines. The frequency baseline for FVGs is five times the

²An anonymous reviewer rightly warns against sketching an overly positive picture of the knowledge captured in the AMMs by only presenting a few clusters. However, the clustering performed here is only secondary to our main goal of improving collocation extraction. The model inspection should thus not be taken as an evaluation of the quality of the models as clustering models.

		# classes										Rnd	Frq
		1	2	4	8	16	32	64	128	256	512		
A-N													
category 1	r_{amm}	45.6	46.4	47.6	47.3	48.3	48.0	47.0	46.1	44.7	41.9	30.1	32.2
	d_{amm}	42.3	42.9	44.4	45.2	46.1	46.5	45.0	46.3	45.5	45.5		
category 1–2	r_{amm}	55.7	56.3	57.4	57.5	58.1	58.1	57.7	56.9	55.7	52.8	43.1	47.0
	d_{amm}	56.3	57.0	58.1	58.4	59.8	60.1	59.3	60.6	59.2	59.3		
category 1–3	r_{amm}	62.3	62.8	63.9	64.0	64.4	62.2	62.2	62.7	62.4	60.0	52.7	56.4
	d_{amm}	64.3	64.7	65.9	66.6	66.7	66.3	66.3	65.4	66.0	64.7		
PP-V													
figurative	r_{amm}	7.5	6.1	6.4	6.0	5.6	5.4	4.5	4.2	3.8	3.5	3.3	10.5
	d_{amm}	14.4	13.0	13.3	13.1	12.2	11.2	9.0	7.7	6.9	5.7		
FVG	r_{amm}	4.1	3.4	3.4	3.0	2.9	2.7	2.2	2.1	2.0	2.0	3.0	14.7
	d_{amm}	15.3	12.7	12.6	10.7	9.0	7.7	3.4	3.2	2.5	2.3		
VPC													
intransitive	r_{amm}	9.3	9.2	9.0	8.3	5.5	5.3					4.8	14.7
	d_{amm}	12.2	12.2	14.0	16.3	6.9	5.8						
transitive	r_{amm}	16.4	14.8	15.2	14.5	11.3	10.0					10.1	20.1
	d_{amm}	19.6	17.3	20.7	23.8	12.8	10.1						

Table 1: Average precision for AMM-based association measures and baselines on three datasets.

Cl	Adjective	Noun
1	<i>dritt</i> ‘third’, <i>erst</i> ‘first’, <i>fünft</i> ‘fifth’, <i>halb</i> ‘half’, <i>kommend</i> ‘next’, <i>laufend</i> ‘current’, <i>letzt</i> ‘last’, <i>nah</i> ‘near’, <i>paar</i> ‘pair’, <i>vergangen</i> ‘last’, <i>viert</i> ‘fourth’, <i>wenig</i> ‘few’, <i>zweit</i> ‘second’	<i>Jahr</i> ‘year’, <i>Klasse</i> ‘class’, <i>Linie</i> ‘line’, <i>Mal</i> ‘time’, <i>Monat</i> ‘month’, <i>Platz</i> ‘place’, <i>Rang</i> ‘grade’, <i>Runde</i> ‘round’, <i>Saison</i> ‘season’, <i>Satz</i> ‘sentence’, <i>Schritt</i> ‘step’, <i>Sitzung</i> ‘session’, <i>Sonntag</i> ‘Sunday’, <i>Spiel</i> ‘game’, <i>Stunde</i> ‘hour’, <i>Tag</i> ‘day’, <i>Woche</i> ‘week’, <i>Wochenende</i> ‘weekend’
2	<i>aktiv</i> ‘active’, <i>alt</i> ‘old’, <i>ausländisch</i> ‘foreign’, <i>betroffen</i> ‘concerned’, <i>jung</i> ‘young’, <i>lebend</i> ‘alive’, <i>meist</i> ‘most’, <i>unbekannt</i> ‘unknown’, <i>viel</i> ‘many’	<i>Besucher</i> ‘visitor’, <i>Bürger</i> ‘citizens’, <i>Deutsche</i> ‘German’, <i>Frau</i> ‘woman’, <i>Gast</i> ‘guest’, <i>Jugendliche</i> ‘youth’, <i>Kind</i> ‘child’, <i>Leute</i> ‘people’, <i>Mädchen</i> ‘girl’, <i>Mann</i> ‘man’, <i>Mensch</i> ‘human’, <i>Mitglied</i> ‘member’
3	<i>deutsch</i> ‘German’, <i>europäisch</i> ‘European’, <i>ganz</i> ‘whole’, <i>gesamt</i> ‘whole’, <i>international</i> ‘international’, <i>national</i> ‘national’, <i>örtlich</i> ‘local’, <i>ostdeutsch</i> ‘East-German’, <i>privat</i> ‘private’, <i>rein</i> ‘pure’, <i>sogenannt</i> ‘so-called’, <i>sonstig</i> ‘other’, <i>westlich</i> ‘western’	<i>Betrieb</i> ‘company’, <i>Familie</i> ‘family’, <i>Firma</i> ‘firm’, <i>Gebiet</i> ‘area’, <i>Gesellschaft</i> ‘society’, <i>Land</i> ‘country’, <i>Mannschaft</i> ‘team’, <i>Markt</i> ‘market’, <i>Organisation</i> ‘organisation’, <i>Staat</i> ‘state’, <i>Stadtteil</i> ‘city district’, <i>System</i> ‘system’, <i>Team</i> ‘team’, <i>Unternehmen</i> ‘enterprise’, <i>Verein</i> ‘club’, <i>Welt</i> ‘world’
4	<i>blau</i> ‘blue’, <i>dick</i> ‘fat’, <i>gelb</i> ‘yellow’, <i>grün</i> ‘green’, <i>linke</i> ‘left’, <i>recht</i> ‘right’, <i>rot</i> ‘red’, <i>schwarz</i> ‘black’, <i>white</i> ‘weiß’	<i>Auge</i> ‘eye’, <i>Auto</i> ‘car’, <i>Haar</i> ‘hair’, <i>Hand</i> ‘hand’, <i>Karte</i> ‘card’, <i>Stimme</i> ‘voice/vote’

Table 2: Selected adjective-noun clusters from a 16-class AMM.

random baseline, and MI does not outperform it by much. Since the AMMs provide a better fit for the more frequent pairs in the training data, they might end up providing too good p_{exp} -estimates for the true collocations from the beginning.

Further investigation is needed to find out whether this situation can be ameliorated and, if not, whether we can systematically identify for what kind of collocation extraction tasks using better p_{exp} s is simply not a good idea.

English Verb-Particle constructions The last gold standard is the English VPC dataset, shown in the bottom slice of Table 1. We have only used class-sizes up to 32, as there are only 35 particle types. We can clearly see the effect of the largest AMMs approaching the full bigramme model as

average precision here approaches the random baseline. The VPC extraction task shows a difference between the two AMM-based measures: AMM-ratio does not improve at all, remaining below the frequency baseline. AMM-divergence, however, shows a slight decrease in precision first, but ends up performing above the frequency baseline for the 8-class AMMs in both subtasks.

Table 3 shows four clusters of verbs and particles. The large first cluster contains verbs that involve motion/displacement of the subject or object and associated particles, for instance *walk about* or *push away*. Interestingly, the description of the gold standard gives exactly such cases as negatives, since they constitute compositional verb-particle constructions (Baldwin, 2008). Classes 2 and 3 show syntactic dependencies, which helps

Cl	Verb	Particle
1	<i>break, bring, come, cut, drive, fall, get, go, lay, look, move, pass, push, put, run, sit, throw, turn, voice, walk</i>	<i>across, ahead, along, around, away, back, backward, down, forward, into, over, through, together</i>
2	<i>accord, add, apply, give, happen, lead, listen, offer, pay, present, refer, relate, return, rise, say, sell, send, speak, write</i>	<i>astray, to</i>
3	<i>know, talk, tell, think</i>	<i>about</i>
4	<i>accompany, achieve, affect, cause, create, follow, hit, increase, issue, mean, produce, replace, require, sign, support</i>	<i>by</i>

Table 3: Selected verb-particle clusters from an 8-class AMM on transitive data.

collocation extraction by decreasing the impact of verb-preposition associations that are due to PP-selecting verbs. Class 4 shows a third type of distributional generalization: the verbs in this class are all frequently used in the passive.

3.3 Variation due to local optima

We start each EM run with a random initialization of the model parameters. Since EM finds local rather than global optima, each run may lead to different AMMs, which in turn will affect AMM-based collocation extraction. To gain insight into this variation, we have trained 40 16-class AMMs on the A-N dataset. Table 4 gives five point summaries of the average precision of the resulting 40 ‘association measures’. Performance varies considerably, spanning 2–3 percentage points in each case. The models consistently outperform (P)MI in Table 1, though.

Several techniques might help to address this variation. One might try to find a good fixed way of initializing EM or to use EM variants that reduce the impact of the initial state (Smith and Eisner, 2004, a.o.), so that a run with the same data and the same number of classes will always learn (almost) the same model. On the assumption that an average over several runs will vary less than individual runs, we have also constructed a combined p_{exp} by averaging over 40 p_{exp} s. The last column

		Variation in avg precision					Comb
		min	q1	med	q3	max	
A-N							
cat 1	r_{amm}	46.5	47.3	47.9	48.4	49.1	48.4
	d_{amm}	44.4	45.4	45.8	46.1	47.1	46.4
cat 1–2	r_{amm}	56.7	57.2	57.9	58.2	59.0	58.2
	d_{amm}	58.1	58.8	59.2	59.4	60.4	60.0
cat 1–3	r_{amm}	63.0	63.7	64.2	64.6	65.3	64.6
	d_{amm}	65.2	66.0	66.4	66.6	67.6	66.9

Table 4: Variation on A-N data over 40 EM runs and result of combining p_{exp} s.

in Table 4 shows this combined estimator leads to good extraction results.

4 Conclusions

In this paper, we have started to explore collocation extraction beyond the assumption of independence. We have introduced two new association measures that do away with this assumption in the estimation of expected probabilities. The success of using these association measures varies. It remains to be investigated whether they can be improved more.

A possible obstacle in the adoption of AMMs in collocation extraction is that we have not provided any heuristic for setting the number of classes for the AMMs. We hope to be able to look into this question in future research. Luckily, for the AN and VPC data, the best models are not that large (in the order of 8–32 classes), which means that model fitting is fast enough to experiment with different settings. In general, considering these smaller models might suffice for tasks that have a fairly restricted definition of collocation candidate, like the tasks in our evaluation do. Because AMM fitting is unsupervised, selecting a class size is in this respect no different from selecting a suitable association measure from the canon of existing measures.

Future research into association measures that are not based on the independence assumption will also include considering different EM variants and other automatically learnable models besides the AMMs used in this paper. Finally, the idea of using an informed estimate of expected probability in an association measure need not be confined to (P)MI, as there are many other measures that employ expected probabilities.

Acknowledgements

This research was carried out in the context of the SFB 632 *Information Structure*, subproject D4: *Methoden zur interaktiven linguistischen Korpusanalyse von Informationsstruktur*.

References

- Timothy Baldwin. 2005. The deep lexical acquisition of english verb-particle constructions. *Computer Speech and Language, Special Issue on Multiword Expressions*, 19(4):398–414.
- Timothy Baldwin. 2008. A resource for evaluating the deep lexical acquisition of English verb-particle constructions. In *Proceedings of the LREC 2008 Workshop Towards a Shared Task for Multiword Expressions (MWE 2008)*, pages 1–2, Marrakech.
- John Blitzer, Amir Globerson, and Fernando Pereira. 2005. Distributed latent variable models of lexical co-occurrences. In *Tenth International Workshop on Artificial Intelligence and Statistics*.
- Kenneth W. Church and Patrick Hanks. 1990. Word association norms, mutual information, and lexicography. *Computational Linguistics*, 16(1):22–29.
- Ted Dunning. 1993. Accurate methods for the statistics of surprise and coincidence. *Computational Linguistics*, 19(1):61–74.
- Stefan Evert. 2007. Corpora and collocations. Extended Manuscript of Chapter 58 of A. Lüdeling and M. Kytö, 2008, *Corpus Linguistics. An International Handbook*, Mouton de Gruyter, Berlin.
- Stefan Evert. 2008. A lexicographic evaluation of German adjective-noun collocations. In *Proceedings of the LREC 2008 Workshop Towards a Shared Task for Multiword Expressions (MWE 2008)*, pages 3–6, Marrakech.
- Thomas Hofmann and Jan Puzicha. 1998. Statistical models for co-occurrence data. Technical report, MIT. AI Memo 1625, CBCL Memo 159.
- Brigitte Krenn and Stefan Evert. 2001. Can we do better than frequency? a case study on extracting PP-verb collocations. In *Proceedings of the ACL Workshop on Collocations*, Toulouse.
- Brigitte Krenn. 2008. Description of evaluation resource – German PP-verb data. In *Proceedings of the LREC 2008 Workshop Towards a Shared Task for Multiword Expressions (MWE 2008)*, pages 7–10, Marrakech.
- Chris Manning and Hinrich Schütze. 1999. *Foundations of Statistical Natural Language Processing*. MIT Press, Cambridge, MA.
- Pavel Pecina. 2008. A machine learning approach to multiword expression extraction. In *Proceedings of the LREC 2008 Workshop Towards a Shared Task for Multiword Expressions (MWE 2008)*, pages 54–57, Marrakech.
- Mats Rooth, Stefan Riester, Detlef Prescher, Glenn Carroll, and Franz Beil. 1999. Inducing a semantically annotated lexicon via em-based clustering. In *Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics*, College Park, MD.
- Lawrence Saul and Fernando Pereira. 1997. Aggregate and mixed-order markov models for statistical language processing. In *Proceedings of the Second Conference on Empirical Methods in Natural Language Processing*, pages 81–89.
- Noah A. Smith and Jason Eisner. 2004. Annealing techniques for unsupervised statistical language learning. In *Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics*.

Automatic Collocation Suggestion in Academic Writing

Jian-Cheng Wu¹

Yu-Chia Chang^{1,*}

Teruko Mitamura²

Jason S. Chang¹

¹National Tsing Hua University
Hsinchu, Taiwan

{wujc86, richtrf, jason.jschang}
@gmail.com

²Carnegie Mellon University
Pittsburgh, United States

teruko@cs.cmu.edu

Abstract

In recent years, collocation has been widely acknowledged as an essential characteristic to distinguish native speakers from non-native speakers. Research on academic writing has also shown that collocations are not only common but serve a particularly important discourse function within the academic community. In our study, we propose a machine learning approach to implementing an online collocation writing assistant. We use a data-driven classifier to provide collocation suggestions to improve word choices, based on the result of classification. The system generates and ranks suggestions to assist learners' collocation usages in their academic writing with satisfactory results.

1 Introduction

The notion of collocation has been widely discussed in the field of language teaching for decades. It has been shown that collocation, a successive common usage of words in a chain, is important in helping language learners achieve native-like fluency. In the field of English for Academic Purpose, more and more researchers are also recognizing this important feature in academic writing. It is often argued that collocation can influence the effectiveness of a piece of writing and the lack of such knowledge might cause cumulative loss of precision (Howarth, 1998).

Many researchers have discussed the function of collocations in the highly conventionalized and specialized writing used within academia. Research also identified noticeable increases in the quantity and quality of collocational usage by

native speakers (Howarth, 1998). Granger (1998) reported that learners underuse native-like collocations and overuse atypical word combinations. This disparity in collocation usage between native and non-native speakers is clear and should receive more attention from the language technology community.

To tackle such word usage problems, traditional language technology often employs a database of the learners' common errors that are manually tagged by teachers or specialists (e.g. Shei and Pain, 2000; Liu, 2002). Such system then identifies errors via string or pattern matching and offer only pre-stored suggestions. Compiling the database is time-consuming and not easily maintainable, and the usefulness is limited by the manual collection of pre-stored suggestions. Therefore, it is beneficial if a system can mainly use untagged data from a corpus containing correct language usages rather than the error-tagged data from a learner corpus. A large corpus of correct language usages is more readily available and useful than a small labeled corpus of incorrect language usages.

For this suggestion task, the large corpus not only provides us with a rich set of common collocations but also provides the context within which these collocations appear. Intuitively, we can take account of such context of collocation to generate more suitable suggestions. Contextual information in this sense often entails more linguistic clues to provide suggestions within sentences or paragraph. However, the contextual information is messy and complex and thus has long been overlooked or ignored. To date, most fashionable suggestion methods still rely upon the linguistic components within collocations as well as the linguistic relationship between misused words and their correct counterparts (Chang et al., 2008; Liu, 2009).

In contrast to other research, we employ contextual information to automate suggestions for verb-noun lexical collocation. Verb-noun collocations are recognized as presenting the most

* Corresponding author: Yu-chia Chang (Email address: richtrf@gmail.com)

challenge to students (Howarth, 1996; Liu, 2002). More specifically, in this preliminary study we start by focusing on the word choice of verbs in collocations which are considered as the most difficult ones for learners to master (Liu, 2002; Chang, 2008). The experiment confirms that our collocation writing assistant proves the feasibility of using machine learning methods to automatically prompt learners with collocation suggestions in academic writing.

2 Collocation Checking and Suggestion

This study aims to develop a web service, *Collocation Inspector* (shown in Figure 1) that accepts sentences as input and generates the related candidates for learners.

In this paper, we focus on automatically providing academic collocation suggestions when users are writing up their abstracts. After an abstract is submitted, the system extracts linguistic features from the user’s text for machine learning model. By using a corpus of published academic texts, we hope to match contextual linguistic clues from users’ text to help elicit the most relevant suggestions. We now formally state the problem that we are addressing:

Problem Statement: Given a sentence S written by a learner and a reference corpus RC , our goal is to output a set of most probable suggestion candidates c_1, c_2, \dots, c_m . For this, we train a classifier MC to map the context (represented as feature set f_1, f_2, \dots, f_n) of each sentence in RC to the collocations. At run-time, we predict these collocations for S as suggestions.

2.1 Academic Collocation Checker Training Procedures

Sentence Parsing and Collocation Extraction: We start by collecting a large number of abstracts from the Web to develop a reference corpus for collocation suggestion. And we continue to identify collocations in each sentence for the subsequent processing.

Collocation extraction is an essential step in preprocessing data. We only expect to extract the collocation which comprises components having a syntactic relationship with one another. However, this extraction task can be complicated. Take the following scholarly sentence from the reference corpus as an example (example (1)):

(1) We introduce a novel method for learning to find documents on the web.



Figure 1. The interface for the collocation suggestion

```
nsubj (introduce-2, We-1)
det (method-5, a-3)
amod (method-5, novel-4)
dobj (introduce-2, method-5)
prepc_for (introduce-2, learning-7)
aux (find-9, to-8)
... ..
```

Figure 2. Dependency parsing of Example (1)

Traditionally, through part-of-speech tagging, we can obtain a tagged sentence as follows (example (2)). We can observe that the desired collocation “introduce method”, conforming to “VERB+NOUN” relationship, exists within the sentence. However, the distance between these two words is often flexible, not necessarily rigid. Heuristically writing patterns to extract such verb and noun might not be effective. The patterns between them can be tremendously varied. In addition, some verbs and nouns are adjacent, but they might be intervened by clause and thus have no syntactic relation with one another (e.g. “propose model” in example (3)).

(2) We/PRP introduce/VB a/DT
 novel/JJ method/NN for/IN
 learning/VBG to/TO find/VB
 documents/NNS on/IN the/DT
 web/NN ./.

(3) We proposed that the web-based model would be more effective than corpus-based one.

A natural language parser can facilitate the extraction of the target type of collocations. Such parser is a program that works out the grammatical structure of sentences, for instance, by identifying which group of words go together or which

word is the subject or object of a verb. In our study, we take advantage of a dependency parser, *Stanford Parser*, which extracts typed dependencies for certain grammatical relations (shown in Figure 2). Within the parsed sentence of example (1), we can notice that the extracted dependency “dobj (introduce-2, method-4)” meets the criterion.

Using a Classifier for the Suggestion task: A classifier is a function generally to take a set of attributes as an input and to provide a tagged class as an output. The basic way to build a classifier is to derive a regression formula from a set of tagged examples. And this trained classifier can thus make predication and assign a tag to any input data.

The suggestion task in this study will be seen as a classification problem. We treat the collocation extracted from each sentence as the class tag (see examples in Table 1). Hopefully, the system can learn the rules between tagged classes (i.e. collocations) and example sentences (i.e. scholarly sentences) and can predict which collocation is the most appropriate one given attributes extracted from the sentences.

Another advantage of using a classifier to automate suggestion is to provide alternatives with regard to the similar attributes shared by sentences. In Table 1, we can observe that these collocations exhibit a similar discourse function and can thus become interchangeable in these sentences. Therefore, based on the outputs along with the probability from the classifier, we can provide more than one adequate suggestions.

Feature Selection for Machine Learning: In the final stage of training, we build a statistical machine-learning model. For our task, we can use a supervised method to automatically learn the relationship between collocations and example sentences.

We choose Maximum Entropy (ME) as our training algorithm to build a collocation suggestion classifier. One advantage of an ME classifier is that in addition to assigning a classification it can provide the probability of each assignment. The ME framework estimates probabilities based on the principle of making as few assumptions as possible. Such constraints are derived from the training data, expressing relationships between features and outcomes.

Moreover, an effective feature selection can increase the precision of machine learning. In our study, we employ the contextual features which

Table 1. Example sentences and class tags (collocations)

Example Sentence	Class tag
We introduce a novel method for learning to find documents on the web.	introduce
We presented a method of improving Japanese dependency parsing by using large-scale statistical information.	present
In this paper, we will describe a method of identifying the syntactic role of antecedents, which consists of two phases	describe
In this paper, we suggest a method that automatically constructs an NE tagged corpus from the web to be used for learning of NER systems.	suggest

consist of two elements, the *head* and the *ngram* of context words:

Head: Each collocation comprises two parts, collocate and head. For example, in a given verb-noun collocation, the verb is the collocate as well as the target for which we provide suggestions; the noun serves as the head of collocation and convey the essential meaning of the collocation. We use the head as a feature to condition the classifier to generate candidates relevant to a given head.

Ngram: We use the context words around the target collocation by considering the corresponding unigrams and bigrams words within the sentence. Moreover, to ensure the relevance, those context words, before and after the punctuation marks enclosing the collocation in question, will be excluded. We use the parsed sentence from previous step (example (2)) to show the extracted context features¹ (example (4)):

```
(4) CN=method UniV_L=we
UniV_R=a UniV_R=novel UniN_L=a
UniN_L=novel UniN_R=for
UniN_R=learn BiV_R=a_novel
BiN_L=a_novel BiN_R=for_learn
BiV_I=we_a BiN_I=novel_for
```

¹ CN refers to the head within collocation. Uni and Bi indicate the unigram and bigram context words of window size two respectively. V and N differentiate the contexts related to verb or noun. The ending alphabets L, R, I show the position of the words in context, L = left, R = right, and I = in between.

2.2 Automatic Collocation Suggestion at Run-time

After the ME classifier is automatically trained, the model is used to find out the best collocation suggestion. Figure 3 shows the algorithm of producing suggestions for a given sentence. The input is a learner’s sentence in an abstract, along with an ME model trained from the reference corpus.

In Step (1) of the algorithm, we parse the sentence for data preprocessing. Based on the parser output, we extract the collocation from a given sentence as well as generate features sets in Step (2) and (3). After that in Step (4), with the trained machine-learning model, we obtain a set of likely collocates with probability as predicted by the ME model. In Step (5), SuggestionFilter singles out the valid collocation and returns the best collocation suggestion as output in Step (6). For example, if a learner inputs the sentence like Example (5), the features and output candidates are shown in Table 2.

(5) There are many investigations about wireless network communication, especially it is important to add Internet transfer calculation speeds.

3 Experiment

From an online research database, *CiteSeer*, we have collected a corpus of 20,306 unique abstracts, which contained 95,650 sentences. To train a Maximum Entropy classifier, 46,255 collocations are extracted and 790 verbal collocates are identified as tagged classes for collocation suggestions. We tested the classifier on scholarly sentences in place of authentic student writings which were not available at the time of this pilot study. We extracted 364 collocations among 600 randomly selected sentences as the held out test data not overlapping with the training set. To automate the evaluation, we blank out the verb collocates within these sentences and treat these verbs directly as the only correct suggestions in question, although two or more suggestions may be interchangeable or at least appropriate. In this sense, our evaluation is an underestimate of the performance of the proposed method.

While evaluating the quality of the suggestions provided by our system, we used the mean reciprocal rank (MRR) of the first relevant suggestions returned so as to assess whether the suggestion list contains an answer and how far up the answer is in the list as a quality metric of the sys-

Procedure CollocationSuggestion(<i>sent</i> , <i>MEmodel</i>) (1) <i>parsedSen</i> = Parsing(<i>sent</i>) (2) <i>extractedColl</i> = CollocationExtraction(<i>parsedSent</i>) (3) <i>features</i> = AssignFeature(<i>ParsedSent</i>) (4) <i>probCollection</i> = MEprob(<i>features</i> , <i>MEmodel</i>) (5) <i>candidate</i> = SuggestionFilter(<i>probCollection</i>) (6) Return candidate

Figure 3. Collocation Suggestion at Run-time

Table 2. An example from learner’s sentence

Extracted Collocation	Features	Ranked Candidates
add speed	CN=speed	
	UniV_L=important	
	UniV_L=to	
	UniV_R=internet	improve
	UniV_R=transfer	increase
	UniN_L=transfer	determine
	UniN_L=calculation	maintain
	BiV_L=important to	...
	BiV_R=internet_transfer	...
	BiN_L=transfer_calca- tion	
	BiV_I=to_intenet	

Table 3. MRR for different feature sets

Feature Sets Included In Classifier	MRR
Features of HEAD	0.407
Features of CONTEXT	0.469
Features of HEAD+CONTEXT	0.518

tem output. Table 3 shows that the best MRR of our prototype system is 0.518. The results indicate that on average users could easily find answers (exactly reproduction of the blanked out collocates) in the first two to three ranking of suggestions. It is very likely that we get a much higher MMR value if we would go through the lists and evaluate each suggestion by hand. Moreover, in Table 3, we can further notice that contextual features are quite informative in comparison with the baseline feature set containing merely the feature of HEAD. Also the integrated feature set of HEAD and CONTEXT together achieves a more satisfactory suggestion result.

4 Conclusion

Many avenues exist for future research that are important for improving the proposed method. For example, we need to carry out the experiment on authentic learners’ texts. We will conduct a user study to investigate whether our system would improve a learner’s writing in a real setting. Additionally, adding classifier features based on the translation of misused words in learners’ text could be beneficial (Chang et al.,

2008). The translation can help to resolve prevalent collocation misuses influenced by a learner's native language. Yet another direction of this research is to investigate if our methodology is applicable to other types of collocations, such as AN and PN in addition to VN dealt with in this paper.

In summary, we have presented an unsupervised method for suggesting collocations based on a corpus of abstracts collected from the Web. The method involves selecting features from the reference corpus of the scholarly texts. Then a classifier is automatically trained to determine the most probable collocation candidates with regard to the given context. The preliminary results show that it is beneficial to use classifiers for identifying and ranking collocation suggestions based on the context features.

Reference

- Y. Chang, J. Chang, H. Chen, and H. Liou. 2008. An automatic collocation writing assistant for Taiwanese EFL learners: A case of corpus-based NLP technology. *Computer Assisted Language Learning*, 21(3), pages 283-299.
- S. Granger. 1998. Prefabricated patterns in advanced EFL writing: collocations and formulae. In Cowie, A. (ed.) *Phraseology: theory, analysis and applications*. Oxford University Press, Oxford, pages 145-160.
- P. Howarth. 1996. *Phraseology in English Academic Writing*. Tübingen: Max Niemeyer Verlag.
- P. Howarth. 1998. The phraseology of learner's academic writing. In Cowie, A. (ed.) *Phraseology: theory, analysis and applications*. Oxford University Press, Oxford, pages 161-186.
- D. Hawking and N. Craswell. 2002. Overview of the TREC-2001 Web track. In *Proceedings of the 10th Text Retrieval Conference (TREC 2001)*, pages 25-31.
- L. E. Liu. 2002. A corpus-based lexical semantic investigation of verb-noun miscollocations in Taiwan learners' English. *Unpublished master's thesis*, Tamkang University, Taipei, January.
- A. L. Liu, D. Wible, and N. L. Tsao. 2009. Automated suggestions for miscollocations. In *Proceedings of the Fourth Workshop on Innovative Use of NLP for Building Educational Applications*, pages 47-50.
- C. C. Shei and H. Pain. 2000. An ESL writer's collocational aid. *Computer Assisted Language Learning*, 13, pages 167-182.

Event-based Hyperspace Analogue to Language for Query Expansion

Tingxu Yan

Tianjin University
Tianjin, China

sunriser2008@gmail.com

Tamsin Maxwell

University of Edinburgh
Edinburgh, United Kingdom

t.maxwell@ed.ac.uk

Dawei Song

Robert Gordon University
Aberdeen, United Kingdom

d.song@rgu.ac.uk

Yuexian Hou

Tianjin University
Tianjin, China

yxhou@tju.edu.cn

Peng Zhang

Robert Gordon University
Aberdeen, United Kingdom.

p.zhang1@rgu.ac.uk

Abstract

Bag-of-words approaches to information retrieval (IR) are effective but assume independence between words. The Hyperspace Analogue to Language (HAL) is a cognitively motivated and validated semantic space model that captures statistical dependencies between words by considering their co-occurrences in a surrounding window of text. HAL has been successfully applied to query expansion in IR, but has several limitations, including high processing cost and use of distributional statistics that do not exploit syntax. In this paper, we pursue two methods for incorporating syntactic-semantic information from textual ‘events’ into HAL. We build the HAL space directly from events to investigate whether processing costs can be reduced through more careful definition of word co-occurrence, and improve the quality of the pseudo-relevance feedback by applying event information as a constraint during HAL construction. Both methods significantly improve performance results in comparison with original HAL, and interpolation of HAL and relevance model expansion outperforms either method alone.

1 Introduction

Despite its intuitive appeal, the incorporation of linguistic and semantic word dependencies in IR has not been shown to significantly improve over a bigram language modeling approach (Song and Croft, 1999) that encodes word dependencies assumed from mere syntactic adjacency. Both the

dependence language model for IR (Gao et al., 2004), which incorporates linguistic relations between non-adjacent words while limiting the generation of meaningless phrases, and the Markov Random Field (MRF) model, which captures short and long range term dependencies (Metzler and Croft, 2005; Metzler and Croft, 2007), consistently outperform a unigram language modelling approach but are closely approximated by a bigram language model that uses no linguistic knowledge. Improving retrieval performance through application of semantic and syntactic information beyond proximity and co-occurrence features is a difficult task but remains a tantalising prospect.

Our approach is like that of Gao et al. (2004) in that it considers semantic-syntactically determined relationships between words at the sentence level, but allows words to have more than one role, such as predicate and argument for different events, while link grammar (Sleator and Temperley, 1991) dictates that a word can only satisfy one connector in a disjunctive set. Compared to the MRF model, our approach is unsupervised where MRFs require the training of parameters using relevance judgments that are often unavailable in practical conditions.

Other work incorporating syntactic and linguistic information into IR includes early research by (Smeaton, O’Donnell and Kellely, 1995), who employed tree structured analytics (TSAs) resembling dependency trees, the use of syntax to detect paraphrases for question answering (QA) (Lin and Pantel, 2001), and semantic role labelling in QA (Shen and Lapata, 2007).

Independent from IR, Pado and Lapata (2007) proposed a general framework for the construction of a semantic space endowed with syntactic

information. This was represented by an undirected graph, where nodes stood for words, dependency edges stood for syntactical relations, and sequences of dependency edges formed paths that were weighted for each target word. Our work is in line with Pado and Lapata (2007) in constructing a semantic space with syntactic information, but builds our space from events, states and attributions as defined linguistically by Bach (1986). We call these simply *events*, and extract them automatically from predicate-argument structures and a dependency parse. We will use this space to perform query expansion in IR, a task that aims to find additional words related to original query terms, such that an expanded query including these words better expresses the information need. To our knowledge, the notion of events has not been applied to query expansion before.

This paper will outline the original HAL algorithm which serves as our baseline, and the event extraction process. We then propose two methods to arm HAL with event information: direct construction of HAL from events (eHAL-1), and treating events as constraints on HAL construction from the corpus (eHAL-2). Evaluation will compare results using original HAL, eHAL-1 and eHAL-2 with a widely used unigram language model (LM) for IR and a state of the art query expansion method, namely the Relevance Model (RM) (Lavrenko and Croft, 2001). We also explore whether a complementary effect can be achieved by combining HAL-based dependency modelling with the unigram-based RM.

2 HAL Construction

Semantic space models aim to capture the meanings of words using co-occurrence information in a text corpus. Two examples are the Hyperspace Analogue to Language (HAL) (Lund and Burgess, 1996), in which a word is represented by a vector of other words co-occurring with it in a sliding window, and Latent Semantic Analysis (LSA) (Deerwester, Dumais, Furnas, Landauer and Harshman, 1990; Landauer, Foltz and Laham, 1998), in which a word is expressed as a vector of documents (or any other syntactical units such as sentences) containing the word. In these semantic spaces, vector-based representations facilitate measurement of similarities between words. Semantic space models have been validated through various studies and demonstrate

compatibility with human information processing. Recently, they have also been applied in IR, such as LSA for latent semantic indexing, and HAL for query expansion. For the purpose of this paper, we focus on HAL, which encodes word co-occurrence information explicitly and thus can be applied to query expansion in a straightforward way.

HAL is premised on context surrounding a word providing important information about its meaning (Harris, 1968). To be specific, an L -size sliding window moves across a large text corpus word-by-word. Any two words in the same window are treated as co-occurring with each other with a weight that is inversely proportional to their separation distance in the text. By accumulating co-occurrence information over a corpus, a word-by-word matrix is constructed, a simple illustration of which is given in Table 1. A single word is represented by a row vector and a column vector that capture the information before and after the word, respectively. In some applications, direction sensitivity is ignored to obtain a single vector representation of a word by adding corresponding row and column vectors (Bai et al., 2005).

	w_1	w_2	w_3	w_4	w_5	w_6
w_1						
w_2	5					
w_3	4	5				
w_4	3	4	5			
w_5	2	3	4	5		
w_6		2	3	4	5	

Table 1: A HAL space for the text “ $w_1 w_2 w_3 w_4 w_5 w_6$ ” using a 5-word sliding window ($L = 5$).

HAL has been successfully applied to query expansion and can be incorporated into this task directly (Bai et al., 2005) or indirectly, as with the Information Flow method based on HAL (Bruza and Song, 2002). However, to date it has used only statistical information from co-occurrence patterns. We extend HAL to incorporate syntactic-semantic information.

3 Event Extraction

Prior to event extraction, predicates, arguments, part of speech (POS) information and syntactic dependencies are annotated using the best-performing joint syntactic-semantic parser from the CoNLL 2008 Shared Task (Johansson and

Nugues, 2008), trained on PropBank and NomBank data. The event extraction algorithm then instantiates the template *REL [modREL] Arg0 [modArg0] ...ArgN [modArgN]*, where *REL* is the predicate relation (or root verb if no predicates are identified), and *Arg0...ArgN* are its arguments. Modifiers (*mod*) are identified by tracing from predicate and argument heads along the dependency tree. All predicates are associated with at least one event unless both *Arg0* and *Arg1* are not identified, or the only argument is not a noun.

The algorithm checks for modifiers based on POS tag¹, tracing up and down the dependency tree, skipping over prepositions, coordinating conjunctions and words indicating appositionment, such as ‘*sample (of)*’. However, to constrain output the search is limited to a depth of one (with the exception of skipping). For example, given the phrase ‘*apples from the store nearby*’ and an argument head *apples*, the first dependent, *store*, will be extracted but not *nearby*, which is the dependent of *store*. This can be detrimental when encountering compound nouns but does focus on core information. For verbs, modal dependents are not included in output.

Available paths up and down the dependency tree are followed until all branches are exhausted, given the rules outlined above. Tracing can result in multiple extracted events for one predicate and predicates may also appear as arguments in a different event, or be part of argument phrases. For this reason, events are constrained to cover only detail appearing above subsequent predicates in the tree, which simplifies the event structure. For example, the sentence “*Baghdad already has the facilities to continue producing massive quantities of its own biological and chemical weapons*” results in the event output: (1) *has Baghdad already facilities continue producing*; (2) *continue quantities producing massive*; (3) *producing quantities massive weapons biological*; (4) *quantities weapons biological massive*.

4 HAL With Events

4.1 eHAL-1: Construction From Events

Since events are extracted from documents, they form a reduced text corpus from which HAL can

¹To be specific, the modifiers include negation, as well as adverbs or particles for verbal heads, adjectives and nominal modifiers for nominal heads, and verbal or nominal dependents of modifiers, provided modifiers are not also identified as arguments elsewhere in the event.

be built in a similar manner to the original HAL. We ignore the parameter of window length (*L*) and treat every event as a single window of length equal to the number of words in the event. Every pair of words in an event is considered to be co-occurrent with each other. The weight assigned to the association between each pair is simply set to one. With this scheme, all the events are traversed and the event-based HAL is constructed.

The advantage of this method is that it substantially reduces the processing time during HAL construction because only events are involved and there is no need to calculate weights per occurrence. Additional processing time is incurred in semantic role labelling (SRL) during event identification. However, the naive approach to extraction might be simulated with a combination of less costly chunking and dependency parsing, given that the word ordering information available with SRL is not utilised.

eHAL-1 combines syntactical and statistical information, but has a potential drawback in that only events are used during construction so some information existing in the co-occurrence patterns of the original text may be lost. This motivates the second method.

4.2 eHAL-2: Event-Based Filtering

This method attempts to include more statistical information in eHAL construction. The key idea is to decide whether a text segment in a corpus should be used for the HAL construction, based on how much event information it covers. Given a corpus of text and the events extracted from it, the eHAL-2 method runs as follows:

1. Select the events of length *M* or more and discard the others for efficiency;
2. Set an “inclusion criterion”, which decides if a text segment, defined as a word sequence within an *L*-size sliding window, contains an event. For example, if 80% of the words in an event are contained in a text segment, it could be considered to “include” the event;
3. Move across the whole corpus word-by-word with an *L*-size sliding window. For each window, complete Steps 4-7;
4. For the current *L*-size text segment, check whether it includes an event according to the “inclusion criterion” (Step 2);

- If an event is included in the current text segment, check the following segments for a consecutive sequence of segments that also include this event. If the current segment includes more than one event, find the longest sequence of related text segments. An illustration is given in Figure 1 in which dark nodes stand for the words in a specific event and an 80% inclusion criterion is used.



Figure 1: Consecutive segments for an event

- Extract the full span of consecutive segments just identified and go to the next available text segment. Repeat Step 3;
- When the scanning is done, construct HAL using the original HAL method over all extracted sequences.

With the guidance of event information, the procedure above keeps only those segments of text that include at least one event and discards the rest. It makes use of more statistical co-occurrence information than eHAL-1 by applying weights that are proportional to word separation distance. It also alleviates the identified drawback of eHAL-1 by using the full text surrounding events. A trade-off is that not all the events are included by the selected text segments, and thus some syntactical information may be lost. In addition, the parametric complexity and computational complexity are also higher than eHAL-1.

5 Evaluation

We empirically test whether our event-based HALs perform better than the original HAL, and standard LM and RM, using three TREC² collections: AP89 with Topics 1-50 (*title* field), AP8889 with Topics 101-150 (*title* field) and WSJ9092 with Topics 201-250 (*description* field). All the collections are stemmed, and stop words are removed, prior to retrieval using the Lemur Toolkit Version 4.11³. Initial retrieval is identical for all models evaluated: KL-divergence

²TREC stands for the Text REtrieval Conference series run by NIST. Please refer to <http://trec.nist.gov/> for details.

³Available at <http://www.lemurproject.org/>

based LM smoothed using Dirichlet prior with μ set to 1000 as appropriate for TREC style title queries (Lavrenko, 2004). The top 50 returned documents form the basis for all pseudo-relevance feedback, with other parameters tuned separately for the RM and HAL methods.

For each dataset, the number of feedback terms for each method is selected optimally among 20, 40, 60, 80⁴ and the interpolation and smoothing coefficient is set to be optimal in [0,1] with interval 0.1. For RM, we choose the first relevance model in Lavrenko and Croft (2001) with the document model smoothing parameter optimally set at 0.8. The number of feedback terms is fixed at 60 (for AP89 and WSJ9092) and 80 (for AP8889), and interpolation between the query and relevance models is set at 0.7 (for WSJ9092) and 0.9 (for AP89 and AP8889). The HAL-based query expansion methods add the top 80 expansion terms to the query with interpolation coefficient 0.9 for WSJ9092 and 1 (that is, no interpolation) for AP89 and AP8889. The other HAL-based parameters are set as follows: shortest event length $M = 5$, for eHAL-2 the “inclusion criterion” is 75% of words in an event, and for HAL and eHAL-2, window size $L = 8$. Top expansion terms are selected according to the formula:

$$P_{HAL}(t_j | \oplus t) = \frac{HAL(t_j | \oplus q)}{\sum_{t_i} HAL(t_i | \oplus q)}$$

where $HAL(t_j | \oplus q)$ is the weight of t_j in the combined HAL vector $\oplus q$ (Bruza and Song, 2002) of original query terms. Mean Average Precision (MAP) is the performance indicator, and t-test (at the level of 0.05) is performed to measure the statistical significance of results.

Table 2 lists the experimental results⁵. It can be observed that all the three HAL-based query expansion methods improve performance over the LM and both eHALs achieve better performance than original HAL, indicating that the incorporation of event information is beneficial. In addition, eHAL-2 leads to better performance than eHAL-1, suggesting that use of linguistic information as a constraint on statistical processing, rather than the focus of extraction, is a more effective strategy. The results are still short of those achieved

⁴For RM, feedback terms were also tested on larger numbers up to 1000 but only comparable result was observed.

⁵In Table 2, brackets show percent improvement of eHALs / RM over HAL / eHAL-2 respectively and * and # indicate the corresponding statistical significance.

Method	AP89	AP8889	WSJ9092
LM	0.2015	0.2290	0.2242
HAL	0.2299	0.2738	0.2346
eHAL-1	0.2364 (+2.83%)	0.2829 (+3.32%*)	0.2409 (+2.69%)
eHAL-2	0.2427 (+5.57%*)	0.2850 (+4.09%*)	0.2460 (+4.86%*)
RM	0.2611 (+7.58%#)	0.3178 (+11.5%#)	0.2676 (+8.78%#)

Table 2: Performance (MAP) comparison of query expansion using different HALs

with RM, but the gap is significantly reduced by incorporating event information here, suggesting this is a promising line of work. In addition, as shown in (Bai et al., 2005), the Information Flow method built upon the original HAL largely outperformed RM. We expect that eHAL would provide an even better basis for Information Flow, but this possibility is yet to be explored.

As is known, RM is a pure unigram model while HAL methods are dependency-based. They capture different information, hence it is natural to consider if their strengths might complement each other in a combined model. For this purpose, we design the following two schemes:

1. Apply RM to the feedback documents (original RM), the events extracted from these documents (eRM-1), and the text segments around each event (eRM-2), where the three sources are the same as used to produce HAL, eHAL-1 and eHAL-2 respectively;
2. Interpolate the expanded query model by RM with the ones generated by each HAL, represented by HAL+RM, eHAL-1+RM and eHAL-2+RM. The interpolation coefficient is again selected to achieve the optimal MAP.

The MAP comparison between the original RM and these new models are demonstrated in Table 3⁶. From the first three lines (Scheme 1), we can observe that in most cases the performance generally deteriorates when RM is directly run over the events and the text segments. The event information is more effective to express the information about the term dependencies while the unigram RM ignores this information and only takes

⁶For rows in Table 3, brackets show percent difference from original RM.

Method	AP89	AP8889	WSJ9092
RM	0.2611	0.3178	0.2676
eRM-1	0.2554 (-2.18%)	0.3150 (-0.88%)	0.2555 (-4.52%)
eRM-2	0.2605 (-0.23%)	0.3167 (-0.35%)	0.2626 (-1.87%)
HAL +RM	0.2640 (+1.11%)	0.3186 (+0.25%)	0.2727 (+1.19%)
eHAL-1 +RM	0.2600 (-0.42%)	0.3210 (+1.01%)	0.2734 (+2.17%)
eHAL-2 +RM	0.2636 (+0.96%)	0.3191 (+0.41%)	0.2735 (+2.20%)

Table 3: Performance (MAP) comparison of query expansion using the combination of RM and term dependencies

the occurrence frequencies of individual words into account, which is not well-captured by the events. In contrast, the performance of Scheme 2 is more promising. The three methods outperform the original RM in most cases, but the improvement is not significant and it is also observed that there is little difference shown between RM with HAL and eHALs. The phenomenon implies more effective methods may be invented to complement the unigram models with the syntactical and statistical dependency information.

6 Conclusions

The application of original HAL to query expansion attempted to incorporate statistical word association information, but did not take into account the syntactical dependencies and had a high processing cost. By utilising syntactic-semantic knowledge from event modelling of pseudo-relevance feedback documents prior to computing the HAL space, we showed that processing costs might be reduced through more careful selection of word co-occurrences and that performance may be enhanced by effectively improving the quality of pseudo-relevance feedback documents. Both methods improved over original HAL query expansion. In addition, interpolation of HAL and RM expansion improved results over those achieved by either method alone.

Acknowledgments

This research is funded in part by the UK’s Engineering and Physical Sciences Research Council, grant number: EP/F014708/2.

References

- Bach E. The Algebra of Events. 1986. *Linguistics and Philosophy*, 9(1): pp. 5–16.
- Bai J. and Song D. and Bruza P. and Nie J.-Y. and Cao G. Query Expansion using Term Relationships in Language Models for Information Retrieval 2005. In: *Proceedings of the 14th International ACM Conference on Information and Knowledge Management*, pp. 688–695.
- Bruza P. and Song D. Inferring Query Models by Computing Information Flow. 2002. In: *Proceedings of the 11th International ACM Conference on Information and Knowledge Management*, pp. 206–269.
- Deerwester S., Dumais S., Furnas G., Landauer T. and Harshman R. Indexing by latent semantic analysis. 1990. *Journal of the American Society for Information Science*, 41(6): pp. 391–407.
- Gao J. and Nie J. and Wu G. and Cao G. Dependence Language Model for Information Retrieval. 2004. In: *Proceedings of the 27th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 170–177.
- Harris Z. 1968. *Mathematical Structures of Language*. Wiley, New York.
- Johansson R. and Nugues P. Dependency-based Syntactic-semantic Analysis with PropBank and NomBank. 2008. In: *CoNLL '08: Proceedings of the Twelfth Conference on Computational Natural Language Learning*, pp. 183–187.
- Landauer T., Foltz P. and Laham D. Introduction to Latent Semantic Analysis. 1998. *Discourse Processes*, 25: pp. 259–284.
- Lavrenko V. 2004. *A Generative Theory of Relevance*, PhD thesis, University of Massachusetts, Amherst.
- Lavrenko V. and Croft W. B. Relevance Based Language Models. 2001. In: *SIGIR '01: Proceedings of the 24th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 120–127, New York, NY, USA, 2001. ACM.
- Lin D. and Pantel P. DIRT - Discovery of Inference Rules from Text. 2001. In: *KDD '01: Proceedings of the Seventh ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 323–328, New York, NY, USA.
- Lund K. and Burgess C. Producing High-dimensional Semantic Spaces from Lexical Co-occurrence. 1996. *Behavior Research Methods, Instruments & Computers*, 28: pp. 203–208. Prentice-Hall, Englewood Cliffs, NJ.
- Metzler D. and Bruce W. B. A Markov Random Field Model for Term Dependencies 2005. In: *SIGIR '05: Proceedings of the 28th annual international ACM SIGIR conference on Research and development in information retrieval*, pp. 472–479, New York, NY, USA. ACM.
- Metzler D. and Bruce W. B. Latent Concept Expansion using Markov Random Fields 2007. In: *SIGIR '07: Proceedings of the 30th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 311–318, ACM, New York, NY, USA.
- Pado S. and Lapata M. Dependency-Based Construction of Semantic Space Models. 2007. *Computational Linguistics*, 33: pp. 161–199.
- Shen D. and Lapata M. Using Semantic Roles to Improve Question Answering. 2007. In: *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pp. 12–21.
- Sleator D. D. and Temperley D. Parsing English with a Link Grammar 1991. *Technical Report CMU-CS-91-196*, Department of Computer Science, Carnegie Mellon University.
- Smeaton A. F., O'Donnell R. and Kellely F. Indexing Structures Derived from Syntax in TREC-3: System Description. 1995. In: *The Third Text REtrieval Conference (TREC-3)*, pp. 55–67.
- Song F. and Croft W. B. A General Language Model for Information Retrieval. 1999. In: *CIKM '99: Proceedings of the Eighth International Conference on Information and Knowledge Management*, pp. 316–321, New York, NY, USA, ACM.

Automatically Generating Term-frequency-induced Taxonomies

Karin Murthy

Tanveer A Faruque

L Venkata Subramaniam

K Hima Prasad

Mukesh Mohania

IBM Research - India

{karinmur|ftanveer|lvsubram|hkaranam|mkmukesh}@in.ibm.com

Abstract

We propose a novel method to automatically acquire a term-frequency-based taxonomy from a corpus using an unsupervised method. A term-frequency-based taxonomy is useful for application domains where the frequency with which terms occur on their own and in combination with other terms imposes a natural term hierarchy. We highlight an application for our approach and demonstrate its effectiveness and robustness in extracting knowledge from real-world data.

1 Introduction

Taxonomy deduction is an important task to understand and manage information. However, building taxonomies manually for specific domains or data sources is time consuming and expensive. Techniques to automatically deduce a taxonomy in an unsupervised manner are thus indispensable. Automatic deduction of taxonomies consist of two tasks: extracting relevant terms to represent concepts of the taxonomy and discovering relationships between concepts. For unstructured text, the extraction of relevant terms relies on information extraction methods (Etzioni et al., 2005).

The relationship extraction task can be classified into two categories. Approaches in the first category use lexical-syntactic formulation to define patterns, either manually (Kozareva et al., 2008) or automatically (Girju et al., 2006), and apply those patterns to mine instances of the patterns. Though producing accurate results, these approaches usually have low coverage for many domains and suffer from the problem of inconsistency between terms when connecting the instances as chains to form a taxonomy. The second category of approaches uses clustering to discover terms and the relationships between them (Roy

and Subramaniam, 2006), even if those relationships do not explicitly appear in the text. Though these methods tackle inconsistency by addressing taxonomy deduction globally, the relationships extracted are often difficult to interpret by humans.

We show that for certain domains, the frequency with which terms appear in a corpus on their own and in conjunction with other terms induces a natural taxonomy. We formally define the concept of a term-frequency-based taxonomy and show its applicability for an example application. We present an unsupervised method to generate such a taxonomy from scratch and outline how domain-specific constraints can easily be integrated into the generation process. An advantage of the new method is that it can also be used to extend an existing taxonomy.

We evaluated our method on a large corpus of real-life addresses. For addresses from emerging geographies no standard postal address scheme exists and our objective was to produce a postal taxonomy that is useful in standardizing addresses (Kothari et al., 2010). Specifically, the experiments were designed to investigate the effectiveness of our approach on noisy terms with lots of variations. The results show that our method is able to induce a taxonomy without using any kind of lexical-semantic patterns.

2 Related Work

One approach for taxonomy deduction is to use explicit expressions (Iwaska et al., 2000) or lexical and semantic patterns such as *is a* (Snow et al., 2004), *similar usage* (Kozareva et al., 2008), *synonyms and antonyms* (Lin et al., 2003), *purpose* (Cimiano and Wenderoth, 2007), and *employed by* (Bunescu and Mooney, 2007) to extract and organize terms. The quality of extraction is often controlled using statistical measures (Pantel and Pennacchiotti, 2006) and external resources such as wordnet (Girju et al., 2006). However, there are

domains (such as the one introduced in Section 3.2) where the text does not allow the derivation of linguistic relations.

Supervised methods for taxonomy induction provide training instances with global semantic information about concepts (Fleischman and Hovy, 2002) and use bootstrapping to induce new seeds to extract further patterns (Cimiano et al., 2005). Semi-supervised approaches start with known terms belonging to a category, construct context vectors of classified terms, and associate categories to previously unclassified terms depending on the similarity of their context (Tanev and Magnini, 2006). However, providing training data and hand-crafted patterns can be tedious. Moreover in some domains (such as the one presented in Section 3.2) it is not possible to construct a context vector or determine the replacement fit.

Unsupervised methods use clustering of word-context vectors (Lin, 1998), co-occurrence (Yang and Callan, 2008), and conjunction features (Carballo, 1999) to discover implicit relationships. However, these approaches do not perform well for small corpora. Also, it is difficult to label the obtained clusters which poses challenges for evaluation. To avoid these problems, incremental clustering approaches have been proposed (Yang and Callan, 2009). Recently, lexical entailment has been used where the term is assigned to a category if its occurrence in the corpus can be replaced by the lexicalization of the category (Giuliano and Gliozzo, 2008). In our method, terms are incrementally added to the taxonomy based on their support and context.

Association rule mining (Agrawal and Srikant, 1994) discovers interesting relations between terms, based on the frequency with which terms appear together. However, the amount of patterns generated is often huge and constructing a taxonomy from all the patterns can be challenging. In our approach, we employ similar concepts but make taxonomy construction part of the relationship discovery process.

3 Term-frequency-induced Taxonomies

For some application domains, a taxonomy is induced by the frequency in which terms appear in a corpus on their own and in combination with other terms. We first introduce the problem formally and then motivate it with an example application.

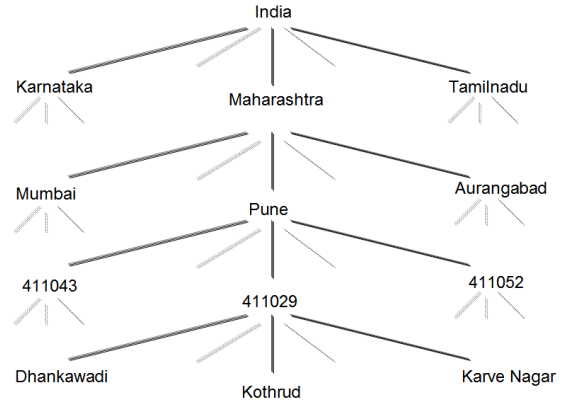


Figure 1: Part of an address taxonomy

3.1 Definition

Let C be a corpus of records r . Each record is represented as a set of terms t . Let $T = \{t \mid t \in r \wedge r \in C\}$ be the set of all terms of C . Let $f(t)$ denote the frequency of term t , that is the number of records in C that contain t . Let $F(t, T^+, T^-)$ denote the *frequency* of term t given a set of must-also-appear terms T^+ and a set of cannot-also-appear terms T^- . $F(t, T^+, T^-) = |\{r \in C \mid t \in r \wedge \forall t' \in T^+ : t' \in r \wedge \forall t' \in T^- : t' \notin r\}|$.

A term-frequency-induced taxonomy (TFIT), is an ordered tree over terms in T . For a node n in the tree, $n.t$ is the term at n , $A(n)$ the ancestors of n , and $P(n)$ the predecessors of n .

A TFIT has a root node with the special term \perp and the conditional frequency ∞ . The following condition is true for any other node n :

$$\forall t \in T, F(n.t, A(n), P(n)) \geq F(t, A(n), P(n)).$$

That is, each node's term has the highest conditional frequency in the context of the node's ancestors and predecessors. Only terms with a conditional frequency above zero are added to a TFIT.

We show in Section 4 how a TFIT taxonomy can be automatically induced from a given corpus. But before that, we show that TFITs are useful in practice and reflect a natural ordering of terms for application domains where the concept hierarchy is expressed through the frequency in which terms appear.

3.2 Example Domain: Address Data

An address taxonomy is a key enabler for address standardization. Figure 1 shows part of such an address taxonomy where the root contains the most generic term and leaf-level nodes contain the most specific terms. For emerging economies building a standardized address taxonomy is a huge chal-

Row	Term	Part of address	Category
1	D-15	house number	alphanumeric
2	Rawal	building name	proper noun
3	Complex	building name	proper noun
4	Behind	landmark	marker
5	Hotel	landmark	marker
6	Ruchira	landmark	proper noun
7	Katre	street	proper noun
8	Road	street	marker
9	Jeevan	area	proper noun
10	Nagar	area	marker
11	Andheri	city (taluk)	proper noun
12	East	city (taluk)	direction
13	Mumbai	district	proper noun
14	Maharashtra	state	proper noun
15	400069	ZIP code	6 digit string

Table 1: Example of a tokenized address

lenge. First, new areas and with it new addresses constantly emerge. Second, there are very limited conventions for specifying an address (Faruque et al., 2010). However, while many developing countries do not have a postal taxonomy, there is often no lack of address data to learn a taxonomy from.

Column 2 of Table 1 shows an example of an Indian address. Although Indian addresses tend to follow the general principal that more specific information is mentioned earlier, there is no fixed order for different elements of an address. For example, the ZIP code of an address may be mentioned before or after the state information and, although ZIP code information is more specific than city information, it is generally mentioned later in the address. Also, while ZIP codes often exist, their use by people is very limited. Instead, people tend to mention copious amounts of landmark information (see for example rows 4-6 in Table 1).

Taking all this into account, there is often not enough structure available to automatically infer a taxonomy purely based on the structural or semantic aspects of an address. However, for address data, the general-to-specific concept hierarchy is reflected in the frequency with which terms appear on their own and together with other terms.

It mostly holds that $f(s) > f(d) > f(c) > f(z)$ where s is a state name, d is a district name, c is a city name, and z is a ZIP code. However, sometimes the name of a large city may be more frequent than the name of a small state. For example, in a given corpus, the term 'Houston' (a populous US city) may appear more frequent than the term 'Vermont' (a small US state). To avoid that 'Houston' is picked as a node at the first level of the taxonomy (which should only contain

states), the conditional-frequency constraint introduced in Section 3.1 is enforced for each node in a TFIT. 'Houston's state 'Texas' (which is more frequent) is picked before 'Houston'. After 'Texas' is picked it appears in the "cannot-also-appear" list for all further siblings on the first level, thus giving 'Houston' has a conditional frequency of zero.

We show in Section 5 that an address taxonomy can be inferred by generating a TFIT taxonomy.

4 Automatically Generating TFITs

We describe a basic algorithm to generate a TFIT and then show extensions to adapt to different application domains.

4.1 Base Algorithm

Algorithm 1 Algorithm for generating a TFIT.

```

// For initialization  $T^+, T^-$  are empty
// For initialization  $l, w$  are zero
genTFIT( $T^+, T^-, C, l, w$ )
// select most frequent term
 $t_{next} = t_j$  with  $F(t_j, T^+, T^-)$  is maximal amongst all
 $t_j \in C$ ;
 $f_{next} = F(t_{next}, T^+, T^-)$ ;
if  $f_{next} \geq \text{support}$  then
  //Output node  $(t_j, l, w)$ 
  ...
  // Generate child node
  genTFIT( $T^+ \cup \{t_{next}\}, T^-, C, l + 1, w$ )
  // Generate sibling node
  genTFIT( $T^+, T^- \cup \{t_{next}\}, C, l, w + 1$ )
end if

```

To generate a TFIT taxonomy as defined in Section 3.1 we recursively pick the most frequent term given previously chosen terms. The basic algorithm *genTFIT* is sketched out in Algorithm 1. When *genTFIT* is called the first time, T^+ and T^- are empty and both level l and width w are zero. With each call of *genTFIT* a new node n in the taxonomy is created with (t, l, w) where t is the most frequent term given T^+ and T^- and l and w capture the position in the taxonomy. *genTFIT* is recursively called to generate a child of n and a sibling for n .

The only input parameter required by our algorithm is *support*. Instead of adding all terms with a conditional frequency above zero, we only add terms with a conditional frequency equal to or higher than *support*. The *support* parameter controls the precision of the resulting TFIT and also the runtime of the algorithm. Increasing *support* increases the precision but also lowers the recall.

4.2 Integrating Constraints

Structural as well as semantic constraints can easily be integrated into the TFIT generation.

We distinguish between taxonomy-level and node-level structural constraints. For example, limiting the depth of the taxonomy by introducing a *maxLevel* constraint and checking before each recursive call if *maxLevel* is reached, is a taxonomy-level constraint. A node-level constraint applies to each node and affects the way the frequency of terms is determined.

For our example application, we introduce the following node-level constraint: at each node we only count terms that appear at specific positions in records with respect to the current level of the node. Specifically, we slide (or incrementally increase) a window over the address records starting from the end. For example, when picking the term 'Washington' as a state name, occurrences of 'Washington' as city or street name are ignored. Using a window instead of an exact position accounts for positional variability. Also, to accommodate varying amounts of landmark information we length-normalize the position of terms. That is, we divide all positions in an address by the average length of an address (which is 10 for our 40 Million addresses). Accordingly, we adjust the size of the window and use increments of 0.1 for sliding (or increasing) the window.

In addition to syntactical constraints, semantic constraints can be integrated by classifying terms for use when picking the next frequent term. In our example application, markers tend to appear much more often than any proper noun. For example, the term 'Road' appears in almost all addresses, and might be picked up as the most frequent term very early in the process. Thus, it is beneficial to ignore marker terms during taxonomy generation and adding them as a post-processing step.

4.3 Handling Noise

The approach we propose naturally handles noise by ignoring it, unless the noise level exceeds the support threshold. Misspelled terms are generally infrequent and will as such not become part of the taxonomy. The same applies to incorrect addresses. Incomplete addresses partially contribute to the taxonomy and only cause a problem if the same information is missing too often. For example, if more than *support* addresses with the city 'Houston' are missing the state 'Texas', then

'Houston' may become a node at the first level and appear to be a state. Generally, such cases only appear at the far right of the taxonomy.

5 Evaluation

We present an evaluation of our approach for address data from an emerging economy. We implemented our algorithm in Java and store the records in a DB2 database. We rely on the DB2 optimizer to efficiently retrieve the next frequent term.

5.1 Dataset

The results are based on 40 Million Indian addresses. Each address record was given to us as a single string and was first tokenized into a sequence of terms as shown in Table 1. In a second step, we addressed spelling variations. There is no fixed way of transliterating Indian alphabets to English and most Indian proper nouns have various spellings in English. We used tools to detect synonyms with the same context to generate a list of rules to map terms to a standard form (Lin, 1998). For example, in Table 1 'Maharashtra' can also be spelled 'Maharastra'. We also used a list of keywords to classify some terms as markers such as 'Road' and 'Nagar' shown in Table 1.

Our evaluation consists of two parts. First, we show results for constructing a TFIT from scratch. To evaluate the precision and recall we also retrieved post office addresses from India Post¹, cleaned them, and organized them in a tree.

Second, we use our approach to enrich the existing hierarchy created from post office addresses with additional area terms. To validate the result, we also retrieved data about which area names appear within a ZIP code.² We also verified whether Google Maps shows an area on its map.³

5.2 Taxonomy Generation

We generated a taxonomy O using all 40 million addresses. We compare the terms assigned to category levels *district* and *taluk*⁴ in O with the tree P constructed from post office addresses. Each district and taluk has at least one post office. Thus P covers all districts and taluks and allows us to test coverage and precision. We compute the precision and recall for each category level CL as

¹<http://www.indiapost.gov.in/Pin/pinsearch.aspx>

²<http://www.whereincity.com/india/pincode/search>

³maps.google.com

⁴Administrative division in some South-Asian countries.

Support		Recall %	Precision %
100	District	93.9	57.4
	Taluk	50.9	60.5
200	District	87.9	64.4
	Taluk	49.6	66.1

Table 2: Precision and recall for categorizing terms belonging to the state Maharashtra

$$Recall_{CL} = \frac{\# \text{ correct paths from root to } CL \text{ in } O}{\# \text{ paths from root to } CL \text{ in } P}$$

$$Precision_{CL} = \frac{\# \text{ correct paths from root to } CL \text{ in } O}{\# \text{ paths from root to } CL \text{ in } O}$$

Table 2 shows precision and recall for district and taluk for the large state Maharashtra. Recall is good for district. For taluk it is lower because a major part of the data belongs to urban areas where taluk information is missing. The precision seems to be low but it has to be noted that in almost 75% of the addresses either district or taluk information is missing or noisy. Given that, we were able to recover a significant portion of the knowledge structure.

We also examined a branch for a smaller state (Kerala). Again, both districts and taluks appear at the next level of the taxonomy. For a support of 200 there are 19 entries in O of which all but two appear in P as district or taluk. One entry is a taluk that actually belongs to Maharashtra and one entry is a name variation of a taluk in P . There were not enough addresses to get a good coverage of all districts and taluks.

5.3 Taxonomy Augmentation

We used P and ran our algorithm for each branch in P to include area information. We focus our evaluation on the city Mumbai. The recall is low because many addresses do not mention a ZIP code or use an incorrect ZIP code. However, the precision is good implying that our approach works even in the presence of large amounts of noise.

Table 3 shows the results for ZIP code 400002 and 400004 for a support of 100. We get similar results for other ZIP codes. For each detected area we compared whether the area is also listed on whereincity.com, part of a post office name (PO), or shown on google maps. All but four areas found are confirmed by at least one of the three external sources. Out of the unconfirmed terms *Fanaswadi* and *MarineDrive* seem to be genuine area names but we could not confirm *DhakurdwarRoad*. The term *th* is due to our

Area	Whereincity	PO	Google
Bhuleshwar	yes	no	yes
Chira Bazar	yes	no	yes
Dhobi Talao	no	no	yes
Fanaswadi	no	no	no
Kalbadevi Road	yes	yes	yes
Marine Drive	no	no	no
Marine Lines	yes	yes	yes
Princess Street	no	no	yes
th	no	no	no
Thakurdwar Road	no	no	no
Zaveri Bazar	yes	no	yes
Charni Road	no	yes	no
Girgaon	yes	yes	yes
Khadilkar Road	yes	no	yes
Khetwadi Road	yes	no	no
Kumbharwada	no	no	yes
Opera House	no	yes	no
Prathna Samaj	yes	no	no

Table 3: Areas found for ZIP code 400002 (top) and 400004 (bottom)

tokenization process. 16 correct terms out of 18 terms results in a precision of 89%.

We also ran experiments to measure the coverage of area detection for Mumbai without using ZIP codes. Initializing our algorithm with *Maharashtra* and *Mumbai* yielded over 100 areas with a support of 300 and more. However, again the precision is low because quite a few of those areas are actually taluk names.

Using a large number of addresses is necessary to achieve good recall and precision.

6 Conclusion

In this paper, we presented a novel approach to generate a taxonomy for data where terms exhibit an inherent frequency-based hierarchy. We showed that term frequency can be used to generate a meaningful taxonomy from address records. The presented approach can also be used to extend an existing taxonomy which is a big advantage for emerging countries where geographical areas evolve continuously.

While we have evaluated our approach on address data, it is applicable to all data sources where the inherent hierarchical structure is encoded in the frequency with which terms appear on their own and together with other terms. Preliminary experiments on real-time analyst’s stock market tips⁵ produced a taxonomy of (TV station, Analyst, Affiliation) with decent precision and recall.

⁵See Live Market voices at: http://money.rediff.com/money/jsp/markets_home.jsp

References

- Rakesh Agrawal and Ramakrishnan Srikant. 1994. Fast algorithms for mining association rules in large databases. In *Proceedings of the 20th International Conference on Very Large Data Bases*, pages 487–499.
- Razvan C. Bunescu and Raymond J. Mooney. 2007. Learning to extract relations from the web using minimal supervision. In *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, pages 576–583.
- Sharon A. Caraballo. 1999. Automatic construction of a hypernym-labeled noun hierarchy from text. In *Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics on Computational Linguistics*, pages 120–126.
- Philipp Cimiano and Johanna Wenderoth. 2007. Automatic acquisition of ranked qualia structures from the web. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics*, pages 888–895.
- Philipp Cimiano, Günter Ladwig, and Steffen Staab. 2005. Gimme’ the context: context-driven automatic semantic annotation with c-pankow. In *Proceedings of the 14th International Conference on World Wide Web*, pages 332–341.
- Oren Etzioni, Michael Cafarella, Doug Downey, Ana-Maria Popescu, Tal Shaked, Stephen Soderland, Daniel S. Weld, and Alexander Yates. 2005. Unsupervised named-entity extraction from the web: an experimental study. *Artificial Intelligence*, 165(1):91–134.
- Tanveer A. Faruque, K. Hima Prasad, L. Venkata Subramaniam, Mukesh K. Mohania, Girish Venkatachaliah, Shrinivas Kulkarni, and Prमित Basu. 2010. Data cleansing as a transient service. In *Proceedings of the 26th International Conference on Data Engineering*, pages 1025–1036.
- Michael Fleischman and Eduard Hovy. 2002. Fine grained classification of named entities. In *Proceedings of the 19th International Conference on Computational Linguistics*, pages 1–7.
- Roxana Girju, Adriana Badulescu, and Dan Moldovan. 2006. Automatic discovery of part-whole relations. *Computational Linguistics*, 32(1):83–135.
- Claudio Giuliano and Alfio Gliozzo. 2008. Instance-based ontology population exploiting named-entity substitution. In *Proceedings of the 22nd International Conference on Computational Linguistics*, pages 265–272.
- Lucja M. Iwaska, Naveen Mata, and Kellyn Kruger. 2000. Fully automatic acquisition of taxonomic knowledge from large corpora of texts. In Lucja M. Iwaska and Stuart C. Shapiro, editors, *Natural Language Processing and Knowledge Representation: Language for Knowledge and Knowledge for Language*, pages 335–345.
- Govind Kothari, Tanveer A Faruque, L V Subramaniam, K H Prasad, and Mukesh Mohania. 2010. Transfer of supervision for improved address standardization. In *Proceedings of the 20th International Conference on Pattern Recognition*.
- Zornitsa Kozareva, Ellen Riloff, and Eduard Hovy. 2008. Semantic class learning from the web with hyponym pattern linkage graphs. In *Proceedings of the 46th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 1048–1056.
- Dekang Lin, Shaojun Zhao, Lijuan Qin, and Ming Zhou. 2003. Identifying synonyms among distributionally similar words. In *Proceedings of the 18th International Joint Conference on Artificial Intelligence*, pages 1492–1493.
- Dekang Lin. 1998. Automatic retrieval and clustering of similar words. In *Proceedings of the 17th International Conference on Computational Linguistics*, pages 768–774.
- Patrick Pantel and Marco Pennacchiotti. 2006. Espresso: leveraging generic patterns for automatically harvesting semantic relations. In *Proceedings of the 21st International Conference on Computational Linguistics and the 44th Annual Meeting of the Association for Computational Linguistics*, pages 113–120.
- Shourya Roy and L Venkata Subramaniam. 2006. Automatic generation of domain models for call centers from noisy transcriptions. In *Proceedings of the 21st International Conference on Computational Linguistics and the 44th Annual Meeting of the Association for Computational Linguistics*, pages 737–744.
- Rion Snow, Daniel Jurafsky, and Andrew Y. Ng. 2004. Learning syntactic patterns for automatic hypernym discovery. In *Advances in Neural Information Processing Systems*, pages 1297–1304.
- Hristo Tanev and Bernardo Magnini. 2006. Weakly supervised approaches for ontology population. In *Proceedings of the 11th Conference of the European Chapter of the Association for Computational Linguistics*, pages 3–7.
- Hui Yang and Jamie Callan. 2008. Learning the distance metric in a personal ontology. In *Proceeding of the 2nd International Workshop on Ontologies and Information Systems for the Semantic Web*, pages 17–24.
- Hui Yang and Jamie Callan. 2009. A metric-based framework for automatic taxonomy induction. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, pages 271–279.

Complexity assumptions in ontology verbalisation

Richard Power

Department of Computing
Open University, UK
r.power@open.ac.uk

Abstract

We describe the strategy currently pursued for verbalising OWL ontologies by sentences in Controlled Natural Language (i.e., combining *generic* rules for realising logical patterns with *ontology-specific* lexicons for realising atomic terms for individuals, classes, and properties) and argue that its success depends on assumptions about the complexity of terms and axioms in the ontology. We then show, through analysis of a corpus of ontologies, that although these assumptions could in principle be violated, they are overwhelmingly respected in practice by ontology developers.

1 Introduction

Since OWL (Web Ontology Language) was adopted as a standard in 2004, researchers have sought ways of mediating between the (decidedly cumbersome) raw code and the human users who aspire to view or edit it. Among the solutions that have been proposed are more readable coding formats such as Manchester OWL Syntax (Horridge et al., 2006), and graphical interfaces such as Protégé (Knublauch et al., 2004); more speculatively, several research groups have explored ways of mapping between OWL and controlled English, with the aim of presenting ontologies (both for viewing and editing) in natural language (Schwitter and Tilbrook, 2004; Sun and Mellish, 2006; Kaljurand and Fuchs, 2007; Hart et al., 2008). In this paper we uncover and test some assumptions on which this latter approach is based.

Historically, ontology verbalisation evolved from a more general tradition (predating OWL and the Semantic Web) that aimed to support knowledge formation by automatic interpretation of texts authored in Controlled Natural Languages

(Fuchs and Schwitter, 1995). The idea is to establish a mapping from a formal language to a natural subset of English, so that any sentence conforming to the Controlled Natural Language (CNL) can be assigned a single interpretation in the formal language — and conversely, any well-formed statement in the formal language can be realised in the CNL. With the advent of OWL, some of these CNLs were rapidly adapted to the new opportunity: part of *Attempto Controlled English* (ACE) was mapped to OWL (Kaljurand and Fuchs, 2007), and *Processable English* (PENG) evolved to *Sydney OWL Syntax* (SOS) (Cregan et al., 2007). In addition, new CNLs were developed specifically for editing OWL ontologies, such as *Rabbit* (Hart et al., 2008) and *Controlled Language for Ontology Editing* (CLOnE) (Funk et al., 2007).

In detail, these CNLs display some variations: thus an inclusion relationship between the classes `Admiral` and `Sailor` would be expressed by the pattern ‘Admirals are a type of sailor’ in CLOnE, ‘Every admiral is a kind of sailor’ in Rabbit, and ‘Every admiral is a sailor’ in ACE and SOS. However, at the level of general strategy, all the CNLs rely on the same set of assumptions concerning the mapping from natural to formal language; for convenience we will refer to these assumptions as the *consensus model*. In brief, the consensus model assumes that when an ontology is verbalised in natural language, axioms are expressed by sentences, and atomic terms are expressed by entries from the lexicon. Such a model may fail in two ways: (1) an ontology might contain axioms that cannot be described transparently by a sentence (for instance, because they contain complex Boolean expressions that lead to structural ambiguity); (2) it might contain atomic terms for which no suitable lexical entry can be found. In the remainder of this paper we first describe the consensus model in more detail, then show that although

Logic	OWL
$C \sqcap D$	IntersectionOf(C D)
$\exists P.C$	SomeValuesFrom(P C)
$C \sqsubseteq D$	SubClassOf(C D)
$a \in C$	ClassAssertion(C a)
$[a, b] \in P$	PropertyAssertion(P a b)

Table 1: Common OWL expressions

in principle it is vulnerable to both the problems just mentioned, in practice these problems almost never arise.

2 Consensus model

Atomic terms in OWL (or any other language implementing description logic) are principally of three kinds, denoting either individuals, classes or properties¹. Individuals denote entities in the domain, such as Horatio Nelson or the Battle of Trafalgar; classes denote sets of entities, such as people or battles; and properties denote relations between individuals, such as the relation *victor of* between a person and a battle.

From these basic terms, a wide range of complex expressions may be constructed for classes, properties and axioms, of which some common examples are shown in table 1. The upper part of the table presents two class constructors (C and D denote any classes; P denotes any property); by combining them we could build the following expression denoting the class of persons that command fleets²:

$$Person \sqcap \exists CommanderOf.Fleet$$

The lower half of the table presents three axiom patterns for making statements about classes and individuals (a, b denote individuals); examples of their usage are as follows:

1. $Admiral \sqsubseteq \exists CommanderOf.Fleet$
2. $Nelson \in Admiral$
3. $[Nelson, Trafalgar] \in VictorOf$

Note that since class expressions contain classes as constituents, they can become indefinitely complex. For instance, given the intersection $A \sqcap B$

¹If data properties are used, there will also be terms for data types and literals (e.g., numbers and strings), but for simplicity these are not considered here.

²In description logic notation, the constructor $C \sqcap D$ forms the intersection of two classes and corresponds to Boolean conjunction, while the existential restriction $\exists P.C$ forms the class of individuals having the relation P to one or more members of class C . Thus $Person \sqcap \exists CommanderOf.Fleet$ denotes the set of individuals x such that x is a person and x commands one or more fleets.

we could replace atomic class A by a constructed class, thus obtaining perhaps $(A_1 \sqcap A_2) \sqcap B$, and so on *ad infinitum*. Moreover, since most axiom patterns contain classes as constituents, they too can become indefinitely complex.

This sketch of knowledge representation in OWL illustrates the central distinction between logical functors (e.g., *IntersectionOf*, *SubClassOf*), which belong to the W3C standard (Motik et al., 2010), and atomic terms for individuals, classes and properties (e.g., *Nelson*, *Admiral*, *VictorOf*). Perhaps the fundamental design decision of the Semantic Web is that *all domain terms remain unstandardised*, leaving ontology developers free to conceptualise the domain in any way they see fit. In the consensus verbalisation model, this distinction is reflected by dividing linguistic resources into a *generic* grammar for realising logical patterns, and an *ontology-specific* lexicon for realising atomic terms.

Consider for instance $C \sqsubseteq D$, the axiom pattern for class inclusion. This purely logical pattern can often be mapped (following ACE and SOS) to the sentence pattern ‘Every [C] is a [D]’, where C and D will be realised by count nouns from the lexicon if they are atomic, or further grammatical rules if they are complex. The more specific pattern $C \sqsubseteq \exists P.D$ can be expressed better by a sentence pattern based on a verb frame (‘Every [C] [P]s a [D]’). All these mappings depend entirely on the OWL logical functors, and will work with any lexicalisation of atomic terms that respects the syntactic constraints of the grammar, to yield verbalisations such as the following (for axioms 1-3 above):

1. Every admiral commands a fleet.
2. Nelson is an admiral.
3. Nelson is the victor of Trafalgar.

The CNLs we have cited are more sophisticated than this, allowing a wider range of linguistic patterns (e.g., adjectives for classes), but the basic assumptions are the same. The model provides satisfactory verbalisations for the simple examples considered so far, but what happens when the axioms and atomic terms become more complex?

3 Complex terms and axioms

The distribution of content among axioms depends to some extent on stylistic decisions by ontology developers, in particular with regard to ax-

iom size. This freedom is possible because description logics (including OWL) allow equivalent formulations using a large number of short axioms at one extreme, and a small number of long ones at the other. For many logical patterns, rules can be stated for amalgamating or splitting axioms while leaving overall content unchanged (thus ensuring that exactly the same inferences are drawn by a reasoning engine); such rules are often used in reasoning algorithms. For instance, any set of `SubClassOf` axioms can be amalgamated into a single ‘metaconstraint’ (Horrocks, 1997) of the form $\top \sqsubseteq M$, where \top is the class containing all individuals in the domain, and M is a class to which any individual respecting the axiom set must belong³. Applying this transformation even to only two axioms (verbalised by 1 and 2 below) will yield an outcome (verbalised by 3) that strains human comprehension:

1. Every admiral is a sailor.
2. Every admiral commands a fleet.
3. Everything is (a) either a non-admiral or a sailor, and (b) either a non-admiral or something that commands a fleet.

An example of axiom-splitting rules is found in a computational complexity proof for the description logic $\mathcal{EL}+$ (Baader et al., 2005), which requires class inclusion axioms to be rewritten to a maximally simple ‘normal form’ permitting only four patterns: $A_1 \sqsubseteq A_2$, $A_1 \sqcap A_2 \sqsubseteq A_3$, $A_1 \sqsubseteq \exists P.A_2$, and $\exists P.A_1 \sqsubseteq A_2$, where P and all A_N are atomic terms. However, this simplification of axiom structure can be achieved only by introducing new atomic terms. For example, to simplify an axiom of the form $A_1 \sqsubseteq \exists P.(A_2 \sqcap A_3)$, the rewriting rules must introduce a new term $A_{23} \equiv A_2 \sqcap A_3$, through which the axiom may be rewritten as $A_1 \sqsubseteq \exists P.A_{23}$ (along with some further axioms expressing the definition of A_{23}); depending on the expressions that they replace, the content of such terms may become indefinitely complex.

A trade-off therefore results. We can often find rules for refactoring an overcomplex axiom by a number of simpler ones, but only at the cost of introducing atomic terms for which no satisfactory lexical realisation may exist. In principle, therefore, there is no guarantee that OWL ontologies

³For an axiom set $C_1 \sqsubseteq D_1, C_2 \sqsubseteq D_2 \dots, M$ will be $(\neg C_1 \sqcup D_1) \sqcap (\neg C_2 \sqcup D_2) \dots$, where the class constructors $\neg C$ (complement of C) and $C \sqcup D$ (union of C and D) correspond to Boolean negation and disjunction.

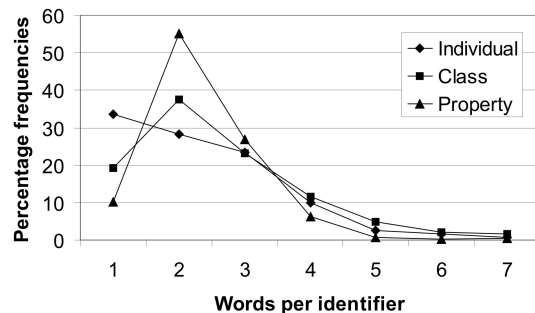


Figure 1: Identifier content

can be verbalised transparently within the assumptions of the consensus model.

4 Empirical studies of usage

We have shown that OWL syntax will permit atomic terms that cannot be lexicalised, and axioms that cannot be expressed clearly in a sentence. However, it remains possible that in practice, ontology developers use OWL in a constrained manner that favours verbalisation by the consensus model. This could happen either because the relevant constraints are psychologically intuitive to developers, or because they are somehow built into the editing tools that they use (e.g., Protégé). To investigate this possibility, we have carried out an exploratory study using a corpus of 48 ontologies mostly downloaded from the University of Manchester TONES repository (TONES, 2010). The corpus covers ontologies of varying expressivity and subject-matter, including some well-known tutorial examples (pets, pizzas) and topics of general interest (photography, travel, heraldry, wine), as well as some highly technical scientific material (mosquito anatomy, worm ontogeny, periodic table). Overall, our sample contains around 45,000 axioms and 25,000 atomic terms.

Our first analysis concerns identifier length, which we measure simply by counting the number of words in the identifying phrase. The program recovers the phrase by the following steps: (1) read an identifier (or label if one is provided⁴); (2) strip off the namespace prefix; (3) segment the resulting string into words. For the third step we

⁴Some ontology developers use ‘non-semantic’ identifiers such as #000123, in which case the meaning of the identifier is indicated in an annotation assertion linking the identifier to a label.

Pattern	Frequency	Percentage
$C_A \sqsubseteq C_A$	18961	42.3%
$C_A \sqcap C_A \sqsubseteq \perp$	8225	18.3%
$C_A \sqsubseteq \exists P_A.C_A$	6211	13.9%
$[I, I] \in P_A$	4383	9.8%
$[I, L] \in D_A$	1851	4.1%
$I \in C_A$	1786	4.0%
$C_A \equiv C_A \sqcap \exists P_A.C_A$	500	1.1%
Other	2869	6.4%
Total	44786	100%

Table 2: Axiom pattern frequencies

assume that word boundaries are marked either by underline characters or by capital letters (e.g., `battle_of.trafalgar`, `BattleOfTrafalgar`), a rule that holds (in our corpus) almost without exception. The analysis (figure 1) reveals that phrase lengths are typically between one and four words (this was true of over 95% of individuals, over 90% of classes, and over 98% of properties), as in the following random selections:

Individuals: beaujolais region, beringer, blue mountains, bondi beach

Classes: abi graph plot, amps block format, abat-toir, abbey church

Properties: has activity, has address, has amino acid, has aunt in law

Our second analysis concerns axiom patterns, which we obtain by replacing all atomic terms with a symbol meaning either individual, class, property, datatype or literal. Thus for example the axioms $Admiral \sqsubseteq Sailor$ and $Dog \sqsubseteq Animal$ are both reduced to the form $C_A \sqsubseteq C_A$, where the symbol C_A means ‘any atomic class term’. In this way we can count the frequencies of all the logical patterns in the corpus, abstracting from the domain-specific identifier names. The results (table 2) show an overwhelming focus on a small number of simple logical patterns⁵. Concerning class constructors, the most common by far were intersection ($C \sqcap C$) and existential restriction ($\exists P.C$); universal restriction ($\forall P.C$) was relatively rare, so that for example the pattern $C_A \sqsubseteq \forall P_A.C_A$ occurred only 54 times (0.1%)⁶.

⁵Most of these patterns have been explained already; the others are disjoint classes ($C_A \sqcap C_A \sqsubseteq \perp$), equivalent classes ($C_A \equiv C_A \sqcap \exists P_A.C_A$) and data property assertion ($[I, L] \in D_A$). In the latter pattern, D_A denotes a data property, which differs from an object property (P_A) in that it ranges over literals (L) rather than individuals (I).

⁶If $C \sqsubseteq \exists P.D$ means ‘Every admiral commands a fleet’, $C \sqsubseteq \forall P.D$ will mean ‘Every admiral commands only fleets’ (this will remain true if some admirals do not command anything at all).

The preference for simple patterns was confirmed by an analysis of argument structure for the OWL functors (e.g., `SubClassOf`, `IntersectionOf`) that take classes as arguments. Overall, 85% of arguments were atomic terms rather than complex class expressions. Interestingly, there was also a clear effect of argument *position*, with the first argument of a functor being atomic rather than complex in as many as 99.4% of cases⁷.

5 Discussion

Our results indicate that although in principle the consensus model cannot guarantee transparent realisations, in practice these are almost always attainable, since ontology developers overwhelmingly favour terms and axioms with relatively simple content. In an analysis of around 50 ontologies we have found that over 90% of axioms fit a mere seven patterns (table 2); the following examples show that each of these patterns can be verbalised by a clear unambiguous sentence – provided, of course, that no problems arise in lexicalising the atomic terms:

1. Every admiral is a sailor
2. No sailor is a landlubber
3. Every admiral commands a fleet
4. Nelson is the victor of Trafalgar
5. Trafalgar is dated 1805
6. Nelson is an admiral
7. An admiral is defined as a person that commands a fleet

However, since identifiers containing 3-4 words are fairly common (figure 1), we need to consider whether these formulations will remain transparent when combined with more complex lexical entries. For instance, a travel ontology in our corpus contains an axiom (fitting pattern 4) which our prototype verbalises as follows:

- 4’. West Yorkshire has as boundary the West Yorkshire Greater Manchester Boundary Fragment

The lexical entries here are far from ideal: ‘has as boundary’ is clumsy, and ‘the West Yorkshire Greater Manchester Boundary Fragment’ has as

⁷One explanation for this result could be that developers (or development tools) treat axioms as having a topic-comment structure, where the topic is usually the first argument; we intend to investigate this possibility in a further study.

many as six content words (and would benefit from hyphens). We assess the sentence as ugly but understandable, but to draw more definite conclusions one would need to perform a different kind of empirical study using human readers.

6 Conclusion

We conclude (a) that existing ontologies can be mostly verbalised using the consensus model, and (b) that an editing tool based on relatively simple linguistic patterns would not inconvenience ontology developers, but merely enforce constraints that they almost always respect anyway. These conclusions are based on analysis of identifier and axiom patterns in a corpus of ontologies; they need to be complemented by studies showing that the resulting verbalisations are understood by ontology developers and other users.

Acknowledgments

The research described in this paper was undertaken as part of the SWAT project (Semantic Web Authoring Tool), which is supported by the UK Engineering and Physical Sciences Research Council (EPSRC) grants G033579/1 (Open University) and G032459/1 (University of Manchester). Thanks are due to the anonymous ACL reviewers and to colleagues on the SWAT project for their comments and suggestions.

References

- F. Baader, I. R. Horrocks, and U. Sattler. 2005. Description logics as ontology languages for the semantic web. *Lecture Notes in Artificial Intelligence*, 2605:228–248.
- Anne Cregan, Rolf Schwitter, and Thomas Meyer. 2007. Sydney OWL Syntax - towards a Controlled Natural Language Syntax for OWL 1.1. In *OWLED*.
- Norbert Fuchs and Rolf Schwitter. 1995. Specifying logic programs in controlled natural language. In *CLNLP-95*.
- Adam Funk, Valentin Tablan, Kalina Bontcheva, Hamish Cunningham, Brian Davis, and Siegfried Handschuh. 2007. CLOnE: Controlled Language for Ontology Editing. In *6th International and 2nd Asian Semantic Web Conference (ISWC2007+ASWC2007)*, pages 141–154, November.
- Glen Hart, Martina Johnson, and Catherine Dolbear. 2008. Rabbit: Developing a control natural language for authoring ontologies. In *ESWC*, pages 348–360.
- Matthew Horridge, Nicholas Drummond, John Goodwin, Alan Rector, Robert Stevens, and Hai Wang. 2006. The Manchester OWL syntax. In *OWL: Experiences and Directions (OWLED'06)*, Athens, Georgia. CEUR.
- Ian Horrocks. 1997. *Optimising Tableaux Decision Procedures for Description Logics*. Ph.D. thesis, University of Manchester.
- K. Kaljurand and N. Fuchs. 2007. Verbalizing OWL in Attempto Controlled English. In *Proceedings of OWL: Experiences and Directions*, Innsbruck, Austria.
- Holger Knublauch, Ray W. Ferguson, Natalya Fridman Noy, and Mark A. Musen. 2004. The Protégé OWL Plugin: An Open Development Environment for Semantic Web Applications. In *International Semantic Web Conference*, pages 229–243.
- Boris Motik, Peter F. Patel-Schneider, and Bijan Parsia. 2010. OWL 2 web ontology language: Structural specification and functional-style syntax. <http://www.w3.org/TR/owl2-syntax/>. 21st April 2010.
- R. Schwitter and M. Tilbrook. 2004. Controlled natural language meets the semantic web. In *Proceedings of the Australasian Language Technology Workshop*, pages 55–62, Macquarie University.
- X. Sun and C. Mellish. 2006. Domain Independent Sentence Generation from RDF Representations for the Semantic Web. In *Proceedings of the Combined Workshop on Language-Enabled Educational Technology and Development and Evaluation of Robust Spoken Dialogue Systems (ECAI06)*, Riva del Garda, Italy.
- TONES. 2010. The TONES ontology repository. <http://owl.cs.manchester.ac.uk/repository/browser>. Last accessed: 21st April 2010.

Word Alignment with Synonym Regularization

Hiroyuki Shindo, Akinori Fujino, and Masaaki Nagata

NTT Communication Science Laboratories, NTT Corp.
2-4 Hikaridai Seika-cho Soraku-gun Kyoto 619-0237 Japan
{shindo, a.fujino}@cslab.kecl.ntt.co.jp
nagata.masaaki@lab.ntt.co.jp

Abstract

We present a novel framework for word alignment that incorporates synonym knowledge collected from monolingual linguistic resources in a bilingual probabilistic model. Synonym information is helpful for word alignment because we can expect a synonym to correspond to the same word in a different language. We design a generative model for word alignment that uses synonym information as a regularization term. The experimental results show that our proposed method significantly improves word alignment quality.

1 Introduction

Word alignment is an essential step in most phrase and syntax based statistical machine translation (SMT). It is an inference problem of word correspondences between different languages given parallel sentence pairs. Accurate word alignment can induce high quality phrase detection and translation probability, which leads to a significant improvement in SMT performance. Many word alignment approaches based on generative models have been proposed and they learn from bilingual sentences in an unsupervised manner (Vogel et al., 1996; Och and Ney, 2003; Fraser and Marcu, 2007).

One way to improve word alignment quality is to add linguistic knowledge derived from a monolingual corpus. This monolingual knowledge makes it easier to determine corresponding words correctly. For instance, functional words in one language tend to correspond to functional words in another language (Deng and Gao, 2007), and the syntactic dependency of words in each language can help the alignment process (Ma et al., 2008). It has been shown that such *grammatical*

information works as a constraint in word alignment models and improves word alignment quality.

A large number of monolingual *lexical* semantic resources such as WordNet (Miller, 1995) have been constructed in more than fifty languages (Sagot and Fiser, 2008). They include word-level relations such as synonyms, hypernyms and hyponyms. Synonym information is particularly helpful for word alignment because we can expect a synonym to correspond to the same word in a different language. In this paper, we explore a method for using synonym information effectively to improve word alignment quality.

In general, synonym relations are defined in terms of word sense, not in terms of word form. In other words, synonym relations are usually context or domain dependent. For instance, ‘head’ and ‘chief’ are synonyms in contexts referring to working environment, while ‘head’ and ‘forefront’ are synonyms in contexts referring to physical positions. It is difficult, however, to imagine a context where ‘chief’ and ‘forefront’ are synonyms. Therefore, it is easy to imagine that simply replacing all occurrences of ‘chief’ and ‘forefront’ with ‘head’ do sometimes harm with word alignment accuracy, and we have to model either the context or senses of words.

We propose a novel method that incorporates synonyms from monolingual resources in a bilingual word alignment model. We formulate a synonym pair generative model with a topic variable and use this model as a regularization term with a bilingual word alignment model. The topic variable in our synonym model is helpful for disambiguating the meanings of synonyms. We extend HM-BiTAM, which is a HMM-based word alignment model with a latent topic, with a novel synonym pair generative model. We applied the proposed method to an English-French word alignment task and successfully improved the word

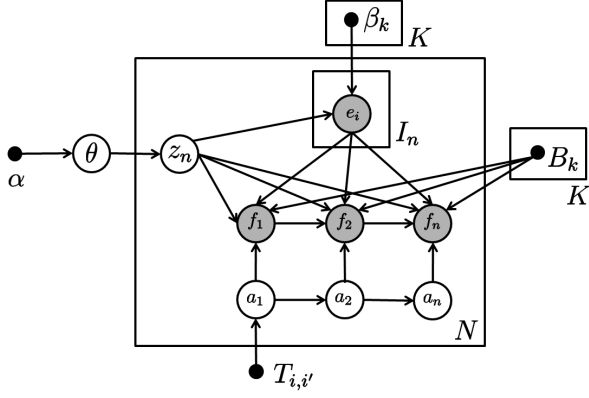


Figure 1: Graphical model of HM-BiTAM

alignment quality.

2 Bilingual Word Alignment Model

In this section, we review a conventional generative word alignment model, HM-BiTAM (Zhao and Xing, 2008).

HM-BiTAM is a bilingual generative model with topic z , alignment a and topic weight vector θ as latent variables. Topic variables such as ‘science’ or ‘economy’ assigned to individual sentences help to disambiguate the meanings of words. HM-BiTAM assumes that the n th bilingual sentence pair, (E_n, F_n) , is generated under a given latent topic $z_n \in \{1, \dots, k, \dots, K\}$, where K is the number of latent topics. Let N be the number of sentence pairs, and I_n and J_n be the lengths of E_n and F_n , respectively. In this framework, all of the bilingual sentence pairs $\{E, F\} = \{(E_n, F_n)\}_{n=1}^N$ are generated as follows.

1. $\theta \sim \text{Dirichlet}(\alpha)$: sample topic-weight vector
2. For each sentence pair (E_n, F_n)
 - (a) $z_n \sim \text{Multinomial}(\theta)$: sample the topic
 - (b) $e_{n,i:I_n} | z_n \sim p(E_n | z_n; \beta)$: sample English words from a monolingual unigram model given topic z_n
 - (c) For each position $j_n = 1, \dots, J_n$
 - i. $a_{j_n} \sim p(a_{j_n} | a_{j_n-1}; T)$: sample an alignment link a_{j_n} from a first order Markov process
 - ii. $f_{j_n} \sim p(f_{j_n} | E_n, a_{j_n}, z_n; B)$: sample a target word f_{j_n} given an aligned source word and topic

where alignment $a_{j_n} = i$ denotes source word e_i and target word f_{j_n} are aligned. α is a parameter over the topic weight vector θ , $\beta = \{\beta_{k,e}\}$ is the source word probability given the k th topic: $p(e | z = k)$. $B = \{B_{f,e,k}\}$ represents the word

translation probability from e to f under the k th topic: $p(f | e, z = k)$. $T = \{T_{i,i'}\}$ is a state transition probability of a first order Markov process. Fig. 1 shows a graphical model of HM-BiTAM.

The total likelihood of bilingual sentence pairs $\{E, F\}$ can be obtained by marginalizing out latent variables z , a and θ ,

$$p(F, E; \Psi) = \sum_z \sum_a \int p(F, E, z, a, \theta; \Psi) d\theta, \quad (1)$$

where $\Psi = \{\alpha, \beta, T, B\}$ is a parameter set. In this model, we can infer word alignment a by maximizing the likelihood above.

3 Proposed Method

3.1 Synonym Pair Generative Model

We design a generative model for synonym pairs $\{f, f'\}$ in language F , which assumes that the synonyms are collected from monolingual linguistic resources. We assume that each synonym pair (f, f') is generated independently given the same ‘sense’ s . Under this assumption, the probability of synonym pair (f, f') can be formulated as,

$$p(f, f') \propto \sum_s p(f | s) p(f' | s) p(s). \quad (2)$$

We define a pair (e, k) as a representation of the sense s , where e and k are a word in a different language E and a latent topic, respectively. It has been shown that a word e in a different language is an appropriate representation of s in synonym modeling (Bannard and Callison-Burch, 2005). We assume that adding a latent topic k for the sense is very useful for disambiguating word meaning, and thus that (e, k) gives us a good approximation of s . Under this assumption, the synonym pair generative model can be defined as follows.

$$p(\{f, f'\}; \tilde{\Psi}) \propto \prod_{(f, f')} \sum_{e, k} p(f | e, k; \tilde{\Psi}) p(f' | e, k; \tilde{\Psi}) p(e, k; \tilde{\Psi}), \quad (3)$$

where $\tilde{\Psi}$ is the parameter set of our model.

3.2 Word Alignment with Synonym Regularization

In this section, we extend the bilingual generative model (HM-BiTAM) with our synonym pair model. Our expectation is that synonym pairs

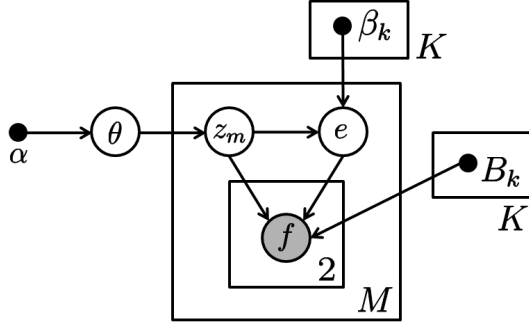


Figure 2: Graphical model of synonym pair generative process

correspond to the same word in a different language, thus they make it easy to infer accurate word alignment. HM-BiTAM and the synonym model share parameters in order to incorporate monolingual synonym information into the bilingual word alignment model. This can be achieved via reparameterizing $\tilde{\Psi}$ in eq. 3 as,

$$p(f | e, k; \tilde{\Psi}) \equiv p(f | e, k; B), \quad (4)$$

$$p(e, k; \tilde{\Psi}) \equiv p(e | k; \beta) p(k; \alpha). \quad (5)$$

Overall, we re-define the synonym pair model with the HM-BiTAM parameter set Ψ ,

$$p(\{f, f'\}; \Psi) \propto \frac{1}{\sum_{k'} \alpha_{k'}} \prod_{(f, f')} \sum_{k, e} \alpha_k \beta_{k, e} B_{f, e, k} B_{f', e, k}. \quad (6)$$

Fig. 2 shows a graphical model of the synonym pair generative process. We estimate the parameter values to maximize the likelihood of HM-BiTAM with respect to bilingual sentences and that of the synonym model with respect to synonym pairs collected from monolingual resources. Namely, the parameter estimate, $\hat{\Psi}$, is computed as

$$\hat{\Psi} = \arg \max_{\Psi} \{ \log p(F, E; \Psi) + \zeta \log p(\{f, f'\}; \Psi) \}, \quad (7)$$

where ζ is a regularization weight that should be set for training. We can expect that the second term of eq. 7 to constrain parameter set Ψ and avoid overfitting for the bilingual word alignment model. We resort to the variational EM approach (Bernardo et al., 2003) to infer $\hat{\Psi}$ following HM-BiTAM. We omit the parameter update equation due to lack of space.

4 Experiments

4.1 Experimental Setting

For an empirical evaluation of the proposed method, we used a bilingual parallel corpus of English-French Hansards (Mihalcea and Pedersen, 2003). The corpus consists of over 1 million sentence pairs, which include 447 manually word-aligned sentences. We selected 100 sentence pairs randomly from the manually word-aligned sentences as development data for tuning the regularization weight ζ , and used the 347 remaining sentence pairs as evaluation data. We also randomly selected 10k, 50k, and 100k sized sentence pairs from the corpus as additional training data. We ran the unsupervised training of our proposed word alignment model on the additional training data and the 347 sentence pairs of the evaluation data. Note that manual word alignment of the 347 sentence pairs was not used for the unsupervised training. After the unsupervised training, we evaluated the word alignment performance of our proposed method by comparing the manual word alignment of the 347 sentence pairs with the prediction provided by the trained model.

We collected English and French synonym pairs from WordNet 2.1 (Miller, 1995) and WOLF 0.1.4 (Sagot and Fiser, 2008), respectively. WOLF is a semantic resource constructed from the Princeton WordNet and various multilingual resources. We selected synonym pairs where both words were included in the bilingual training set.

We compared the word alignment performance of our model with that of GIZA++ 1.03¹ (Vogel et al., 1996; Och and Ney, 2003), and HM-BiTAM (Zhao and Xing, 2008) implemented by us. GIZA++ is an implementation of IBM-model 4 and HMM, and HM-BiTAM corresponds to $\zeta = 0$ in eq. 7. We adopted $K = 3$ topics, following the setting in (Zhao and Xing, 2006).

We trained the word alignment in two directions: English to French, and French to English. The alignment results for both directions were refined with ‘GROW’ heuristics to yield high precision and high recall in accordance with previous work (Och and Ney, 2003; Zhao and Xing, 2006). We evaluated these results for precision, recall, F-measure and alignment error rate (AER), which are standard metrics for word alignment accuracy (Och and Ney, 2000).

¹<http://fjoch.com/GIZA++.html>

10k		Precision	Recall	F-measure	AER
GIZA++	standard	0.856	0.718	0.781	0.207
	with SRH	0.874	0.720	0.789	0.198
HM-BiTAM	standard	0.869	0.788	0.826	0.169
	with SRH	0.884	0.790	0.834	0.160
Proposed		0.941	0.808	0.870	0.123

(a)

50k		Precision	Recall	F-measure	AER
GIZA++	standard	0.905	0.770	0.832	0.156
	with SRH	0.903	0.759	0.825	0.164
HM-BiTAM	standard	0.901	0.814	0.855	0.140
	with SRH	0.899	0.808	0.853	0.145
Proposed		0.947	0.824	0.881	0.112

(b)

100k		Precision	Recall	F-measure	AER
GIZA++	standard	0.925	0.791	0.853	0.136
	with SRH	0.934	0.803	0.864	0.126
HM-BiTAM	standard	0.898	0.851	0.874	0.124
	with SRH	0.909	0.860	0.879	0.114
Proposed		0.927	0.862	0.893	0.103

(c)

Table 1: Comparison of word alignment accuracy. The best results are indicated in bold type. The additional data set sizes are (a) 10k, (b) 50k, (c) 100k.

4.2 Results and Discussion

Table 1 shows the word alignment accuracy of the three methods trained with 10k, 50k, and 100k additional sentence pairs. For all settings, our proposed method outperformed other conventional methods. This result shows that synonym information is effective for improving word alignment quality as we expected.

As mentioned in Sections 1 and 3.1, the main idea of our proposed method is to introduce *latent topics* for modeling synonym pairs, and then to utilize the synonym pair model for the regularization of word alignment models. We expect the latent topics to be useful for modeling polysemous words included in synonym pairs and to enable us to incorporate synonym information effectively into word alignment models. To confirm the effect of the synonym pair model with latent topics, we also tested GIZA++ and HM-BiTAM with what we call *Synonym Replacement Heuristics (SRH)*, where all of the synonym pairs in the bilingual training sentences were simply replaced with a representative word. For instance, the words ‘sick’ and ‘ill’ in the bilingual sentences

# vocabularies		10k	50k	100k
English	standard	8578	16924	22817
	with SRH	5435	7235	13978
French	standard	10791	21872	30294
	with SRH	9737	20077	27970

Table 2: The number of vocabularies in the 10k, 50k and 100k data sets.

were replaced with the word ‘sick’. As shown in Table 2, the number of vocabularies in the English and French data sets decreased as a result of employing the SRH.

We show the performance of GIZA++ and HM-BiTAM with the SRH in the lines entitled “with SRH” in Table 1. The GIZA++ and HM-BiTAM with the SRH slightly outperformed the *standard* GIZA++ and HM-BiTAM for the 10k and 100k data sets, but underperformed with the 50k data set. We assume that the SRH mitigated the overfitting of these models into low-frequency word pairs in bilingual sentences, and then improved the word alignment performance. The SRH regards all of the different words coupled with the same word in the synonym pairs as synonyms. For instance, the words ‘head’, ‘chief’ and ‘forefront’ in the bilingual sentences are replaced with ‘chief’, since (‘head’, ‘chief’) and (‘head’, ‘forefront’) are synonyms. Obviously, (‘chief’, ‘forefront’) are not synonyms, which is detrimental to word alignment.

The proposed method consistently outperformed GIZA++ and HM-BiTAM with the SRH in 10k, 50k and 100k data sets in F-measure. The synonym pair model in our proposed method can automatically learn that (‘head’, ‘chief’) and (‘head’, ‘forefront’) are individual synonyms with different meanings by assigning these pairs to different topics. By sharing latent topics between the synonym pair model and the word alignment model, the synonym information incorporated in the synonym pair model is used directly for training word alignment model. The experimental results show that our proposed method was effective in improving the performance of the word alignment model by using synonym pairs including such *ambiguous* synonym words.

Finally, we discuss the data set size used for unsupervised training. As shown in Table 1, using a large number of additional sentence pairs improved the performance of all the models. In all our experimental settings, all the additional sen-

tence pairs and the evaluation data were selected from the Hansards data set. These experimental results show that a larger number of sentence pairs was more effective in improving word alignment performance when the sentence pairs were collected from a *homogeneous* data source. However, in practice, it might be difficult to collect a large number of such homogeneous sentence pairs for a specific target domain and language pair. One direction for future work is to confirm the effect of the proposed method when training the word alignment model by using a large number of sentence pairs collected from various data sources including many topics for a specific language pair.

5 Conclusions and Future Work

We proposed a novel framework that incorporates synonyms from monolingual linguistic resources in a word alignment generative model. This approach utilizes both bilingual and monolingual synonym resources effectively for word alignment. Our proposed method uses a latent topic for bilingual sentences and monolingual synonym pairs, which is helpful in terms of word sense disambiguation. Our proposed method improved word alignment quality with both small and large data sets. Future work will involve examining the proposed method for different language pairs such as English-Chinese and English-Japanese and evaluating the impact of our proposed method on SMT performance. We will also apply our proposed method to a larger data sets of multiple domains since we can expect a further improvement in word alignment accuracy if we use more bilingual sentences and more monolingual knowledge.

References

C. Bannard and C. Callison-Burch. 2005. Paraphrasing with bilingual parallel corpora. In *Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics*, pages 597–604. Association for Computational Linguistics Morristown, NJ, USA.

J. M. Bernardo, M. J. Bayarri, J. O. Berger, A. P. Dawid, D. Heckerman, A. F. M. Smith, and M. West. 2003. The variational bayesian EM algorithm for incomplete data: with application to scoring graphical model structures. In *Bayesian Statistics 7: Proceedings of the 7th Valencia International Meeting, June 2-6, 2002*, page 453. Oxford University Press, USA.

Y. Deng and Y. Gao. 2007. Guiding statistical word alignment models with prior knowledge. In *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, pages 1–8, Prague, Czech Republic, June. Association for Computational Linguistics.

A. Fraser and D. Marcu. 2007. Getting the structure right for word alignment: LEAF. In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 51–60, Prague, Czech Republic, June. Association for Computational Linguistics.

Y. Ma, S. Ozdowska, Y. Sun, and A. Way. 2008. Improving word alignment using syntactic dependencies. In *Proceedings of the ACL-08: HLT Second Workshop on Syntax and Structure in Statistical Translation (SSST-2)*, pages 69–77, Columbus, Ohio, June. Association for Computational Linguistics.

R. Mihalcea and T. Pedersen. 2003. An evaluation exercise for word alignment. In *Proceedings of the HLT-NAACL 2003 Workshop on building and using parallel texts: data driven machine translation and beyond-Volume 3*, page 10. Association for Computational Linguistics.

G. A. Miller. 1995. WordNet: a lexical database for English. *Communications of the ACM*, 38(11):41.

F. J. Och and H. Ney. 2000. Improved statistical alignment models. In *Proceedings of the 38th Annual Meeting on Association for Computational Linguistics*, pages 440–447. Association for Computational Linguistics Morristown, NJ, USA.

F. J. Och and H. Ney. 2003. A systematic comparison of various statistical alignment models. *Computational Linguistics*, 29(1):19–51.

B. Sagot and D. Fiser. 2008. Building a free French wordnet from multilingual resources. In *Proceedings of Ontolex*.

S. Vogel, H. Ney, and C. Tillmann. 1996. HMM-based word alignment in statistical translation. In *Proceedings of the 16th Conference on Computational Linguistics-Volume 2*, pages 836–841. Association for Computational Linguistics Morristown, NJ, USA.

B. Zhao and E. P. Xing. 2006. BiTAM: Bilingual topic admixture models for word alignment. In *Proceedings of the COLING/ACL on Main Conference Poster Sessions*, page 976. Association for Computational Linguistics.

B. Zhao and E. P. Xing. 2008. HM-BiTAM: Bilingual topic exploration, word alignment, and translation. In *Advances in Neural Information Processing Systems 20*, pages 1689–1696, Cambridge, MA. MIT Press.

Better Filtration and Augmentation for Hierarchical Phrase-Based Translation Rules

Zhiyang Wang[†] Yajuan Lü[†] Qun Liu[†] Young-Sook Hwang[‡]

[†]Key Lab. of Intelligent Information Processing
Institute of Computing Technology
Chinese Academy of Sciences
P.O. Box 2704, Beijing 100190, China
wangzhiyang@ict.ac.cn

[‡]HILab Convergence Technology Center
C&I Business
SKTelecom
11, Euljiro2-ga, Jung-gu, Seoul 100-999, Korea
yshwang@sktelecom.com

Abstract

This paper presents a novel filtration criterion to restrict the rule extraction for the hierarchical phrase-based translation model, where a bilingual but relaxed well-formed dependency restriction is used to filter out bad rules. Furthermore, a new feature which describes the regularity that the source/target dependency edge triggers the target/source word is also proposed. Experimental results show that, the new criteria weeds out about 40% rules while with translation performance improvement, and the new feature brings another improvement to the baseline system, especially on larger corpus.

1 Introduction

Hierarchical phrase-based (HPB) model (Chiang, 2005) is the state-of-the-art statistical machine translation (SMT) model. By looking for phrases that contain other phrases and replacing the sub-phrases with nonterminal symbols, it gets hierarchical rules. Hierarchical rules are more powerful than conventional phrases since they have better generalization capability and could capture long distance reordering. However, when the training corpus becomes larger, the number of rules will grow exponentially, which inevitably results in slow and memory-consuming decoding.

In this paper, we address the problem of reducing the hierarchical translation rule table resorting to the dependency information of bilingual languages. We only keep rules that both sides are *relaxed-well-formed* (RWF) dependency structure (see the definition in Section 3), and discard others which do not satisfy this constraint. In this way, about 40% bad rules are weeded out from the original rule table. However, the performance is even better than the traditional HPB translation system.

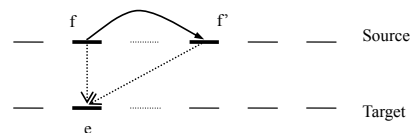


Figure 1: Solid wire reveals the dependency relation pointing from the child to the parent. Target word e is triggered by the source word f and its head word f' , $p(e|f \rightarrow f')$.

Based on the *relaxed-well-formed* dependency structure, we also introduce a new linguistic feature to enhance translation performance. In the traditional phrase-based SMT model, there are always lexical translation probabilities based on IBM model 1 (Brown et al., 1993), i.e. $p(e|f)$, namely, the target word e is triggered by the source word f . Intuitively, however, the generation of e is not only involved with f , sometimes may also be triggered by other context words in the source side. Here we assume that the dependency edge ($f \rightarrow f'$) of word f generates target word e (we call it head word trigger in Section 4). Therefore, two words in one language trigger one word in another, which provides a more sophisticated and better choice for the target word, i.e. Figure 1. Similarly, the dependency feature works well in Chinese-to-English translation task, especially on large corpus.

2 Related Work

In the past, a significant number of techniques have been presented to reduce the hierarchical rule table. He et al. (2009) just used the key phrases of source side to filter the rule table without taking advantage of any linguistic information. Iglesias et al. (2009) put rules into syntactic classes based on the number of non-terminals and patterns, and applied various filtration strategies to improve the rule table quality. Shen et al. (2008) discarded

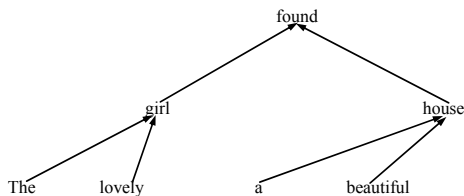


Figure 2: An example of dependency tree. The corresponding plain sentence is *The lovely girl found a beautiful house.*

most entries of the rule table by using the constraint that rules of the target-side are well-formed (WF) dependency structure, but this filtering led to degradation in translation performance. They obtained improvements by adding an additional dependency language model. The basic difference of our method from (Shen et al., 2008) is that we keep rules that both sides should be *relaxed-well-formed* dependency structure, not just the target side. Besides, our system complexity is not increased because no additional language model is introduced.

The feature of head word trigger which we apply to the log-linear model is motivated by the trigger-based approach (Hasan and Ney, 2009). Hasan and Ney (2009) introduced a second word to trigger the target word without considering any linguistic information. Furthermore, since the second word can come from any part of the sentence, there may be a prohibitively large number of parameters involved. Besides, He et al. (2008) built a maximum entropy model which combines rich context information for selecting translation rules during decoding. However, as the size of the corpus increases, the maximum entropy model will become larger. Similarly, In (Shen et al., 2009), context language model is proposed for better rule selection. Taking the dependency edge as condition, our approach is very different from previous approaches of exploring context information.

3 Relaxed-well-formed Dependency Structure

Dependency models have recently gained considerable interest in SMT (Ding and Palmer, 2005; Quirk et al., 2005; Shen et al., 2008). Dependency tree can represent richer structural information. It reveals long-distance relation between words and directly models the semantic structure of a sentence without any constituent labels. Fig-

ure 2 shows an example of a dependency tree. In this example, the word *found* is the root of the tree.

Shen et al. (2008) propose the well-formed dependency structure to filter the hierarchical rule table. A well-formed dependency structure could be either a single-rooted dependency tree or a set of sibling trees. Although most rules are discarded with the constraint that the target side should be well-formed, this filtration leads to degradation in translation performance.

As an extension of the work of (Shen et al., 2008), we introduce the so-called *relaxed-well-formed* dependency structure to filter the hierarchical rule table. Given a sentence $S = w_1 w_2 \dots w_n$. Let $d_1 d_2 \dots d_n$ represent the position of parent word for each word. For example, $d_3 = 4$ means that w_3 depends on w_4 . If w_i is a root, we define $d_i = -1$.

Definition A dependency structure $w_i \dots w_j$ is a *relaxed-well-formed* structure, where there is $h \notin [i, j]$, all the words $w_i \dots w_j$ are directly or indirectly depended on w_h or -1 (here we define $h = -1$). If and only if it satisfies the following conditions

- $d_h \notin [i, j]$
- $\forall k \in [i, j], d_k \in [i, j]$ or $d_k = h$

From the definition above, we can see that the *relaxed-well-formed* structure obviously covers the well-formed one. In this structure, we don't constrain that all the children of the sub-root should be complete. Let's review the dependency tree in Figure 2 as an example. Except for the well-formed structure, we could also extract *girl found a beautiful house*. Therefore, if the modifier *The lovely* changes to *The cute*, this rule also works.

4 Head Word Trigger

(Koehn et al., 2003) introduced the concept of lexical weighting to check how well words of the phrase translate to each other. Source word f aligns with target word e , according to the IBM model 1, the lexical translation probability is $p(e|f)$. However, in the sense of dependency relationship, we believe that the generation of the target word e , is not only triggered by the aligned source word f , but also associated with f 's head word f' . Therefore, the lexical translation probability becomes $p(e|f \rightarrow f')$, which of course allows for a more fine-grained lexical choice of

the target word. More specifically, the probability could be estimated by the maximum likelihood (MLE) approach,

$$p(e|f \rightarrow f') = \frac{\text{count}(e, f \rightarrow f')}{\sum_{e'} \text{count}(e', f \rightarrow f')} \quad (1)$$

Given a phrase pair \bar{f} , \bar{e} and word alignment a , and the dependent relation of the source sentence d_1^J (J is the length of the source sentence, I is the length of the target sentence). Therefore, given the lexical translation probability distribution $p(e|f \rightarrow f')$, we compute the feature score of a phrase pair (\bar{f}, \bar{e}) as

$$p(\bar{e}|\bar{f}, d_1^J, a) = \prod_{i=1}^{|\bar{e}|} \frac{1}{|\{j|(j, i) \in a\}|} \sum_{\forall (j, i) \in a} p(e_i|f_j \rightarrow f_{d_j}) \quad (2)$$

Now we get $p(\bar{e}|\bar{f}, d_1^J, a)$, we could obtain $p(\bar{f}|\bar{e}, d_1^I, a)$ (d_1^I represents dependent relation of the target side) in the similar way. This new feature can be easily integrated into the log-linear model as lexical weighting does.

5 Experiments

In this section, we describe the experimental setting used in this work, and verify the effect of the *relaxed-well-formed* structure filtering and the new feature, head word trigger.

5.1 Experimental Setup

Experiments are carried out on the NIST¹ Chinese-English translation task with two different size of training corpora.

- **FBIS**: We use the FBIS corpus as the first training corpus, which contains 239K sentence pairs with 6.9M Chinese words and 8.9M English words.
- **GQ**: This is manually selected from the LDC² corpora. GQ contains 1.5M sentence pairs with 41M Chinese words and 48M English words. In fact, FBIS is the subset of GQ.

¹www.nist.gov/speech/tests/mt

²It consists of six LDC corpora: LDC2002E18, LDC2003E07, LDC2003E14, Hansards part of LDC2004T07, LDC2004T08, LDC2005T06.

For language model, we use the SRI Language Modeling Toolkit (Stolcke, 2002) to train a 4-gram model on the first 1/3 of the Xinhua portion of GIGAWORD corpus. And we use the NIST 2002 MT evaluation test set as our development set, and NIST 2004, 2005 test sets as our blind test sets. We evaluate the translation quality using *case-insensitive* BLEU metric (Papineni et al., 2002) without dropping OOV words, and the feature weights are tuned by minimum error rate training (Och, 2003).

In order to get the dependency relation of the training corpus, we re-implement a beam-search style monolingual dependency parser according to (Nivre and Scholz, 2004). Then we use the same method suggested in (Chiang, 2005) to extract SCFG grammar rules within dependency constraint on both sides except that unaligned words are allowed at the edge of phrases. Parameters of head word trigger are estimated as described in Section 4. As a default, the maximum initial phrase length is set to 10 and the maximum rule length of the source side is set to 5. Besides, we also re-implement the decoder of Hiero (Chiang, 2007) as our baseline. In fact, we just exploit the dependency structure during the rule extraction phase. Therefore, we don't need to change the main decoding algorithm of the SMT system.

5.2 Results on FBIS Corpus

A series of experiments was done on the FBIS corpus. We first parse the bilingual languages with monolingual dependency parser respectively, and then only retain the rules that both sides are in line with the constraint of dependency structure. In Table 1, the *relaxed-well-formed* structure filtered out 35% of the rule table and the well-formed discarded 74%. *RWF* extracts additional 39% compared to *WF*, which can be seen as some kind of evidence that the rules we additional get seem common in the sense of linguistics. Compared to (Shen et al., 2008), we just use the dependency structure to constrain rules, not to maintain the tree structures to guide decoding.

Table 2 shows the translation result on FBIS. We can see that the *RWF* structure constraint can improve translation quality substantially both at development set and different test sets. On the Test04 task, it gains +0.86% BLEU, and +0.84% on Test05. Besides, we also used Shen et al. (2008)'s *WF* structure to filter both sides. Although it discard about 74% of the rule table, the

System	Rule table size
<i>HPB</i>	30,152,090
<i>RWF</i>	19,610,255
<i>WF</i>	7,742,031

Table 1: Rule table size with different constraint on FBIS. Here *HPB* refers to the baseline hierarchal phrase-based system, *RWF* means *relaxed-well-formed* constraint and *WF* represents the well-formed structure.

System	Dev02	Test04	Test05
<i>HPB</i>	0.3285	0.3284	0.2965
<i>WF</i>	0.3125	0.3218	0.2887
<i>RWF</i>	0.3326	0.3370**	0.3050
<i>RWF+Tri</i>	0.3281	/	0.2965

Table 2: Results of FBIS corpus. Here *Tri* means the feature of head word trigger on both sides. And we don't test the new feature on Test04 because of the bad performance on development set. * or ** = significantly better than baseline ($p < 0.05$ or 0.01 , respectively).

over-all BLEU is decreased by 0.66%-0.78% on the test sets.

As for the feature of head word trigger, it seems not work on the FBIS corpus. On Test05, it gets the same score with the baseline, but lower than *RWF* filtering. This may be caused by the data sparseness problem, which results in inaccurate parameter estimation of the new feature.

5.3 Result on GQ Corpus

In this part, we increased the size of the training corpus to check whether the feature of head word trigger works on large corpus.

We get 152M rule entries from the GQ corpus according to (Chiang, 2007)'s extraction method. If we use the *RWF* structure to constrain both sides, the number of rules is 87M, about 43% of rule entries are discarded. From Table 3, the new

System	Dev02	Test04	Test05
<i>HPB</i>	0.3473	0.3386	0.3206
<i>RWF</i>	0.3539	0.3485**	0.3228
<i>RWF+Tri</i>	0.3540	0.3607**	0.3339*

Table 3: Results of GQ corpus. * or ** = significantly better than baseline ($p < 0.05$ or 0.01 , respectively).

feature works well on two different test sets. The gain is +2.21% BLEU on Test04, and +1.33% on Test05. Compared to the result of the baseline, only using the *RWF* structure to filter performs the same as the baseline on Test05, and +0.99% gains on Test04.

6 Conclusions

This paper proposes a simple strategy to filter the hierarchal rule table, and introduces a new feature to enhance the translation performance. We employ the *relaxed-well-formed* dependency structure to constrain both sides of the rule, and about 40% of rules are discarded with improvement of the translation performance. In order to make full use of the dependency information, we assume that the target word e is triggered by dependency edge of the corresponding source word f . And this feature works well on large parallel training corpus.

How to estimate the probability of head word trigger is very important. Here we only get the parameters in a generative way. In the future, we we are plan to exploit some discriminative approach to train parameters of this feature, such as EM algorithm (Hasan et al., 2008) or maximum entropy (He et al., 2008).

Besides, the quality of the parser is another effect for this method. As the next step, we will try to exploit bilingual knowledge to improve the monolingual parser, i.e. (Huang et al., 2009).

Acknowledgments

This work was partly supported by National Natural Science Foundation of China Contract 60873167. It was also funded by SK Telecom, Korea under the contract 4360002953. We show our special thanks to Wenbin Jiang and Shu Cai for their valuable suggestions. We also thank the anonymous reviewers for their insightful comments.

References

- Peter F. Brown, Vincent J. Della Pietra, Stephen A. Della Pietra, and Robert L. Mercer. 1993. The mathematics of statistical machine translation: parameter estimation. *Comput. Linguist.*, 19(2):263–311.
- David Chiang. 2005. A hierarchical phrase-based model for statistical machine translation. In *ACL*

- '05: *Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics*, pages 263–270.
- David Chiang. 2007. Hierarchical phrase-based translation. *Comput. Linguist.*, 33(2):201–228.
- Yuan Ding and Martha Palmer. 2005. Machine translation using probabilistic synchronous dependency insertion grammars. In *ACL '05: Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics*, pages 541–548.
- Saša Hasan and Hermann Ney. 2009. Comparison of extended lexicon models in search and rescoring for smt. In *NAACL '09: Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics, Companion Volume: Short Papers*, pages 17–20.
- Saša Hasan, Juri Ganitkevitch, Hermann Ney, and Jesús Andrés-Ferrer. 2008. Triplet lexicon models for statistical machine translation. In *EMNLP '08: Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 372–381.
- Zhongjun He, Qun Liu, and Shouxun Lin. 2008. Improving statistical machine translation using lexicalized rule selection. In *COLING '08: Proceedings of the 22nd International Conference on Computational Linguistics*, pages 321–328.
- Zhongjun He, Yao Meng, Yajuan Lü, Hao Yu, and Qun Liu. 2009. Reducing smt rule table with monolingual key phrase. In *ACL-IJCNLP '09: Proceedings of the ACL-IJCNLP 2009 Conference Short Papers*, pages 121–124.
- Liang Huang, Wenbin Jiang, and Qun Liu. 2009. Bilingually-constrained (monolingual) shift-reduce parsing. In *EMNLP '09: Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing*, pages 1222–1231.
- Gonzalo Iglesias, Adrià de Gispert, Eduardo R. Barga, and William Byrne. 2009. Rule filtering by pattern for efficient hierarchical translation. In *EACL '09: Proceedings of the 12th Conference of the European Chapter of the Association for Computational Linguistics*, pages 380–388.
- Philipp Koehn, Franz Josef Och, and Daniel Marcu. 2003. Statistical phrase-based translation. In *NAACL '03: Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology*, pages 48–54.
- Joakim Nivre and Mario Scholz. 2004. Deterministic dependency parsing of english text. In *COLING '04: Proceedings of the 20th international conference on Computational Linguistics*, pages 64–70.
- Franz Josef Och. 2003. Minimum error rate training in statistical machine translation. In *ACL '03: Proceedings of the 41st Annual Meeting on Association for Computational Linguistics*, pages 160–167.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *ACL '02: Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, pages 311–318.
- Chris Quirk, Arul Menezes, and Colin Cherry. 2005. Dependency treelet translation: syntactically informed phrasal smt. In *ACL '05: Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics*, pages 271–279.
- Libin Shen, Jinxi Xu, and Ralph Weischedel. 2008. A new string-to-dependency machine translation algorithm with a target dependency language model. In *Proceedings of ACL-08: HLT*, pages 577–585.
- Libin Shen, Jinxi Xu, Bing Zhang, Spyros Matsoukas, and Ralph Weischedel. 2009. Effective use of linguistic and contextual information for statistical machine translation. In *EMNLP '09: Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing*, pages 72–80.
- Andreas Stolcke. 2002. Srilman extensible language modeling toolkit. In *In Proceedings of the 7th International Conference on Spoken Language Processing (ICSLP 2002)*, pages 901–904.

Fixed Length Word Suffix for Factored Statistical Machine Translation

Narges Sharif Razavian
School of Computer Science
Carnegie Mellon University
Pittsburgh, USA
nsharifr@cs.cmu.edu

Stephan Vogel
School of Computer Science
Carnegie Mellon University
Pittsburgh, USA
stephan.vogel@cs.cmu.edu

Abstract

Factored Statistical Machine Translation extends the Phrase Based SMT model by allowing each word to be a vector of factors. Experiments have shown effectiveness of many factors, including the Part of Speech tags in improving the grammaticality of the output. However, high quality part of speech taggers are not available in open domain for many languages. In this paper we used fixed length word suffix as a new factor in the Factored SMT, and were able to achieve significant improvements in three set of experiments: large NIST Arabic to English system, medium WMT Spanish to English system, and small TRANSTAC English to Iraqi system.

1 Introduction

Statistical Machine Translation(SMT) is currently the state of the art solution to the machine translation. Phrase based SMT is also among the top performing approaches available as of today. This approach is a purely lexical approach, using surface forms of the words in the parallel corpus to generate the translations and estimate probabilities. It is possible to incorporate syntactical information into this framework through different ways. Source side syntax based re-ordering as preprocessing step, dependency based re-ordering models, cohesive decoding features are among many available successful attempts for the integration of syntax into the translation model. Factored translation modeling is another way to achieve this goal. These models allow each word to be represented as a vector of factors rather than a single surface form. Factors can represent richer expression power on each word. Any factors such as word stems, gender, part of speech, tense, etc. can be easily used in this framework.

Previous work in factored translation modeling have reported consistent improvements from Part of Speech(POS) tags, morphology, gender, and case factors (Koehn et. a. 2007). In another work, Birch et. al. 2007 have achieved improvement using Combinational Categorical Grammar (CCG) super-tag factors. Creating the factors is done as a preprocessing step, and so far, most of the experiments have assumed existence of external tools for the creation of these factors (i. e. Part of speech taggers, CCG parsers, etc.). Unfortunately high quality language processing tools, especially for the open domain, are not available for most languages.

While linguistically identifiable representations (i.e. POS tags, CCG supertags, etc) have been very frequently used as factors in many applications including MT, simpler representations have also been effective in achieving the same result in other application areas. Grzymala-Busse and Old 1997, DINCER et.al. 2008, were able to use fixed length suffixes as features for training a POS tagging. In another work Saberi and Perrot 1999 showed that reversing middle chunks of the words while keeping the first and last part intact, does not decrease listeners' recognition ability. This result is very relevant to Machine Translation, suggesting that inaccurate context which is usually modeled with n-gram language models, can still be as effective as accurate surface forms. Another research (Rawlinson 1997) confirms this finding; this time in textual domain, observing that randomization of letters in the middle of words has little or no effect on the ability of skilled readers to understand the text. These results suggest that the inexpensive representational factors which do not need unavailable tools might also be worth investigating.

These results encouraged us to introduce language independent simple factors for machine translation. In this paper, following the work of Grzymala-Busse et. al. we used fixed length suf-

fix as word factor, to lower the perplexity of the language model, and have the factors roughly function as part of speech tags, thus increasing the grammaticality of the translation results. We were able to obtain consistent, significant improvements over our baseline in 3 different experiments, large NIST Arabic to English system, medium WMT Spanish to English system, and small TRANSTAC English to Iraqi system.

The rest of this paper is as follows. Section 2 briefly reviews the Factored Translation Models. In section 3 we will introduce our model, and section 4 will contain the experiments and the analysis of the results, and finally, we will conclude this paper in section 5.

2 Factored Translation Model

Statistical Machine Translation uses the log linear combination of a number of features, to compute the highest probable hypothesis as the translation.

$$e = \operatorname{argmax}_e p(\mathbf{e}|\mathbf{f}) = \operatorname{argmax}_e p \exp \sum_{i=1}^n \lambda_i h_i(\mathbf{e}, \mathbf{f})$$

In phrase based SMT, assuming the source and target phrase segmentation as $\{(f_i, e_i)\}$, the most important features include: the Language Model feature $h_{lm}(\mathbf{e}, \mathbf{f}) = p_{lm}(e)$; the phrase translation feature $h_t(\mathbf{e}, \mathbf{f})$ defined as product of translation probabilities, lexical probabilities and phrase penalty; and the reordering probability, $h_d(\mathbf{e}, \mathbf{f})$, usually defined as $\pi_{i=1}^n d(\text{start}_i, \text{end}_{i-1})$ over the source phrase reordering events.

Factored Translation Model, recently introduced by (Koehn et. al. 2007), allow words to have a vector representation. The model can then extend the definition of each of the features from a uni-dimensional value to an arbitrary joint and conditional combination of features. Phrase based SMT is in fact a special case of Factored SMT.

The factored features are defined as an extension of phrase translation features. The function $\tau(f_j, e_j)$, which was defined for a phrase pair before, can now be extended as a log linear combination $\sum_f \tau(f_{jt}, e_{jt})$. The model also allows for a generation feature, defining the relationship between final surface form and target factors. Other features include additional language model features over individual factors, and factored reordering features.

Figure 1 shows an example of a possible factored model.

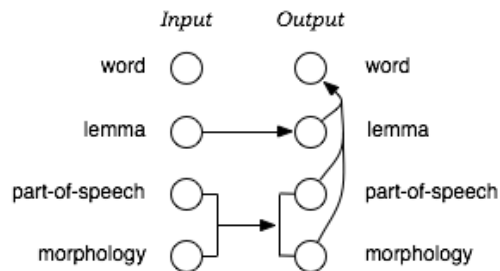


Figure 1: An example of a Factored Translation and Generation Model

In this particular model, words on both source and target side are represented as a vector of four factors: surface form, lemma, part of speech (POS) and the morphology. The target phrase is generated as follows: Source word lemma generates target word lemma. Source word's Part of speech and morphology together generate the target word's part of speech and morphology, and from its lemma, part of speech and morphology the surface form of the target word is finally generated. This model has been able to result in higher translation BLEU score as well as grammatical coherency for English to German, English to Spanish, English to Czech, English to Chinese, Chinese to English and German to English.

3 Fixed Length Suffix Factors for Factored Translation Modeling

Part of speech tagging, constituent and dependency parsing, combinatory categorical grammar super tagging are used extensively in most applications when syntactic representations are needed. However training these tools require medium size treebanks and tagged data, which for most languages will not be available for a while. On the other hand, many simple words features, such as their character n-grams, have in fact proven to be comparably as effective in many applications.

(Keikha et. al. 2008) did an experiment on text classification on noisy data, and compared several word representations. They compared surface form, stemmed words, character n-grams, and semantic relationships, and found that for noisy and open domain text, character-ngrams outperform other representations when used for text classification. In another work (Dincer et al 2009) showed that using fixed length word ending outperforms whole word representation for training a part of speech tagger for Turkish language.

Based on this result, we proposed a suffix factored model for translation, which is shown in Figure 2.

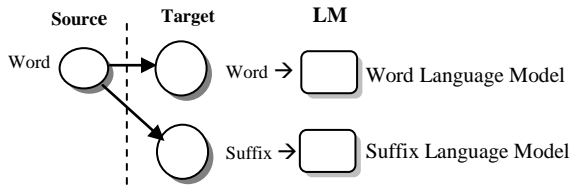


Figure 2: Suffix Factored model: Source word determines factor vectors (target word, target word suffix) and each factor will be associated with its language model.

Based on this model, the final probability of the translation hypothesis will be the log linear combination of phrase probabilities, reordering model probabilities, and each of the language models' probabilities.

$$\begin{aligned}
 P(e|f) &\sim p_{lm-word}(e_{word}) * p_{lm-suffix}(e_{suffix}) \\
 &* \sum_{i=1}^n p(e_{word-j} \& e_{suffix-j} | f_j) \\
 &* \sum_{i=1}^n p(f_j | e_{word-j} \& e_{suffix-j})
 \end{aligned}$$

Where $p_{lm-word}$ is the n-gram language model probability over the word surface sequence, with the language model built from the surface forms. Similarly, $p_{lm-suffix}(e_{suffix})$ is the language model probability over suffix sequences. $p(e_{word-j} \& e_{suffix-j} | f_j)$ and $p(f_j | e_{word-j} \& e_{suffix-j})$ are translation probabilities for each phrase pair i , used in by the decoder. This probability is estimated after the phrase extraction step which is based on grow-diag heuristic at this stage.

4 Experiments and Results

We used Moses implementation of the factored model for training the feature weights, and SRI toolkit for building n-gram language models. The baseline for all systems included the moses system with lexicalized re-ordering, SRI 5-gram language models.

4.1 Small System from Dialog Domain: English to Iraqi

This system was TRANSTAC system, which was built on about 650K sentence pairs with the average sentence length of 5.9 words. After choosing length 3 for suffixes, we built a new parallel corpus, and SRI 5-gram language models for each factor. Vocabulary size for the surface form was 110K whereas the word suffixes had

about 8K distinct words. Table 1 shows the result (BLEU Score) of the system compared to the baseline.

System	Tune on Set- July07	Test on Set- June08	Test on Set- Nov08
Baseline	27.74	21.73	15.62
Factored	28.83	22.84	16.41
Improvement	1.09	1.11	0.79

Table 1: BLEU score, English to Iraqi Transtac system, comparing Factored and Baseline systems.

As you can see, this improvement is consistent over multiple unseen datasets. Arabic cases and numbers show up as the word suffix. Also, verb numbers usually appear partly as word suffix and in some cases as word prefix. Defining a language model over the word endings increases the probability of sequences which have this case and number agreement, favoring correct agreements over the incorrect ones.

4.2 Medium System on Travel Domain: Spanish to English

This system is the WMT08 system, on a corpus of 1.2 million sentence pairs with average sentence length 27.9 words. Like the previous experiment, we defined the 3 character suffix of the words as the second factor, and built the language model and reordering model on the joint event of (surface, suffix) pairs. We built 5-gram language models for each factor. The system had about 97K distinct vocabulary in the surface language model, which was reduced to 8K using the suffix corpus. Having defined the baseline, the system results are as follows.

System	Tune-WMT06	Test set-WMT08
Baseline	33.34	32.53
Factored	33.60	32.84
Improvement	0.26	0.32

Table 2: BLEU score, Spanish to English WMT system, comparing Factored and Baseline systems.

Here, we see improvement with the suffix factors compared to the baseline system. Word endings in English language are major indicators of word's part of speech in the sentence. In fact

most common stemming algorithm, Porter’s Stemmer, works by removing word’s suffix. Having a language model on these suffixes pushes the common patterns of these suffixes to the top, making the more grammatically coherent sentences to achieve a better probability.

4.3 Large NIST 2009 System: Arabic to English

We used NIST2009 system as our baseline in this experiment. The corpus had about 3.8 Million sentence pairs, with average sentence length of 33.4 words. The baseline defined the lexicalized reordering model. As before we defined 3 character long word endings, and built 5-gram SRI language models for each factor. The result of this experiment is shown in table 3.

System	Tune on MT06	Test on Dev07 News Wire	Test on Dev07 Weblog	Test on MT08
Baseline	43.06	48.87	37.84	41.70
Factored	44.20	50.39	39.93	42.74
Improve ment	1.14	1.52	2.09	1.04

Table 3: BLEU score, Arabic to English NIST 2009 system, comparing Factored and Baseline systems.

This result confirms the positive effect of the suffix factors even on large systems. As mentioned before we believe that this result is due to the ability of the suffix to reduce the word into a very simple but rough grammatical representation. Defining language models for this factor forces the decoder to prefer sentences with more probable suffix sequences, which is believed to increase the grammaticality of the result. Future error analysis will show us more insight of the exact effect of this factor on the outcome.

5 Conclusion

In this paper we introduced a simple yet very effective factor: fixed length word suffix, to use in Factored Translation Models. This simple factor has been shown to be effective as a rough replacement for part of speech. We tested our factors in three experiments in a small, English to Iraqi system, a medium sized system of Spanish to English, and a large system, NIST09 Arabic to English. We observed consistent and significant

improvements over the baseline. This result, obtained from the language independent and inexpensive factor, shows promising new opportunities for all language pairs.

References

- Birch, A., Osborne, M., and Koehn, P. CCG supertags in factored statistical machine translation. Proceedings of the Second Workshop on Statistical Machine Translation, pages 9–16, Prague, Czech Republic. Association for Computational Linguistics, 2007.
- Dincer T., Karaoglan B. and Kisla T., A Suffix Based Part-Of-Speech Tagger For Turkish, Fifth International Conference on Information Technology: New Generations, 2008.
- Grzymala-Busse J.W., Old L.J. A machine learning experiment to determine part of speech from word-endings, Lecture Notes in Computer Science, Communications Session 6B Learning and Discovery Systems, 1997.
- Keikha M., Sharif Razavian N, Oroumchian F., and Seyed Razi H., Document Representation and Quality of Text: An Analysis, Chapter 12, Survey of Text Mining II, Springer London, 2008.
- Koehn Ph., Hoang H., Factored Translation Models, Proceedings of 45th Annual Meeting of the Association for Computational Linguistics (ACL), 2007.
- Rawlinson G. E., The significance of letter position in word recognition, PhD Thesis, Psychology Department, University of Nottingham, Nottingham UK, 1976.
- Saberi K and Perrot D R, Cognitive restoration of reversed speech, Nature (London) 1999.

Unsupervised Discourse Segmentation of Documents with Inherently Parallel Structure

Minwoo Jeong and Ivan Titov

Saarland University

Saarbrücken, Germany

{m.jeong|titov}@mmci.uni-saarland.de

Abstract

Documents often have inherently parallel structure: they may consist of a text and commentaries, or an abstract and a body, or parts presenting alternative views on the same problem. Revealing relations between the parts by jointly segmenting and predicting links between the segments, would help to visualize such documents and construct friendlier user interfaces. To address this problem, we propose an unsupervised Bayesian model for joint discourse segmentation and alignment. We apply our method to the “English as a second language” podcast dataset where each episode is composed of two parallel parts: a story and an explanatory lecture. The predicted topical links uncover hidden relations between the stories and the lectures. In this domain, our method achieves competitive results, rivaling those of a previously proposed supervised technique.

1 Introduction

Many documents consist of parts exhibiting a high degree of parallelism: e.g., abstract and body of academic publications, summaries and detailed news stories, etc. This is especially common with the emergence of the Web 2.0 technologies: many texts on the web are now accompanied with comments and discussions. Segmentation of these parallel parts into coherent fragments and discovery of hidden relations between them would facilitate the development of better user interfaces and improve the performance of summarization and information retrieval systems.

Discourse segmentation of the documents composed of parallel parts is a novel and challenging problem, as previous research has mostly focused on the linear segmentation of isolated texts

(e.g., (Hearst, 1994)). The most straightforward approach would be to use a pipeline strategy, where an existing segmentation algorithm finds discourse boundaries of each part independently, and then the segments are aligned. Or, conversely, a sentence-alignment stage can be followed by a segmentation stage. However, as we will see in our experiments, these strategies may result in poor segmentation and alignment quality.

To address this problem, we construct a non-parametric Bayesian model for joint segmentation and alignment of parallel parts. In comparison with the discussed pipeline approaches, our method has two important advantages: (1) it leverages the lexical cohesion phenomenon (Halliday and Hasan, 1976) in modeling the parallel parts of documents, and (2) ensures that the effective number of segments can grow adaptively. Lexical cohesion is an idea that topically-coherent segments display compact lexical distributions (Hearst, 1994; Utiyama and Isahara, 2001; Eisenstein and Barzilay, 2008). We hypothesize that not only isolated fragments but also each group of linked fragments displays a compact and consistent lexical distribution, and our generative model leverages this inter-part cohesion assumption.

In this paper, we consider the dataset of “English as a second language” (ESL) podcast¹, where each episode consists of two parallel parts: a *story* (an example monologue or dialogue) and an explanatory *lecture* discussing the meaning and usage of English expressions appearing in the story. Fig. 1 presents an example episode, consisting of two parallel parts, and their hidden topical relations.² From the figure we may conclude that there is a tendency of word repetition between each pair of aligned segments, illustrating our hypothesis of compactness of their joint distribution. Our goal is

¹<http://www.eslpod.com/>

²Episode no. 232 post on Jan. 08, 2007.

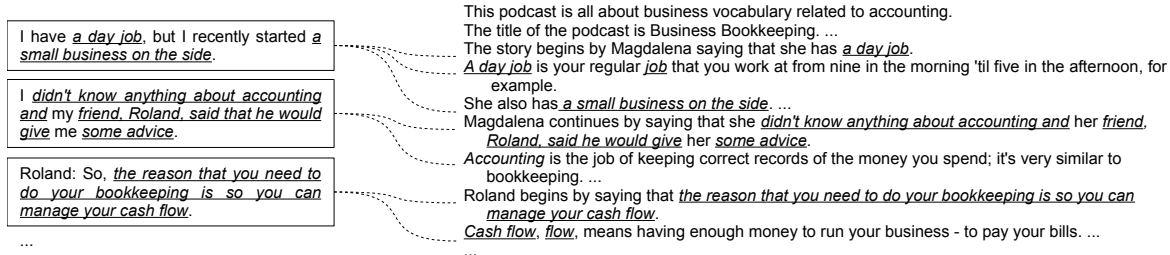


Figure 1: An example episode of ESL podcast. Co-occurred words are represented in *italic* and underline.

to divide the lecture transcript into discourse units and to align each unit to the related segment of the story. Predicting these structures for the ESL podcast could be the first step in development of an e-learning system and a podcast search engine for ESL learners.

2 Related Work

Discourse segmentation has been an active area of research (Hearst, 1994; Utiyama and Isahara, 2001; Galley et al., 2003; Malioutov and Barzilay, 2006). Our work extends the Bayesian segmentation model (Eisenstein and Barzilay, 2008) for isolated texts, to the problem of segmenting parallel parts of documents.

The task of aligning each sentence of an abstract to one or more sentences of the body has been studied in the context of summarization (Marcu, 1999; Jing, 2002; Daumé and Marcu, 2004). Our work is different in that we do not try to extract the most relevant sentence but rather aim to find coherent fragments with maximally overlapping lexical distributions. Similarly, the query-focused summarization (e.g., (Daumé and Marcu, 2006)) is also related but it focuses on sentence extraction rather than on joint segmentation.

We are aware of only one previous work on joint segmentation and alignment of multiple texts (Sun et al., 2007) but their approach is based on similarity functions rather than on modeling lexical cohesion in the generative framework. Our application, the analysis of the ESL podcast, was previously studied in (Noh et al., 2010). They proposed a supervised method which is driven by pairwise classification decisions. The main drawback of their approach is that it neglects the discourse structure and the lexical cohesion phenomenon.

3 Model

In this section we describe our model for discourse segmentation of documents with inherently parallel structure. We start by clarifying our assumptions about their structure.

We assume that a document x consists of K parallel parts, that is, $x = \{x^{(k)}\}_{k=1:K}$, and each part of the document consists of segments, $x^{(k)} = \{s_i^{(k)}\}_{i=1:I}$. Note that the effective number of fragments I is unknown. Each segment can either be specific to this part (drawn from a *part-specific* language model $\phi_i^{(k)}$) or correspond to the entire document (drawn from a *document-level* language model $\phi_i^{(doc)}$). For example, the first and the second sentences of the lecture transcript in Fig. 1 are part-specific, whereas other linked sentences belong to the document-level segments. The document-level language models define topical links between segments in different parts of the document, whereas the part-specific language models define the linear segmentation of the remaining unaligned text.

Each document-level language model corresponds to the set of aligned segments, at most one segment per part. Similarly, each part-specific language model corresponds to a single segment of the single corresponding part. Note that all the documents are modeled independently, as we aim not to discover collection-level topics (as e.g. in (Blei et al., 2003)), but to perform joint discourse segmentation and alignment.

Unlike (Eisenstein and Barzilay, 2008), we cannot make an assumption that the number of segments is known a-priori, as the effective number of part-specific segments can vary significantly from document to document, depending on their size and structure. To tackle this problem, we use Dirichlet processes (DP) (Ferguson, 1973) to de-

fine priors on the number of segments. We incorporate them in our model in a similar way as it is done for the Latent Dirichlet Allocation (LDA) by Yu et al. (2005). Unlike the standard LDA, the topic proportions are chosen not from a Dirichlet prior but from the marginal distribution $GEM(\alpha)$ defined by the *stick breaking* construction (Sethuraman, 1994), where α is the concentration parameter of the underlying DP distribution. $GEM(\alpha)$ defines a distribution of partitions of the unit interval into a countable number of parts.

The formal definition of our model is as follows:

- Draw the document-level topic proportions $\beta^{(doc)} \sim GEM(\alpha^{(doc)})$.
- Choose the document-level language model $\phi_i^{(doc)} \sim Dir(\gamma^{(doc)})$ for $i \in \{1, 2, \dots\}$.
- Draw the part-specific topic proportions $\beta^{(k)} \sim GEM(\alpha^{(k)})$ for $k \in \{1, \dots, K\}$.
- Choose the part-specific language models $\phi_i^{(k)} \sim Dir(\gamma^{(k)})$ for $k \in \{1, \dots, K\}$ and $i \in \{1, 2, \dots\}$.
- For each part k and each sentence n :
 - Draw type $t_n^{(k)} \sim Unif(Doc, Part)$.
 - If $(t_n^{(k)} = Doc)$; draw topic $z_n^{(k)} \sim \beta^{(doc)}$; generate words $\mathbf{x}_n^{(k)} \sim Mult(\phi_{z_n^{(k)}}^{(doc)})$
 - Otherwise; draw topic $z_n^{(k)} \sim \beta^{(k)}$; generate words $\mathbf{x}_n^{(k)} \sim Mult(\phi_{z_n^{(k)}}^{(k)})$.

The priors $\gamma^{(doc)}$, $\gamma^{(k)}$, $\alpha^{(doc)}$ and $\alpha^{(k)}$ can be estimated at learning time using non-informative hyperpriors (as we do in our experiments), or set manually to indicate preferences of segmentation granularity.

At inference time, we enforce each latent topic $z_n^{(k)}$ to be assigned to a contiguous span of text, assuming that coherent topics are not recurring across the document (Halliday and Hasan, 1976). It also reduces the search space and, consequently, speeds up our sampling-based inference by reducing the time needed for Monte Carlo chains to mix. In fact, this constraint can be integrated in the model definition but it would significantly complicate the model description.

4 Inference

As exact inference is intractable, we follow Eisenstein and Barzilay (2008) and instead use a Metropolis-Hastings (MH) algorithm. At each iteration of the MH algorithm, a new potential alignment-segmentation pair (z', t') is drawn from a proposal distribution $Q(z', t' | z, t)$, where (z, t)

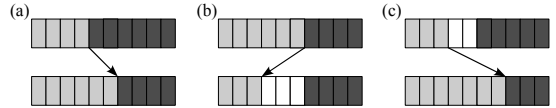


Figure 2: Three types of moves: (a) shift, (b) split and (c) merge.

is the current segmentation and its type. The new pair (z', t') is accepted with the probability

$$\min \left(1, \frac{P(z', t', \mathbf{x})Q(z, t | z', t')}{P(z, t, \mathbf{x})Q(z', t' | z, t)} \right).$$

In order to implement the MH algorithm for our model, we need to define the set of potential *moves* (i.e. admissible changes from (z, t) to (z', t')), and the proposal distribution Q over these moves. If the actual number of segments is known and only a linear discourse structure is acceptable, then a single move, *shift* of the segment border (Fig. 2(a)), is sufficient (Eisenstein and Barzilay, 2008). In our case, however, a more complex set of moves is required.

We make two assumptions which are motivated by the problem considered in Section 5: we assume that (1) we are given the number of document-level segments and also that (2) the aligned segments appear in the same order in each part of the document. With these assumptions in mind, we introduce two additional moves (Fig. 2(b) and (c)):

- *Split move*: select a segment, and split it at one of the spanned sentences; if the segment was a document-level segment then one of the fragments becomes the same document-level segment.
- *Merge move*: select a pair of adjacent segments where at least one of the segments is part-specific, and merge them; if one of them was a document-level segment then the new segment has the same document-level topic.

All the moves are selected with the uniform probability, and the distance c for the shift move is drawn from the proposal distribution proportional to $c^{-1/c_{max}}$. The moves are selected independently for each part.

Although the above two assumptions are not crucial as a simple modification to the set of moves would support both introduction and deletion of document-level fragments, this modification was not necessary for our experiments.

5 Experiment

5.1 Dataset and setup

Dataset We apply our model to the ESL podcast dataset (Noh et al., 2010) of 200 episodes, with an average of 17 sentences per story and 80 sentences per lecture transcript. The gold standard alignments assign each fragment of the story to a segment of the lecture transcript. We can induce segmentations at different levels of granularity on both the story and the lecture side. However, given that the segmentation of the story was obtained by an automatic sentence splitter, there is no reason to attempt to reproduce this segmentation. Therefore, for quantitative evaluation purposes we follow Noh et al. (2010) and restrict our model to alignment structures which agree with the given segmentation of the story. For all evaluations, we apply standard stemming algorithm and remove common stop words.

Evaluation metrics To measure the quality of segmentation of the lecture transcript, we use two standard metrics, P_k (Beeferman et al., 1999) and WindowDiff (WD) (Pevzner and Hearst, 2002), but both metrics disregard the alignment links (i.e. the topic labels). Consequently, we also use the macro-averaged F_1 score on pairs of aligned span, which measures both the segmentation and alignment quality.

Baseline Since there has been little previous research on this problem, we compare our results against two straightforward unsupervised baselines. For the first baseline, we consider the pairwise sentence alignment (SentAlign) based on the unigram and bigram overlap. The second baseline is a pipeline approach (Pipeline), where we first segment the lecture transcript with BayesSeg (Eisenstein and Barzilay, 2008) and then use the pairwise alignment to find their best alignment to the segments of the story.

Our model We evaluate our joint model of segmentation and alignment both with and without the split/merge moves. For the model without these moves, we set the desired number of segments in the lecture to be equal to the actual number of segments in the story I . In this setting, the moves can only adjust positions of the segment borders. For the model with the split/merge moves, we start with the same number of segments I but it can be increased or decreased during inference. For evaluation of our model, we run our inference algorithm from five random states, and

Method	P_k	WD	$1 - F_1$
Uniform	0.453	0.458	0.682
SentAlign	0.446	0.547	0.313
Pipeline (I)	0.250	0.249	0.443
Pipeline ($2I+1$)	0.268	0.289	0.318
Our model (I)	0.193	0.204	0.254
+split/merge	0.181	0.193	0.239

Table 1: Results on the ESL podcast dataset. For all metrics, lower values are better.

take the 100,000th iteration of each chain as a sample. Results are the average over these five runs. Also we perform L-BFGS optimization to automatically adjust the non-informative hyperpriors after each 1,000 iterations of sampling.

5.2 Result

Table 1 summarizes the obtained results. ‘Uniform’ denotes the minimal baseline which uniformly draws a random set of I spans for each lecture, and then aligns them to the segments of the story preserving the linear order. Also, we consider two variants of the pipeline approach: segmenting the lecture on I and $2I + 1$ segments, respectively.³ Our joint model substantially outperforms the baselines. The difference is statistically significant with the level $p < .01$ measured with the paired t-test. The significant improvement over the pipeline results demonstrates benefits of joint modeling for the considered problem. Moreover, additional benefits are obtained by using the DP priors and the split/merge moves (the last line in Table 1). Finally, our model significantly outperforms the previously proposed supervised model (Noh et al., 2010): they report micro-averaged F_1 score 0.698 while our best model achieves 0.778 with the same metric. This observation confirms that lexical cohesion modeling is crucial for successful discourse analysis.

6 Conclusions

We studied the problem of joint discourse segmentation and alignment of documents with inherently parallel structure and achieved favorable results on the ESL podcast dataset outperforming the cascaded baselines. Accurate prediction of these hidden relations would open interesting possibilities

³The use of the DP priors and the split/merge moves on the first stage of the pipeline did not result in any improvement in accuracy.

for construction of friendlier user interfaces. One example being an application which, given a user-selected fragment of the abstract, produces a summary from the aligned segment of the document body.

Acknowledgment

The authors acknowledge the support of the Excellence Cluster on Multimodal Computing and Interaction (MMCI), and also thank Mikhail Kozhevnikov and the anonymous reviewers for their valuable comments, and Hyungjong Noh for providing their data.

References

- Doug Beeferman, Adam Berger, and John Lafferty. 1999. Statistical models for text segmentation. *Computational Linguistics*, 34(1–3):177–210.
- David M. Blei, Andrew Ng, and Michael I. Jordan. 2003. Latent dirichlet allocation. *JMLR*, 3:993–1022.
- Hal Daumé and Daniel Marcu. 2004. A phrase-based hmm approach to document/abstract alignment. In *Proceedings of EMNLP*, pages 137–144.
- Hal Daumé and Daniel Marcu. 2006. Bayesian query-focused summarization. In *Proceedings of ACL*, pages 305–312.
- Jacob Eisenstein and Regina Barzilay. 2008. Bayesian unsupervised topic segmentation. In *Proceedings of EMNLP*, pages 334–343.
- Thomas S. Ferguson. 1973. A Bayesian analysis of some non-parametric problems. *Annals of Statistics*, 1:209–230.
- Michel Galley, Kathleen R. McKeown, Eric Fosler-Lussier, and Hongyan Jing. 2003. Discourse segmentation of multi-party conversation. In *Proceedings of ACL*, pages 562–569.
- M. A. K. Halliday and Ruqaiya Hasan. 1976. *Cohesion in English*. Longman.
- Marti Hearst. 1994. Multi-paragraph segmentation of expository text. In *Proceedings of ACL*, pages 9–16.
- Hongyan Jing. 2002. Using hidden Markov modeling to decompose human-written summaries. *Computational Linguistics*, 28(4):527–543.
- Igor Malioutov and Regina Barzilay. 2006. Minimum cut model for spoken lecture segmentation. In *Proceedings of ACL*, pages 25–32.
- Daniel Marcu. 1999. The automatic construction of large-scale corpora for summarization research. In *Proceedings of ACM SIGIR*, pages 137–144.
- Hyungjong Noh, Minwoo Jeong, Sungjin Lee, Jonghoon Lee, and Gary Geunbae Lee. 2010. Script-description pair extraction from text documents of English as second language podcast. In *Proceedings of the 2nd International Conference on Computer Supported Education*.
- Lev Pevzner and Marti Hearst. 2002. A critique and improvement of an evaluation metric for text segmentation. *Computational Linguistics*, 28(1):19–36.
- Jayaram Sethuraman. 1994. A constructive definition of Dirichlet priors. *Statistica Sinica*, 4:639–650.
- Bingjun Sun, Prasenjit Mitra, C. Lee Giles, John Yen, and Hongyuan Zha. 2007. Topic segmentation with shared topic detection and alignment of multiple documents. In *Proceedings of ACM SIGIR*, pages 199–206.
- Masao Utiyama and Hitoshi Isahara. 2001. A statistical model for domain-independent text segmentation. In *Proceedings of ACL*, pages 491–498.
- Kai Yu, Shipeng Yu, and Vokler Tresp. 2005. Dirichlet enhanced latent semantic analysis. In *Proceedings of AISTATS*.

Coreference Resolution with Reconcile

Veselin Stoyanov
Center for Language
and Speech Processing
Johns Hopkins Univ.
Baltimore, MD
ves@cs.jhu.edu

Claire Cardie
Department of
Computer Science
Cornell University
Ithaca, NY
cardie@cs.cornell.edu

Nathan Gilbert
Ellen Riloff
School of Computing
University of Utah
Salt Lake City, UT
ngilbert@cs.utah.edu
riloff@cs.utah.edu

David Buttler
David Hysom
Lawrence Livermore
National Laboratory
Livermore, CA
buttler1@llnl.gov
hysom1@llnl.gov

Abstract

Despite the existence of several noun phrase coreference resolution data sets as well as several formal evaluations on the task, it remains frustratingly difficult to compare results across different coreference resolution systems. This is due to the high cost of implementing a complete end-to-end coreference resolution system, which often forces researchers to substitute available gold-standard information in lieu of implementing a module that would compute that information. Unfortunately, this leads to inconsistent and often unrealistic evaluation scenarios.

With the aim to facilitate consistent and realistic experimental evaluations in coreference resolution, we present Reconcile, an infrastructure for the development of learning-based noun phrase (NP) coreference resolution systems. Reconcile is designed to facilitate the rapid creation of coreference resolution systems, easy implementation of new feature sets and approaches to coreference resolution, and empirical evaluation of coreference resolvers across a variety of benchmark data sets and standard scoring metrics. We describe Reconcile and present experimental results showing that Reconcile can be used to create a coreference resolver that achieves performance comparable to state-of-the-art systems on six benchmark data sets.

1 Introduction

Noun phrase coreference resolution (or simply coreference resolution) is the problem of identifying all noun phrases (NPs) that refer to the same entity in a text. The problem of coreference resolution is fundamental in the field of natural language processing (NLP) because of its usefulness for other NLP tasks, as well as the theoretical interest in understanding the computational mechanisms involved in government, binding and linguistic reference.

Several formal evaluations have been conducted for the coreference resolution task (e.g., MUC-6 (1995), ACE NIST (2004)), and the data sets created for these evaluations have become standard benchmarks in the field (e.g., MUC and ACE data sets). However, it is still frustratingly difficult to compare results across different coreference resolution systems. Reported coreference resolution scores vary wildly across data sets, evaluation metrics, and system configurations.

We believe that one root cause of these disparities is the high cost of implementing an end-to-end coreference resolution system. Coreference resolution is a complex problem, and successful systems must tackle a variety of non-trivial subproblems that are central to the coreference task — e.g., mention/markable detection, anaphor identification — and that require substantial implementation efforts. As a result, many researchers exploit gold-standard annotations, when available, as a substitute for component technologies to solve these subproblems. For example, many published research results use gold standard annotations to identify NPs (substituting for mention/markable detection), to distinguish anaphoric NPs from non-anaphoric NPs (substituting for anaphoricity determination), to identify named entities (substituting for named entity recognition), and to identify the semantic types of NPs (substituting for semantic class identification). Unfortunately, the use of gold standard annotations for key/critical component technologies leads to an unrealistic evaluation setting, and makes it impossible to directly compare results against coreference resolvers that solve all of these subproblems from scratch.

Comparison of coreference resolvers is further hindered by the use of several competing (and non-trivial) evaluation measures, and data sets that have substantially different task definitions and annotation formats. Additionally, coreference resolution is a pervasive problem in NLP and many NLP applications could benefit from an effective coreference resolver that can be easily configured and customized.

To address these issues, we have created a platform for coreference resolution, called Reconcile, that can serve as a software infrastructure to support the creation of, experimentation with, and evaluation of coreference resolvers. Reconcile was designed with the following seven desiderata in mind:

- implement the basic underlying software ar-

chitecture of contemporary state-of-the-art learning-based coreference resolution systems;

- support experimentation on most of the standard coreference resolution data sets;
- implement most popular coreference resolution scoring metrics;
- exhibit state-of-the-art coreference resolution performance (i.e., it can be configured to create a resolver that achieves performance close to the best reported results);
- can be easily extended with new methods and features;
- is relatively fast and easy to configure and run;
- has a set of pre-built resolvers that can be used as black-box coreference resolution systems.

While several other coreference resolution systems are publicly available (e.g., Poesio and Kabadjov (2004), Qiu et al. (2004) and Versley et al. (2008)), none meets all seven of these desiderata (see Related Work). Reconcile is a modular software platform that abstracts the basic architecture of most contemporary supervised learning-based coreference resolution systems (e.g., Soon et al. (2001), Ng and Cardie (2002), Bengtson and Roth (2008)) and achieves performance comparable to the state-of-the-art on several benchmark data sets. Additionally, Reconcile can be easily reconfigured to use different algorithms, features, preprocessing elements, evaluation settings and metrics.

In the rest of this paper, we review related work (Section 2), describe Reconcile’s organization and components (Section 3) and show experimental results for Reconcile on six data sets and two evaluation metrics (Section 4).

2 Related Work

Several coreference resolution systems are currently publicly available. JavaRap (Qiu et al., 2004) is an implementation of the Lappin and Leass’ (1994) Resolution of Anaphora Procedure (RAP). JavaRap resolves only pronouns and, thus, it is not directly comparable to Reconcile. GuiTaR

(Poesio and Kabadjov, 2004) and BART (Versley et al., 2008) (which can be considered a successor of GuiTaR) are both modular systems that target the full coreference resolution task. As such, both systems come close to meeting the majority of the desiderata set forth in Section 1. BART, in particular, can be considered an alternative to Reconcile, although we believe that Reconcile’s approach is more flexible than BART’s. In addition, the architecture and system components of Reconcile (including a comprehensive set of features that draw on the expertise of state-of-the-art supervised learning approaches, such as Bengtson and Roth (2008)) result in performance closer to the state-of-the-art.

Coreference resolution has received much research attention, resulting in an array of approaches, algorithms and features. Reconcile is modeled after typical supervised learning approaches to coreference resolution (e.g. the architecture introduced by Soon et al. (2001)) because of the popularity and relatively good performance of these systems.

However, there have been other approaches to coreference resolution, including unsupervised and semi-supervised approaches (e.g. Haghighi and Klein (2007)), structured approaches (e.g. McCallum and Wellner (2004) and Finley and Joachims (2005)), competition approaches (e.g. Yang et al. (2003)) and a bell-tree search approach (Luo et al. (2004)). Most of these approaches rely on some notion of pairwise feature-based similarity and can be directly implemented in Reconcile.

3 System Description

Reconcile was designed to be a research testbed capable of implementing most current approaches to coreference resolution. Reconcile is written in Java, to be portable across platforms, and was designed to be easily reconfigurable with respect to subcomponents, feature sets, parameter settings, etc.

Reconcile’s architecture is illustrated in Figure 1. For simplicity, Figure 1 shows Reconcile’s operation during the classification phase (i.e., assuming that a trained classifier is present).

The basic architecture of the system includes five major steps. Starting with a corpus of documents together with a manually annotated coreference resolution answer key¹, Reconcile performs

¹Only required during training.

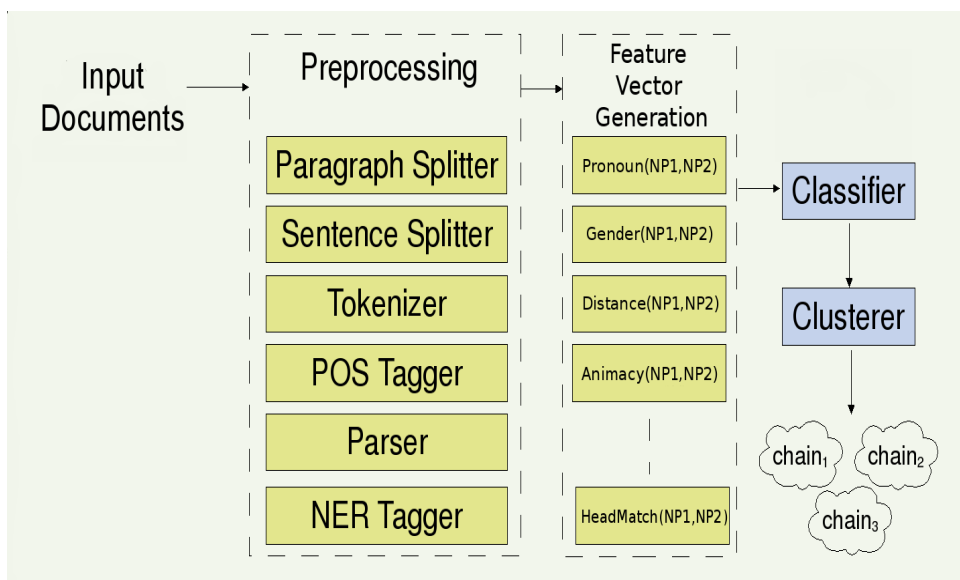


Figure 1: The Reconcile classification architecture.

the following steps, in order:

1. **Preprocessing.** All documents are passed through a series of (external) linguistic processors such as tokenizers, part-of-speech taggers, syntactic parsers, etc. These components produce annotations of the text. Table 1 lists the preprocessors currently interfaced in Reconcile. Note that Reconcile includes several in-house NP detectors, that conform to the different data sets' definitions of what constitutes a NP (e.g., MUC vs. ACE). All of the extractors utilize a syntactic parse of the text and the output of a Named Entity (NE) extractor, but extract different constructs as specialized in the corresponding definition. The NP extractors successfully recognize about 95% of the NPs in the MUC and ACE gold standards.
2. **Feature generation.** Using annotations produced during preprocessing, Reconcile produces feature vectors for pairs of NPs. For example, a feature might denote whether the two NPs agree in number, or whether they have any words in common. Reconcile includes over 80 features, inspired by other successful coreference resolution systems such as Soon et al. (2001) and Ng and Cardie (2002).
3. **Classification.** Reconcile learns a classifier that operates on feature vectors representing

Task	Systems
Sentence splitter	UIUC (CC Group, 2009) OpenNLP (Baldrige, J., 2005)
Tokenizer	OpenNLP (Baldrige, J., 2005)
POS Tagger	OpenNLP (Baldrige, J., 2005) + the two parsers below
Parser	Stanford (Klein and Manning, 2003) Berkeley (Petrov and Klein, 2007)
Dep. parser	Stanford (Klein and Manning, 2003)
NE Recognizer	OpenNLP (Baldrige, J., 2005) Stanford (Finkel et al., 2005)
NP Detector	In-house

Table 1: Preprocessing components available in Reconcile.

pairs of NPs and it is trained to assign a score indicating the likelihood that the NPs in the pair are coreferent.

4. **Clustering.** A clustering algorithm consolidates the predictions output by the classifier and forms the final set of coreference clusters (chains).²
5. **Scoring.** Finally, during testing Reconcile runs scoring algorithms that compare the chains produced by the system to the gold-standard chains in the answer key.

Each of the five steps above can invoke different components. Reconcile's modularity makes it

²Some structured coreference resolution algorithms (e.g., McCallum and Wellner (2004) and Finley and Joachims (2005)) combine the classification and clustering steps above. Reconcile can easily accommodate this modification.

Step	Available modules
Classification	various learners in the Weka toolkit libSVM (Chang and Lin, 2001) SVM _{light} (Joachims, 2002)
Clustering	Single-link Best-First Most Recent First
Scoring	MUC score (Vilain et al., 1995) B^3 score (Bagga and Baldwin, 1998) CEAF score (Luo, 2005)

Table 2: Available implementations for different modules available in Reconcile.

easy for new components to be implemented and existing ones to be removed or replaced. Reconcile’s standard distribution comes with a comprehensive set of implemented components – those available for steps 2–5 are shown in Table 2. Reconcile contains over 38,000 lines of original Java code. Only about 15% of the code is concerned with running existing components in the preprocessing step, while the rest deals with NP extraction, implementations of features, clustering algorithms and scorers. More details about Reconcile’s architecture and available components and features can be found in Stoyanov et al. (2010).

4 Evaluation

4.1 Data Sets

Reconcile incorporates the six most commonly used coreference resolution data sets, two from the MUC conferences (MUC-6, 1995; MUC-7, 1997) and four from the ACE Program (NIST, 2004). For ACE, we incorporate only the newswire portion. When available, Reconcile employs the standard test/train split. Otherwise, we randomly split the data into a training and test set following a 70/30 ratio. Performance is evaluated according to the B^3 and MUC scoring metrics.

4.2 The *Reconcile*₂₀₁₀ Configuration

Reconcile can be easily configured with different algorithms for markable detection, anaphoricity determination, feature extraction, etc., and run against several scoring metrics. For the purpose of this sample evaluation, we create only one particular instantiation of Reconcile, which we will call *Reconcile*₂₀₁₀ to differentiate it from the general platform. *Reconcile*₂₀₁₀ is configured using the following components:

1. **Preprocessing**
 - (a) **Sentence Splitter:** *OpenNLP*

- (b) **Tokenizer:** *OpenNLP*
- (c) **POS Tagger:** *OpenNLP*
- (d) **Parser:** *Berkeley*
- (e) **Named Entity Recognizer:** *Stanford*

2. **Feature Set** - A hand-selected subset of 60 out of the more than 80 features available. The features were selected to include most of the features from Soon et al. (2001), Ng and Cardie (2002) and Bengtson and Roth (2008).
3. **Classifier** - *Averaged Perceptron*
4. **Clustering** - *Single-link* - Positive decision threshold was tuned by cross validation of the training set.

4.3 Experimental Results

The first two rows of Table 3 show the performance of *Reconcile*₂₀₁₀. For all data sets, B^3 scores are higher than MUC scores. The MUC score is highest for the MUC6 data set, while B^3 scores are higher for the ACE data sets as compared to the MUC data sets.

Due to the difficulties outlined in Section 1, results for Reconcile presented here are directly comparable only to a limited number of scores reported in the literature. The bottom three rows of Table 3 list these comparable scores, which show that *Reconcile*₂₀₁₀ exhibits state-of-the-art performance for supervised learning-based coreference resolvers. A more detailed study of Reconcile-based coreference resolution systems in different evaluation scenarios can be found in Stoyanov et al. (2009).

5 Conclusions

Reconcile is a general architecture for coreference resolution that can be used to easily create various coreference resolvers. Reconcile provides broad support for experimentation in coreference resolution, including implementation of the basic architecture of contemporary state-of-the-art coreference systems and a variety of individual modules employed in these systems. Additionally, Reconcile handles all of the formatting and scoring peculiarities of the most widely used coreference resolution data sets (those created as part of the MUC and ACE conferences) and, thus, allows for easy implementation and evaluation across these data sets. We hope that Reconcile will support experimental research in coreference resolution and provide a state-of-the-art coreference resolver for both researchers and application developers. We believe that in this way Reconcile will facilitate meaningful and consistent comparisons of coreference resolution systems. The full Reconcile release is available for download at <http://www.cs.utah.edu/nlp/reconcile/>.

System	Score	Data sets					
		MUC6	MUC7	ACE-2	ACE03	ACE04	ACE05
<i>Reconcile</i> ₂₀₁₀	<i>MUC</i>	68.50	62.80	65.99	67.87	62.03	67.41
	<i>B</i> ³	70.88	65.86	78.29	79.39	76.50	73.71
Soon et al. (2001)	<i>MUC</i>	62.6	60.4	–	–	–	–
Ng and Cardie (2002)	<i>MUC</i>	70.4	63.4	–	–	–	–
Yang et al. (2003)	<i>MUC</i>	71.3	60.2	–	–	–	–

Table 3: Scores for Reconcile on six data sets and scores for comparable coreference systems.

Acknowledgments

This research was supported in part by the National Science Foundation under Grant # 0937060 to the Computing Research Association for the CIFellows Project, Lawrence Livermore National Laboratory subcontract B573245, Department of Homeland Security Grant N0014-07-1-0152, and Air Force Contract FA8750-09-C-0172 under the DARPA Machine Reading Program.

The authors would like to thank the anonymous reviewers for their useful comments.

References

- A. Bagga and B. Baldwin. 1998. Algorithms for scoring coreference chains. In *Linguistic Coreference Workshop at the Language Resources and Evaluation Conference*.
- Baldrige, J. 2005. The OpenNLP project. <http://opennlp.sourceforge.net/>.
- E. Bengtson and D. Roth. 2008. Understanding the value of features for coreference resolution. In *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- CC Group. 2009. Sentence Segmentation Tool. <http://l2r.cs.uiuc.edu/cogcomp/atool.php?tkey=SS>.
- C. Chang and C. Lin. 2001. LIBSVM: a Library for Support Vector Machines. Available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- J. Finkel, T. Grenager, and C. Manning. 2005. Incorporating Non-local Information into Information Extraction Systems by Gibbs Sampling. In *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the ACL*.
- T. Finley and T. Joachims. 2005. Supervised clustering with support vector machines. In *Proceedings of the Twenty-second International Conference on Machine Learning (ICML 2005)*.
- A. Haghighi and D. Klein. 2007. Unsupervised Coreference Resolution in a Nonparametric Bayesian Model. In *Proceedings of the 45th Annual Meeting of the ACL*.
- T. Joachims. 2002. SVM_{Light}, <http://svmlight.joachims.org>.
- D. Klein and C. Manning. 2003. Fast Exact Inference with a Factored Model for Natural Language Parsing. In *Advances in Neural Information Processing (NIPS 2003)*.
- S. Lappin and H. Leass. 1994. An algorithm for pronominal anaphora resolution. *Computational Linguistics*, 20(4):535–561.
- X. Luo, A. Ittycheriah, H. Jing, N. Kambhatla, and S. Roukos. 2004. A mention-synchronous coreference resolution algorithm based on the bell tree. In *Proceedings of the 42nd Annual Meeting of the ACL*.
- X. Luo. 2005. On Coreference Resolution Performance Metrics. In *Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing (HLT/EMNLP)*.
- A. McCallum and B. Wellner. 2004. Conditional Models of Identity Uncertainty with Application to Noun Coreference. In *Advances in Neural Information Processing (NIPS 2004)*.
- MUC-6. 1995. Coreference Task Definition. In *Proceedings of the Sixth Message Understanding Conference (MUC-6)*.
- MUC-7. 1997. Coreference Task Definition. In *Proceedings of the Seventh Message Understanding Conference (MUC-7)*.
- V. Ng and C. Cardie. 2002. Improving Machine Learning Approaches to Coreference Resolution. In *Proceedings of the 40th Annual Meeting of the ACL*.
- NIST. 2004. *The ACE Evaluation Plan*. NIST.
- S. Petrov and D. Klein. 2007. Improved Inference for Unlexicalized Parsing. In *Proceedings of the Joint Meeting of the Human Language Technology Conference and the North American Chapter of the Association for Computational Linguistics (HLT-NAACL 2007)*.
- M. Poesio and M. Kabadjov. 2004. A general-purpose, off-the-shelf anaphora resolution module: implementation and preliminary evaluation. In *Proceedings of the Language Resources and Evaluation Conference*.
- L. Qiu, M.-Y. Kan, and T.-S. Chua. 2004. A public reference implementation of the rap anaphora resolution algorithm. In *Proceedings of the Language Resources and Evaluation Conference*.
- W. Soon, H. Ng, and D. Lim. 2001. A Machine Learning Approach to Coreference of Noun Phrases. *Computational Linguistics*, 27(4):521–541.
- V. Stoyanov, N. Gilbert, C. Cardie, and E. Riloff. 2009. Conundrums in noun phrase coreference resolution: Making sense of the state-of-the-art. In *Proceedings of ACL/IJCNLP*.

- V. Stoyanov, C. Cardie, N. Gilbert, E. Riloff, D. Buttler, and D. Hysom. 2010. Reconcile: A coreference resolution research platform. Technical report, Cornell University.
- Y. Versley, S. Ponzetto, M. Poesio, V. Eidelman, A. Jern, J. Smith, X. Yang, and A. Moschitti. 2008. BART: A modular toolkit for coreference resolution. In *Proceedings of the Language Resources and Evaluation Conference*.
- M. Vilain, J. Burger, J. Aberdeen, D. Connolly, and L. Hirschman. 1995. A Model-Theoretic Coreference Scoring Theme. In *Proceedings of the Sixth Message Understanding Conference (MUC-6)*.
- X. Yang, G. Zhou, J. Su, and C. Tan. 2003. Coreference resolution using competition learning approach. In *Proceedings of the 41st Annual Meeting of the ACL*.

Predicate Argument Structure Analysis using Transformation-based Learning

Hirotohi Taira

Sanae Fujita

Masaaki Nagata

NTT Communication Science Laboratories

2-4, Hikaridai, Seika-cho, Souraku-gun, Kyoto 619-0237, Japan

{taira, sanae}@cslab.kecl.ntt.co.jp nagata.masaaki@lab.ntt.co.jp

Abstract

Maintaining high annotation consistency in large corpora is crucial for statistical learning; however, such work is hard, especially for tasks containing semantic elements. This paper describes predicate argument structure analysis using transformation-based learning. An advantage of transformation-based learning is the readability of learned rules. A disadvantage is that the rule extraction procedure is time-consuming. We present incremental-based, transformation-based learning for semantic processing tasks. As an example, we deal with Japanese predicate argument analysis and show some tendencies of annotators for constructing a corpus with our method.

1 Introduction

Automatic predicate argument structure analysis (PAS) provides information of “who did what to whom” and is an important base tool for such various text processing tasks as machine translation information extraction (Hirschman et al., 1999), question answering (Narayanan and Harabagiu, 2004; Shen and Lapata, 2007), and summarization (Melli et al., 2005). Most recent approaches to predicate argument structure analysis are statistical machine learning methods such as support vector machines (SVMs) (Pradhan et al., 2004). For predicate argument structure analysis, we have the following representative large corpora: FrameNet (Fillmore et al., 2001), PropBank (Palmer et al., 2005), and NomBank (Meyers et al., 2004) in English, the Chinese PropBank (Xue, 2008) in Chinese, the GDA Corpus (Hashida, 2005), Kyoto Text Corpus Ver.4.0 (Kawahara et al., 2002), and the NAIST Text Corpus (Iida et al., 2007) in Japanese.

The construction of such large corpora is strenuous and time-consuming. Additionally, maintaining high annotation consistency in such corpora is crucial for statistical learning; however, such work is hard, especially for tasks containing semantic elements. For example, in Japanese corpora, distinguishing true dative (or indirect object) arguments from time-type argument is difficult because the arguments of both types are often accompanied with the ‘ni’ case marker.

A problem with such statistical learners as SVM is the lack of interpretability; if accuracy is low, we cannot identify the problems in the annotations.

We are focusing on transformation-based learning (TBL). An advantage for such learning methods is that we can easily interpret the learned model. The tasks in most previous research are such simple tagging tasks as part-of-speech tagging, insertion and deletion of parentheses in syntactic parsing, and chunking (Brill, 1995; Brill, 1993; Ramshaw and Marcus, 1995). Here we experiment with a complex task: Japanese PASs. TBL can be slow, so we proposed an incremental training method to speed up the training. We experimented with a Japanese PAS corpus with a graph-based TBL. From the experiments, we interrelated the annotation tendency on the dataset.

The rest of this paper is organized as follows. Section 2 describes Japanese predicate structure, our graph expression of it, and our improved method. The results of experiments using the NAIST Text Corpus, which is our target corpus, are reported in Section 3, and our conclusion is provided in Section 4.

2 Predicate argument structure and graph transformation learning

First, we illustrate the structure of a Japanese sentence in Fig. 1. In Japanese, we can divide a sentence into *bunsetsu* phrases (BP). A BP usually consists of one or more content words and zero,

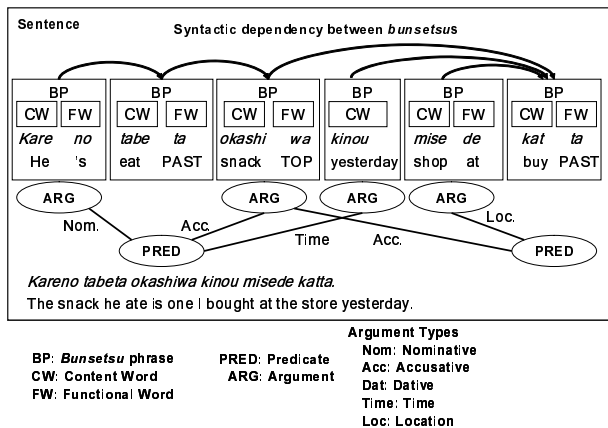


Figure 1: Graph expression for PAS

one, or more than one functional words. Syntactic dependency between *bunsetsu* phrases can be defined. Japanese dependency parsers such as Cabocha (Kudo and Matsumoto, 2002) can extract BPs and their dependencies with about 90% accuracy.

Since predicates and arguments in Japanese are mainly annotated on the head content word in each BP, we can deal with BPs as candidates of predicates or arguments. In our experiments, we mapped each BP to an argument candidate node of graphs. We also mapped each predicate to a predicate node. Each predicate-argument relation is identified by an edge between a predicate and an argument, and the argument type is mapped to the edge label. In our experiments below, we defined five argument types: nominative (subjective), accusative (direct objective), dative (indirect objective), time, and location. We use five transformation types: a) add or b) delete a predicate node, c) add or d) delete an edge between a predicate and an argument node, e) change a label (= an argument type) to another label (Fig. 2). We explain the existence of an edge between a predicate and an argument labeled t candidate node as that the predicate and the argument have a t type relationship.

Transformation-based learning was proposed by (Brill, 1995). Below we explain our learning strategy when we directly adapt the learning method to our graph expression of PASs. First, unstructured texts from the training data are inputted. After pre-processing, each text is mapped to an initial graph. In our experiments, the initial graph has argument candidate nodes with corresponding BPs and no predicate nodes or edges. Next, com-

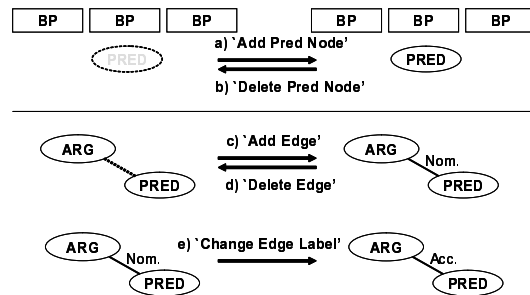


Figure 2: Transform types

paring the current graphs with the gold standard graph structure in the training data, we find the different statuses of the nodes and edges among the graphs. We extract such transformation rule candidates as 'add node' and 'change edge label' with constraints, including 'the corresponding BP includes a verb' and 'the argument candidate and the predicate node have a syntactic dependency.' The extractions are executed based on the rule templates given in advance. Each extracted rule is evaluated for the current graphs, and error reduction is calculated. The best rule for the reduction is selected as a new rule and inserted at the bottom of the current rule list. The new rule is applied to the current graphs, which are transferred to other graph structures. This procedure is iterated until the total errors for the gold standard graphs become zero. When the process is completed, the rule list is the final model. In the test phase, we iteratively transform nodes and edges in the graphs mapped from the test data, based on rules in the model like decision lists. The last graph after all rule adaptations is the system output of the PAS.

In this procedure, the calculation of error reduction is very time-consuming, because we have to check many constraints from the candidate rules for all training samples. The calculation order is $O(MN)$, where M is the number of articles and N is the number of candidate rules. Additionally, an edge rule usually has three types of constraints: 'pred node constraint,' 'argument candidate node constraint,' and 'relation constraint.' The number of combinations and extracted rules are much larger than one of the rules for the node rules. Ramshaw et al. proposed an index-based efficient reduction method for the calculation of error reduction (Ramshaw and Marcus, 1994). However, in PAS tasks, we need to check the exclusiveness of the argument types (for example, a predicate argument structure does not have two nominative ar-

guments), and we cannot directly use the method. Jijkoun et al. only used candidate rules that happen in the current and gold standard graphs and used SVM learning for constraint checks (Jijkoun and de Rijke, 2007). This method is effective for achieving high accuracy; however, it loses the readability of the rules. This is contrary to our aim to extract readable rules.

To reduce the calculations while maintaining readability, we propose an incremental method and describe its procedure below. In this procedure, we first have PAS graphs for only one article. After the total errors among the current and gold standard graphs become zero in the article, we proceed to the next article. For the next article, we first adapt the rules learned from the previous article. After that, we extract new rules from the two articles until the total errors for the articles become zero. We continue these processes until the last article. Additionally, we count the number of rule occurrences and only use the rule candidates that happen more than once, because most such rules harm the accuracy. We save and use these rules again if the occurrence increases.

3 Experiments

3.1 Experimental Settings

We used the articles in the NAIST Text Corpus version 1.4 β (Iida et al., 2007) based on the *Mainichi Shinbun* Corpus (Mainichi, 1995), which were taken from news articles published in the Japanese *Mainichi Shinbun* newspaper. We used articles published on January 1st for training examples and on January 3rd for test examples. Three original argument types are defined in the NAIST Text Corpus: nominative (or subjective), accusative (or direct object), and dative (or indirect object). For evaluation of the difficult annotation cases, we also added annotations for ‘time’ and ‘location’ types by ourselves. We show the dataset distribution in Table 1. We extracted the BP units and dependencies among these BPs from the dataset using Cabocha, a Japanese dependency parser, as pre-processing. After that, we adapted our incremental learning to the training data. We used two constraint templates in Tables 2 and 3 for predicate nodes and edges when extracting the rule candidates.

Table 1: Data distribution

	Training	Test
# of Articles	95	74
# of Sentences	1,129	687
# of Predicates	3,261	2,038
# of Arguments	3,877	2,468
Nom.	1,717	971
Acc.	1,012	701
Dat.	632	376
Time	371	295
Loc.	145	125

Table 4: Total performances (F1-measure (%))

Type	System	P	R	F1
Pred.	Baseline	89.4	85.1	87.2
	Our system	91.8	85.3	88.4
Arg.	Baseline	79.3	59.5	68.0
	Our system	81.9	62.4	70.8

3.2 Results

Our incremental method takes an hour. In comparison, the original TBL cannot even extract one rule in a day. The results of predicate and argument type predictions are shown in Table 4. Here, ‘Baseline’ is the baseline system that predicts the BSs that contain verbs, adjectives, and *da* form nouns (‘to be’ in English) as predicates and predicts argument types for BSs having syntactical dependency with a predicted predicate BS, based on the following rules: 1) BSs containing nominative (*ga*) / accusative (*wo*) / dative (*ni*) case markers are predicted to be nominative, accusative, and dative, respectively. 2) BSs containing a topic case marker (*wa*) are predicted to be nominative. 3) When a word sense category from a Japanese ontology of the head word in BS belongs to a ‘time’ or ‘location’ category, the BS is predicted to be a ‘time’ and ‘location’ type argument. In all precision, recall, and F1-measure, our system outperformed the baseline system.

Next, we show our system’s learning curve in Fig. 3. The number of final rules was 68. This indicates that the first twenty rules are mainly effective rules for the performance. The curve also shows that no overfitting happened. Next, we show the performance for every argument type in Table 5. ‘TBL,’ which stands for ‘transformation-based learning,’ is our system. In this table, the performance of the dative and time types improved, even though they are difficult to distinguish. On the other hand, the performance of the location type argument in our system is very low. Our method learns rules as decreasing errors of

Table 2: Predicate node constraint templates

Pred. Node Constraint Template		Rule Example	
Constraint	Description	Pred. Node Constraint	Operation
pos1 pos2 pos1 & pos2 'da' lemma	noun, verb, adjective, etc. independent, attached word, etc. above two features combination <i>da</i> form (copula) word base form	pos1='ADJECTIVE' pos2='DEPENDENT WORD' pos1='VERB' & pos2='ANCILLARY WORD' 'da form' lemma='%'	add pred node del pred node add pred node add pred node add pred node

Table 3: Edge constraint templates

Edge Constraint Template			Rule Example	
Arg. Cand. Const.	Pred. Node Const.	Relation Const.	Edge Constraint	Operation
FW (=func. word)	*	dep(arg → pred)	FW of Arg. = ' <i>wa</i> (TOP)' & dep(arg → pred)	add NOM edge
*	FW	dep(arg ← pred)	FW of Pred. = ' <i>na</i> (ADNOMINAL)' & dep(arg ← pred)	add NOM edge
SemCat (=semantic category)	*	dep(arg → pred)	SemCat of Arg. = 'TIME' & dep(arg → pred)	add TIME edge
FW	passive form	dep(arg → pred)	FW of Arg. = ' <i>ga</i> (NOM) & Pred.: passive form	chg edge label NOM → ACC
*	kform (= type of inflected form)	*	kform of Pred. = continuative 'ta' form	add NOM edge
SemCat	Pred. SemCat	*	SemCat of Arg. = 'HUMAN' & Pred. SemCat = 'PHYSICAL MOVE'	add NOM edge

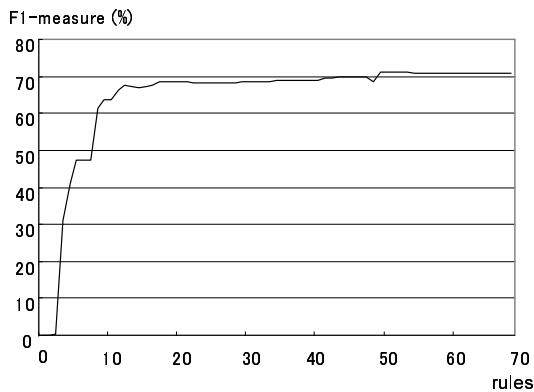


Figure 3: Learning curves: x-axis = number of rules; y-axis: F1-measure (%)

all arguments, and the performance of the location type argument is probably sacrificed for total error reduction because the number of location type arguments is much smaller than the number of other argument types (Table 1), and the improvement of the performance-based learning for location type arguments is relatively low. To confirm this, we performed an experiment in which we gave the rules of the baseline system to our system as initial rules and subsequently performed our incremental learning. 'Base + TBL' shows the experiment. The performance for the location type argument improved drastically. However, the total performance of the arguments was below the original TBL. Moreover, the 'Base + TBL' performance surpassed the baseline system. This indicates that our system learned a reasonable model.

Finally, we show some interesting extracted rules in Fig. 4. The first rule stands for an expression where the sentence ends with the performance of something, which is often seen in Japanese newspaper articles. The second and third rules represent that annotators of this dataset tend to annotate time types for which the semantic category of the argument is time, even if the argument looks like the *dat.* type, and annotators tend to annotate *dat.* type for arguments that have an *dat.*

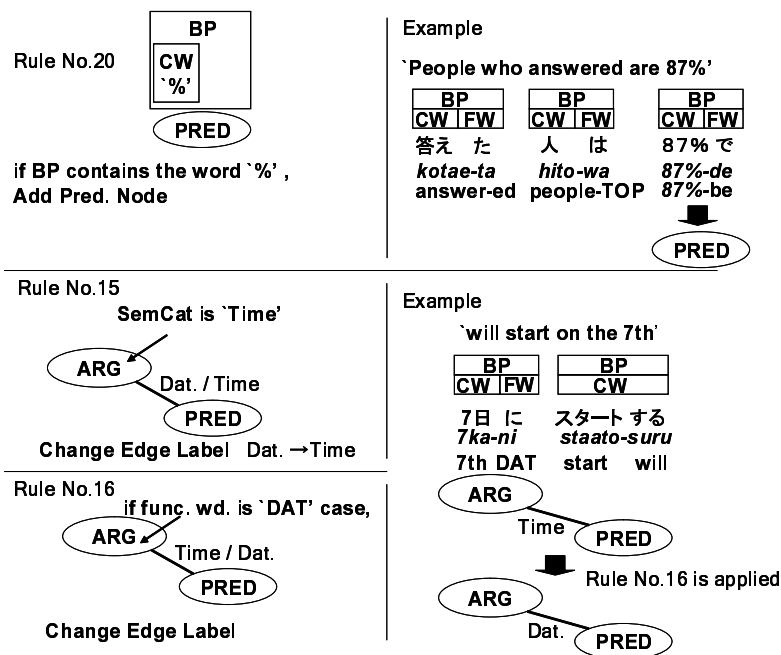


Figure 4: Examples of extracted rules

Table 5: Results for every arg. type (F-measure (%))

System	Args.	Nom.	Acc.	Dat.	Time	Loc.
Base	68.0	65.8	79.6	70.5	51.5	38.0
TBL	70.8	64.9	86.4	74.8	59.6	1.7
Base + TBL	69.5	63.9	85.8	67.8	55.8	37.4

type case marker.

4 Conclusion

We performed experiments for Japanese predicate argument structure analysis using transformation-based learning and extracted rules that indicate the tendencies annotators have. We presented an incremental procedure to speed up rule extraction. The performance of PAS analysis improved, especially, the dative and time types, which are difficult to distinguish. Moreover, when time expressions are attached to the ‘ni’ case, the learned model showed a tendency to annotate them as dative arguments in the used corpus. Our method has potential for dative predictions and interpreting the tendencies of annotator inconsistencies.

Acknowledgments

We thank Kevin Duh for his valuable comments.

References

- Eric Brill. 1993. Transformation-based error-driven parsing. In *Proc. of the Third International Workshop on Parsing Technologies*.
- Eric Brill. 1995. Transformation-based error-driven learning and natural language processing: A case study in part-of-speech tagging. *Computational Linguistics*, 21(4):543–565.
- Charles J. Fillmore, Charles Wooters, and Collin F. Baker. 2001. Building a large lexical databank which provides deep semantics. In *Proc. of the Pacific Asian Conference on Language, Information and Computation (PACLING)*.
- Kouichi Hashida. 2005. Global document annotation (GDA) manual. <http://i-content.org/GDA/>.
- Lynette Hirschman, Patricia Robinson, Lisa Ferro, Nancy Chinchor, Erica Brown, Ralph Grishman, and Beth Sundheim. 1999. Hub-4 Event’99 general guidelines. http://www.itl.nist.gov/iaui/894.02/related_projects/muc/.
- Ryu Iida, Mamoru Komachi, Kentaro Inui, and Yuji Matsumoto. 2007. Annotating a Japanese text corpus with predicate-argument and coreference relations. In *Proc. of ACL 2007 Workshop on Linguistic Annotation*, pages 132–139.
- Valentin Jijkoun and Maarten de Rijke. 2007. Learning to transform linguistic graphs. In *Proc. of the Second Workshop on TextGraphs: Graph-Based Algorithms for Natural Language Processing (TextGraphs-2)*, pages 53–60. Association for Computational Linguistics.

- Daisuke Kawahara, Sadao Kurohashi, and Koichi Hashida. 2002. Construction of a Japanese relevance-tagged corpus (in Japanese). *Proc. of the 8th Annual Meeting of the Association for Natural Language Processing*, pages 495–498.
- Taku Kudo and Yuji Matsumoto. 2002. Japanese dependency analysis using cascaded chunking. In *Proc. of the 6th Conference on Natural Language Learning 2002 (CoNLL 2002)*.
- Mainichi. 1995. *CD Mainichi Shinbun 94*. Nichigai Associates Co.
- Gabor Melli, Yang Wang, Yudong Liu, Mehdi M. Kashani, Zhongmin Shi, Baohua Gu, Anoop Sarkar, and Fred Popowich. 2005. Description of SQUASH, the SFU question answering summary handler for the DUC-2005 summarization task. In *Proc. of DUC 2005*.
- Adam Meyers, Ruth Reeves, Catherine Macleod, Rachel Szekely, Veronika Zielinska, Brian Young, and Ralph Grishman. 2004. The NomBank project: An interim report. In *Proc. of HLT-NAACL 2004 Workshop on Frontiers in Corpus Annotation*.
- Srini Narayanan and Sanda Harabagiu. 2004. Question answering based on semantic structures. In *Proc. of the 20th International Conference on Computational Linguistics (COLING)*.
- M. Palmer, P. Kingsbury, and D. Gildea. 2005. The proposition bank: An annotated corpus of semantic roles. *Computational Linguistics*, 31(1):71–106.
- Sameer Pradhan, Wayne Ward, Kadri Hacioglu, James Martin, and Dan Jurafsky. 2004. Shallow semantic parsing using support vector machines. In *Proc. of the Human Language Technology Conference/North American Chapter of the Association of Computational Linguistics HLT/NAACL 2004*.
- Lance Ramshaw and Mitchell Marcus. 1994. Exploring the statistical derivation of transformational rule sequences for part-of-speech tagging. In *The Balancing Act: Proc. of the ACL Workshop on Combining Symbolic and Statistical Approaches to Language*.
- Lance Ramshaw and Mitchell Marcus. 1995. Text chunking using transformation-based learning. In *Proc. of the third workshop on very large corpora*, pages 82–94.
- Dan Shen and Mirella Lapata. 2007. Using semantic roles to improve question answering. In *Proc. of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP/CoNLL)*, pages 12–21.
- Nianwen Xue. 2008. Labeling Chinese predicates with semantic roles. *Computational Linguistics*, 34(2):224–255.

Improving Chinese Semantic Role Labeling with Rich Syntactic Features

Weiwei Sun*

Department of Computational Linguistics, Saarland University
German Research Center for Artificial Intelligence (DFKI)
D-66123, Saarbrücken, Germany
wsun@coli.uni-saarland.de

Abstract

Developing features has been shown crucial to advancing the state-of-the-art in Semantic Role Labeling (SRL). To improve Chinese SRL, we propose a set of additional features, some of which are designed to better capture structural information. Our system achieves 93.49 F-measure, a significant improvement over the best reported performance 92.0. We are further concerned with the effect of parsing in Chinese SRL. We empirically analyze the two-fold effect, grouping words into constituents and providing syntactic information. We also give some preliminary linguistic explanations.

1 Introduction

Previous work on Chinese Semantic Role Labeling (SRL) mainly focused on how to implement SRL methods which are successful on English. Similar to English, parsing is a *standard* pre-processing for Chinese SRL. Many features are extracted to represent constituents in the input parses (Sun and Jurafsky, 2004; Xue, 2008; Ding and Chang, 2008). By using these features, semantic classifiers are trained to predict whether a constituent fills a semantic role. Developing features that capture the right kind of information encoded in the input parses has been shown crucial to advancing the state-of-the-art. Though there has been some work on feature design in Chinese SRL, information encoded in the syntactic trees is not fully exploited and requires more research effort. In this paper, we propose a set of additional

*The work was partially completed while this author was at Peking University.

features, some of which are designed to better capture structural information of sub-trees in a given parse. With help of these new features, our system achieves 93.49 F-measure with hand-crafted parses. Comparison with the best reported results, 92.0 (Xue, 2008), shows that these features yield a significant improvement of the state-of-the-art.

We further analyze the effect of syntactic parsing in Chinese SRL. The main effect of parsing in SRL is two-fold. First, grouping words into constituents, parsing helps to find argument candidates. Second, parsers provide semantic classifiers plenty of syntactic information, not to only recognize arguments from all candidate constituents but also to classify their detailed semantic types. We empirically analyze each effect in turn. We also give some preliminary linguistic explanations for the phenomena.

2 Chinese SRL

The Chinese PropBank (CPB) is a semantic annotation for the syntactic trees of the Chinese TreeBank (CTB). The arguments of a predicate are labeled with a contiguous sequence of integers, in the form of AN (N is a natural number); the adjuncts are annotated as such with the label AM followed by a secondary tag that represents the semantic classification of the adjunct. The assignment of semantic roles is illustrated in Figure 1, where the predicate is the verb “调查/investigate”. E.g., the NP “事故原因/the cause of the accident” is labeled as *A1*, meaning that it is the *Patient*.

In previous research, SRL methods that are successful on English are adopted to resolve Chinese SRL (Sun and Jurafsky, 2004; Xue, 2008; Ding and Chang, 2008, 2009; Sun et al., 2009; Sun, 2010). Xue (2008) produced complete and systematic research on full parsing based methods.

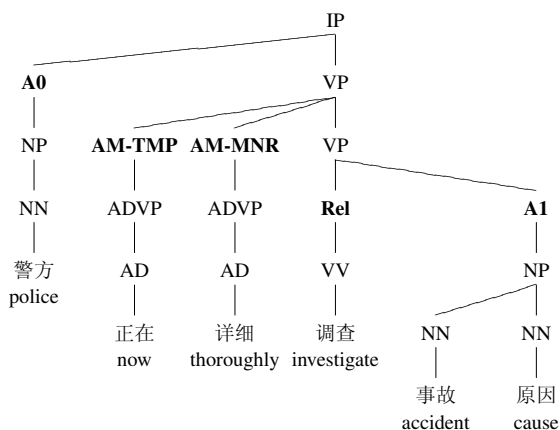


Figure 1: An example sentence: *The police are thoroughly investigating the cause of the accident.*

Their method divided SRL into three sub-tasks: 1) pruning with a heuristic rule, 2) Argument Identification (AI) to recognize arguments, and 3) Semantic Role Classification (SRC) to predict semantic types. The main two sub-tasks, AI and SRC, are formulated as two classification problems. Ding and Chang (2008) divided SRC into two sub-tasks in sequence: Each argument should first be determined whether it is a core argument or an adjunct, and then be classified into fine-grained categories. However, delicately designed features are more important and our experiments suggest that by using rich features, a better SRC solver can be directly trained without using hierarchical architecture. There are also some attempts at relaxing the necessity of using full syntactic parses, and semantic chunking methods have been introduced by (Sun et al., 2009; Sun, 2010; Ding and Chang, 2009).

2.1 Our System

We implement a three-stage (i.e. pruning, AI and SRC) SRL system. In the pruning step, our system keeps all constituents (except punctuations) that c-command¹ current predicate in focus as argument candidates. In the AI step, a lot of syntactic features are extracted to distinguish argument and non-argument. In other words, a binary classifier is trained to classify each argument candidate as either an argument or not. Finally, a multi-class classifier is trained to label each argument recognized in the former stage with a specific semantic role label. In both AI and SRC, the main job is to select strong syntactic features.

¹See (Sun et al., 2008) for detailed definition.

3 Features

A majority of features used in our system are a combination of features described in (Xue, 2008; Ding and Chang, 2008) as well as the word formation and coarse frame features introduced in (Sun et al., 2009), the c-command thread features proposed in (Sun et al., 2008). We give a brief description of features used in previous work, but explain new features in details. For more information, readers can refer to relevant papers and our source codes² that are well commented. To conveniently illustrate, we denote a candidate constituent c_k with a fixed context $w_{i-1}[c_k w_i \dots w_h \dots w_j]w_{j+1}$, where w_h is the head word of c_k , and denote predicate in focus with a context $w_{-2}^v w_{-1}^v w^v w_{+1}^v w_{+2}^v$, where w^v is the predicate in focus.

3.1 Baseline Features

The following features are introduced in previous Chinese SRL systems. We use them as baseline.

Word content of w^v , w_h , w_i , w_j and w_{i+w_j} ; *POS tag* of w^v , w_h ; *subcategorization frame*, *verb class* of w^v ; *position*, *phrase type* c_k , *path* from c_k to w^v (from (Xue, 2008; Ding and Chang, 2008))

First character, *last character* and *word length* of w^v , *first character+length*, *last character+word length*, *first character+position*, *last character+position*, *coarse frame*, *frame+ w^v* , *frame+left character*, *frame+verb class*, *frame+ c_k* (from (Sun et al., 2009)).

Head word POS, *head word* of PP phrases, *category* of c_k 's left and right siblings, *CFG rewrite rule* that expands c_k and c_k 's parent (from (Ding and Chang, 2008)).

3.2 New Word Features

We introduce some new features which can be extracted without syntactic structure. We denote them as word features. They include:

Word content of w_{-1}^v , w_{+1}^v , w_{i-1} and w_{j+1} ; *POS tag* of w_{-1}^v , w_{+1}^v , w_{-2}^v , w_{+2}^v , w_{i-1} , w_i , w_j , w_{j+1} , w_{i+2} and w_{j-2} .

Length of c_k : how many words are there in c^k .

Word before "LC": If the POS of w_j is "LC" (localizer), we use w_{j-1} and its POS tag as two new features.

NT: Does c_k contain a word with POS "NT" (temporal noun)?

²Available at <http://code.google.com/p/csrlr/>.

Combination features: w_i 's POS+ w_j 's POS, w_v +Position

3.3 New Syntactic Features

Taking complex syntax trees as inputs, the classifiers should characterize their structural properties. We put forward a number of new features to encode the structural information.

Category of c_k 's parent; head word and POS of head word of parent, left sibling and right sibling of c_k .

Lexicalized Rewrite rules: Conjunction of rewrite rule and head word of its corresponding RHS. These features of candidate (lrw-c) and its parent (lrw-p) are used. For example, this *lrw-c* feature of the NP “事故原因” in Figure 1 is $NP \rightarrow NN + NN$ (原因).

Partial Path: Path from the c_k or w^v to the lowest common ancestor of c_k and w^v . One *path* feature, hence, is divided into *left path* and *right path*.

Clustered Path: We use the manually created clusters (see (Sun and Sui, 2009)) of categories of all nodes in the *path* (cpath) and *right path*.

C-commander thread between c_k and w^v (cct): (proposed by (Sun et al., 2008)). For example, this feature of the NP “警方” in Figure 1 is $NP + ADVP + ADVP + VV$.

Head Trace: The sequential container of the head down upon the phrase (from (Sun and Sui, 2009)). We design two kinds of traces (htr-p, htr-w): one uses POS of the head word; the other uses the head word word itself. E.g., the head word of 事故原因 is “原因” therefore these feature of this NP are $NP \downarrow NN$ and $NP \downarrow$ 原因.

Combination features: *verb class*+ c_k , w_h + w^v , w_h +Position, w_h + w^v +Position, *path*+ w^v , w_h +*right path*, w^v +*left path*, *frame*+ w^v + w_h , and w^v +cct.

4 Experiments and Analysis

4.1 Experimental Setting

To facilitate comparison with previous work, we use CPB 1.0 and CTB 5.0, the same data setting with (Xue, 2008). The data is divided into three parts: files from 081 to 899 are used as training set; files from 041 to 080 as development set; files from 001 to 040, and 900 to 931 as test set. Nearly all previous research on constituency based SRL evaluation use this setting, also including (Ding and Chang, 2008, 2009; Sun

et al., 2009; Sun, 2010). All parsing and SRL experiments use this data setting. To resolve classification problems, we use a linear SVM classifier SVM_{lin}^3 , along with One-Vs-All approach for multi-class classification. To evaluate SRL with automatic parsing, we use a state-of-the-art parser, Bikel parser⁴ (Bikel, 2004). We use gold segmentation and POS as input to the Bikel parser and use its parsing results as input to our SRL system. The overall LP/LR/F performance of Bikel parser is 79.98%/82.95%/81.43.

4.2 Overall Performance

Table 1 summarizes precision, recall and F-measure of AI, SRC and the whole task (AI+SRC) of our system respectively. The fourth line is the best published SRC performance reported in (Ding and Chang, 2008), and the sixth line is the best SRL performance reported in (Xue, 2008). Other lines show the performance of our system. These results indicate a significant improvement over previous systems due to the new features.

Test	P(%)	R(%)	F/A
AI	98.56	97.91	98.24
SRC	--	--	95.04
(Ding and Chang, 2008)	--	--	94.68
AI + SRC	93.80	93.18	93.49
(Xue, 2008)	93.0	91.0	92.0

Table 1: SRL performance on the test data with gold standard parses.

4.3 Two-fold Effect of Parsing in SRL

The effect of parsing in SRL is two-fold. On the one hand, SRL systems should group words as argument candidates, which are also constituents in a given sentence. Full parsing provides boundary information of all constituents. As arguments should c-command the predicate, a full parser can further prune a majority of useless constituents. In other words, parsing can effectively supply SRL with argument candidates. Unfortunately, it is very hard to rightly produce full parses for Chinese text. On the other hand, given a constituent, SRL systems should identify whether it is an argument and further predict detailed semantic types if

³<http://people.cs.uchicago.edu/~vikass/svmlin.html>

⁴<http://www.cis.upenn.edu/~dbikel/software.html>

Task	Parser	Bracket	Feat	P(%)	R(%)	F/A
AI	--	Gold	W	82.44	86.78	84.55
	CTB	Gold	W+S	98.69	98.11	98.40
	Bikel	Bikel	W+S	77.54	71.62	74.46
SRC	--	Gold	W	--	--	93.93
	CTB	Gold	W+S	--	--	95.80
	Bikel	Gold	W+S	--	--	92.62

Table 2: Classification performance on development data. In the *Feat* column, *W* means word features; *W+S* means word and syntactic features.

it is an argument. For the two classification problems, parsing can provide complex syntactic information such as *path* features.

4.3.1 The Effect of Parsing in AI

In AI, full parsing is very important for both grouping words and classification. Table 2 summarizes relative experimental results. Line 2 is the AI performance when gold candidate boundaries and word features are used; Line 3 is the performance with additional syntactic features. Line 4 shows the performance by using automatic parses generated by Bikel parser. We can see that: 1) word features only cannot train good classifiers to identify arguments; 2) it is very easy to recognize arguments with good enough syntactic parses; 3) there is a severe performance decline when automatic parses are used. The third observation is a similar conclusion in English SRL. However this problem in Chinese is much more serious due to the state-of-the-art of Chinese parsing.

Information theoretic criteria are popular criteria in variable selection (Guyon and Elisseeff, 2003). This paper uses empirical mutual information between each variable and the target, $I(X, Y) = \sum_{x \in X, y \in Y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}$, to roughly rank the importance of features. Table 3 shows the ten most useful features in AI. We can see that the most important features all based on full parsing information. Nine of these top 10 useful features are our new features.

Rank	Feature	Rank	Feature
1	w^v_cct	2	$\ddagger w_h + w^v + \text{Position}$
3	htr-w	4	htr-p
5	path	6	$\ddagger w_h + w^v$
7	cpath	8	cct
9	path+ w^v	10	lrw-p

Table 3: Top 10 useful features for AI. \ddagger means *word features*.

4.3.2 The Effect of Parsing in SRC

The second block in Table 2 summarizes the SRC performance with gold argument boundaries. Line 5 is the accuracy when word features are used; Line 6 is the accuracy when additional syntactic features are added; The last row is the accuracy when syntactic features used are extracted from automatic parses (*Bikel+Gold*). We can see that different from AI, word features only can train reasonable good semantic classifiers. The comparison between Line 5 and 7 suggests that with parsing errors, automatic parsed syntactic features cause noise to the semantic role classifiers.

4.4 Why Word Features Are Effective for SRC?

Rank	Feature	Rank	Feature
1	$\ddagger \text{frame} + w_h + w^v$	2	$\ddagger w_h + w^v + \text{position}$
3	$\ddagger w_h + w^v$	4	$w^v + \text{cct}$
5	lrw-p	6	$\ddagger w_i + w_j$
7	lrw-c	8	$\ddagger w_h + \text{Position}$
9	$\ddagger \text{frame} + w^v$	10	htr-p

Table 4: Top 10 useful features for SRC.

Table 4 shows the ten most useful features in SRC. We can see that two of these ten features are word features (denoted by \ddagger). Namely, word features play a more important role in SRC than in AI. Though the other eight features are based on full parsing, four of them (denoted by \ddagger) use the head word which can be well approximated by word features, according to some language specific properties. The head rules described in (Sun and Jurafsky, 2004) are very popular in Chinese parsing research, such as in (Duan et al., 2007; Zhang and Clark, 2008). From these head rules, we can see that head words of most phrases in Chinese are located at the first or the last position. We implement these rules on Chinese Tree Bank and find that 84.12%⁵ nodes realize their heads as either their first or last word. Head position suggests that boundary words are good approximation of head word features. If head words have good approximation word features, then it is not strange that the four features denoted by \ddagger can be effectively represented by word features. Similar with feature effect in AI, most of most useful features in SRC are our new features.

⁵This statistics excludes all empty categories in CTB.

5 Conclusion

This paper proposes an additional set of features to improve Chinese SRL. These new features yield a significant improvement over the best published performance. We further analyze the effect of parsing in Chinese SRL, and linguistically explain some phenomena. We found that (1) full syntactic information plays an essential role only in AI and that (2) due to the head word position distribution, SRC is easy to resolve in Chinese SRL.

Acknowledgments

The author is funded both by German Academic Exchange Service (DAAD) and German Research Center for Artificial Intelligence (DFKI).

The author would like to thank the anonymous reviewers for their helpful comments.

References

- Daniel M. Bikel. 2004. A distributional analysis of a lexicalized statistical parsing model. In Dekang Lin and Dekai Wu, editors, *Proceedings of EMNLP 2004*, pages 182–189. Association for Computational Linguistics, Barcelona, Spain.
- Weiwei Ding and Baobao Chang. 2008. Improving Chinese semantic role classification with hierarchical feature selection strategy. In *Proceedings of the EMNLP 2008*, pages 324–333. Association for Computational Linguistics, Honolulu, Hawaii.
- Weiwei Ding and Baobao Chang. 2009. Fast semantic role labeling for Chinese based on semantic chunking. In *ICCPOL '09: Proceedings of the 22nd International Conference on Computer Processing of Oriental Languages. Language Technology for the Knowledge-based Economy*, pages 79–90. Springer-Verlag, Berlin, Heidelberg.
- Xiangyu Duan, Jun Zhao, and Bo Xu. 2007. Probabilistic models for action-based Chinese dependency parsing. In *ECML '07: Proceedings of the 18th European conference on Machine Learning*, pages 559–566. Springer-Verlag, Berlin, Heidelberg.
- Isabelle Guyon and André Elisseeff. 2003. An introduction to variable and feature selection. *Journal of Machine Learning Research*, 3:1157–1182.
- Honglin Sun and Daniel Jurafsky. 2004. Shallow semantic parsing of Chinese. In Daniel Marcu, Susan Dumais and Salim Roukos, editors, *HLT-NAACL 2004: Main Proceedings*.
- Weiwei Sun. 2010. Semantics-driven shallow parsing for Chinese semantic role labeling. In *Proceedings of the ACL 2010*.
- Weiwei Sun and Zhifang Sui. 2009. Chinese function tag labeling. In *Proceedings of the 23rd Pacific Asia Conference on Language, Information and Computation*. Hong Kong.
- Weiwei Sun, Zhifang Sui, and Haifeng Wang. 2008. Prediction of maximal projection for semantic role labeling. In *Proceedings of the 22nd International Conference on Computational Linguistics*.
- Weiwei Sun, Zhifang Sui, Meng Wang, and Xin Wang. 2009. Chinese semantic role labeling with shallow parsing. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing*, pages 1475–1483. Association for Computational Linguistics, Singapore.
- Nianwen Xue. 2008. Labeling Chinese predicates with semantic roles. *Comput. Linguist.*, 34(2):225–255.
- Yue Zhang and Stephen Clark. 2008. A tale of two parsers: Investigating and combining graph-based and transition-based dependency parsing. In *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*, pages 562–571. Association for Computational Linguistics, Honolulu, Hawaii.

Balancing User Effort and Translation Error in Interactive Machine Translation Via Confidence Measures

Jesús González-Rubio

Inst. Tec. de Informática
Univ. Politéc. de Valencia
46021 Valencia, Spain
jegonzalez@iti.upv.es

Daniel Ortiz-Martínez

Dpto. de Sist Inf. y Comp.
Univ. Politéc. de Valencia
46021 Valencia, Spain
dortiz@dsic.upv.es

Francisco Casacuberta

Dpto. de Sist Inf. y Comp.
Univ. Politéc. de Valencia
46021 Valencia, Spain
fcn@dsic.upv.es

Abstract

This work deals with the application of confidence measures within an interactive-predictive machine translation system in order to reduce human effort. If a small loss in translation quality can be tolerated for the sake of efficiency, user effort can be saved by interactively translating only those initial translations which the confidence measure classifies as incorrect. We apply confidence estimation as a way to achieve a balance between user effort savings and final translation error. Empirical results show that our proposal allows to obtain almost perfect translations while significantly reducing user effort.

1 Introduction

In *Statistical Machine Translation* (SMT), the translation is modelled as a decision process. For a given source string $f_1^J = f_1 \dots f_j \dots f_J$, we seek for the target string $e_1^I = e_1 \dots e_i \dots e_I$ which maximises posterior probability:

$$\hat{e}_1^I = \operatorname{argmax}_{I, e_1^I} Pr(e_1^I | f_1^J). \quad (1)$$

Within the *Interactive-predictive Machine Translation* (IMT) framework, a state-of-the-art SMT system is employed in the following way: For a given source sentence, the SMT system fully automatically generates an initial translation. A human translator checks this translation from left to right, correcting the first error. The SMT system then proposes a new extension, taking the correct prefix $e_1^i = e_1 \dots e_i$ into account. These steps are repeated until the whole input sentence has been correctly translated. In the resulting decision rule, we maximise over all possible extensions e_{i+1}^I of e_1^i :

$$\hat{e}_{i+1}^I = \operatorname{argmax}_{I, e_{i+1}^I} Pr(e_{i+1}^I | e_1^i, f_1^J). \quad (2)$$

An implementation of the IMT framework was performed in the TransType project (Foster et al., 1997; Langlais et al., 2002) and further improved within the TransType2 project (Esteban et al., 2004; Barrachina et al., 2009).

IMT aims at reducing the effort and increasing the productivity of translators, while preserving high-quality translation. In this work, we integrate *Confidence Measures* (CMs) within the IMT framework to further reduce the user effort. As will be shown, our proposal allows to balance the ratio between user effort and final translation error.

1.1 Confidence Measures

Confidence estimation have been extensively studied for speech recognition. Only recently have researchers started to investigate CMs for MT (Gandrabur and Foster, 2003; Blatz et al., 2004; Ueffing and Ney, 2007).

Different TransType-style MT systems use confidence information to improve translation prediction accuracy (Gandrabur and Foster, 2003; Ueffing and Ney, 2005). In this work, we propose a focus shift in which CMs are used to modify the interaction between the user and the system instead of modify the IMT translation predictions.

To compute CMs we have to select suitable confidence features and define a binary classifier. Typically, the classification is carried out depending on whether the confidence value exceeds a given threshold or not.

2 IMT with Sentence CMs

In the conventional IMT scenario a human translator and a SMT system collaborate in order to obtain the translation the user has in mind. Once the user has interactively translated the source sentences, the output translations are error-free. We propose an alternative scenario where not all the source sentences are interactively translated by the user. Specifically, only those source sentences

whose initial fully automatic translation are incorrect, according to some quality criterion, are interactively translated. We propose to use CMs as the quality criterion to classify those initial translations.

Our approach implies a modification of the user-machine interaction protocol. For a given source sentence, the SMT system generates an initial translation. Then, if the CM classifies this translation as correct, we output it as our final translation. On the contrary, if the initial translation is classified as incorrect, we perform a conventional IMT procedure, validating correct prefixes and generating new suffixes, until the sentence that the user has in mind is reached.

In our scenario, we allow the final translations to be different from the ones the user has in mind. This implies that the output may contain errors. If a small loss in translation can be tolerated for the sake of efficiency, user effort can be saved by interactively translating only those sentences that the CMs classify as incorrect.

It is worth of notice that our proposal can be seen as a generalisation of the conventional IMT approach. Varying the value of the CM classification threshold, we can range from a fully automatic SMT system where all sentences are classified as correct to a conventional IMT system where all sentences are classified as incorrect.

2.1 Selecting a CM for IMT

We compute sentence CMs by combining the scores given by a word CM based on the IBM model 1 (Brown et al., 1993), similar to the one described in (Blatz et al., 2004). We modified this word CM by replacing the *average* by the *maximal* lexicon probability, because the average is dominated by this maximum (Ueffing and Ney, 2005). We choose this word CM because it can be calculated very fast during search, which is crucial given the time constraints of the IMT systems. Moreover, its performance is similar to that of other word CMs as results presented in (Blatz et al., 2003; Blatz et al., 2004) show. The word confidence value of word e_i , $c_w(e_i)$, is given by

$$c_w(e_i) = \max_{0 \leq j \leq J} p(e_i | f_j), \quad (3)$$

where $p(e_i | f_j)$ is the IBM model 1 lexicon probability, and f_0 is the empty source word.

From this word CM, we compute two sentence CMs which differ in the way the word confidence

		Spanish	English
Train	Sentences	214.5K	
	Running words	5.8M	5.2M
	Vocabulary	97.4K	83.7K
Dev.	Sentences	400	
	Running words	11.5K	10.1K
	Perplexity (trigrams)	46.1	59.4
Test	Sentences	800	
	Running words	22.6K	19.9K
	Perplexity (trigrams)	45.2	60.8

Table 1: Statistics of the Spanish–English EU corpora. K and M denote thousands and millions of elements respectively.

scores $c_w(e_i)$ are combined:

MEAN CM ($c_M(e_1^I)$) is computed as the geometric mean of the confidence scores of the words in the sentence:

$$c_M(e_1^I) = \sqrt[I]{\prod_{i=1}^I c_w(e_i)}. \quad (4)$$

RATIO CM ($c_R(e_1^I)$) is computed as the percentage of words classified as correct in the sentence. A word is classified as correct if its confidence exceeds a word classification threshold τ_w .

$$c_R(e_1^I) = \frac{|\{e_i / c_w(e_i) > \tau_w\}|}{I} \quad (5)$$

After computing the confidence value, each sentence is classified as either correct or incorrect, depending on whether its confidence value exceeds or not a sentence classification threshold τ_s . If $\tau_s = 0.0$ then all the sentences will be classified as correct whereas if $\tau_s = 1.0$ all the sentences will be classified as incorrect.

3 Experimentation

The aim of the experimentation was to study the possibly trade-off between saved user effort and translation error obtained when using sentence CMs within the IMT framework.

3.1 System evaluation

In this paper, we report our results as measured by *Word Stroke Ratio* (WSR) (Barrachina et al., 2009). WSR is used in the context of IMT to measure the effort required by the user to generate her

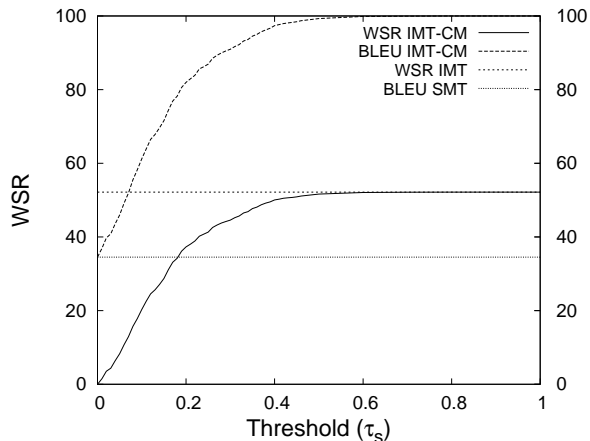


Figure 1: BLEU translation scores versus WSR for different values of the sentence classification threshold using the MEAN CM.

translations. WSR is computed as the ratio between the number of word-strokes a user would need to achieve the translation she has in mind and the total number of words in the sentence. In this context, a word-stroke is interpreted as a single action, in which the user types a complete word, and is assumed to have constant cost.

Additionally, and because our proposal allows differences between its output and the reference translation, we will also present translation quality results in terms of *BiLingual Evaluation Understudy* (BLEU) (Papineni et al., 2002). BLEU computes a geometric mean of the precision of n -grams multiplied by a factor to penalise short sentences.

3.2 Experimental Setup

Our experiments were carried out on the EU corpora (Barrachina et al., 2009). The EU corpora were extracted from the Bulletin of the European Union. The EU corpora is composed of sentences given in three different language pairs. Here, we will focus on the Spanish–English part of the EU corpora. The corpus is divided into training, development and test sets. The main figures of the corpus can be seen in Table 1.

As a first step, we built a SMT system to translate from Spanish into English. This was done by means of the Thot toolkit (Ortiz et al., 2005), which is a complete system for building phrase-based SMT models. This toolkit involves the estimation, from the training set, of different statistical models, which are in turn combined in a log-linear fashion by adjusting a weight for each of them by means of the MERT (Och, 2003) procedure,

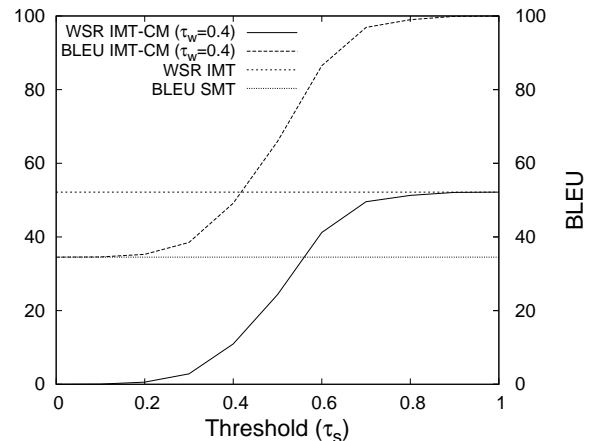


Figure 2: BLEU translation scores versus WSR for different values of the sentence classification threshold using the RATIO CM with $\tau_w = 0.4$.

optimising the BLEU score on the development set.

The IMT system which we have implemented relies on the use of word graphs (Ueffing et al., 2002) to efficiently compute the suffix for a given prefix. A word graph has to be generated for each sentence to be interactively translated. For this purpose, we used a multi-stack phrase-based decoder which will be distributed in the near future together with the Thot toolkit. We discarded to use the state-of-the-art Moses toolkit (Koehn et al., 2007) because preliminary experiments performed with it revealed that the decoder by Ortiz-Martínez et al. (2005) performs better in terms of WSR when used to generate word graphs for their use in IMT (Sanchis-Trilles et al., 2008). Moreover, the performance difference in regular SMT is negligible. The decoder was set to only consider monotonic translation, since in real IMT scenarios considering non-monotonic translation leads to excessive response time for the user.

Finally, the obtained word graphs were used within the IMT procedure to produce the reference translations in the test set, measuring WSR and BLEU.

3.3 Results

We carried out a series of experiments ranging the value of the sentence classification threshold τ_s , between 0.0 (equivalent to a fully automatic SMT system) and 1.0 (equivalent to a conventional IMT system), for both the MEAN and RATIO CMs. For each threshold value, we calculated the effort of the user in terms of WSR, and the translation quality of the final output as measured by BLEU.

src-1	DECLARACIÓN (No 17) relativa al derecho de acceso a la información
ref-1	DECLARATION (No 17) on the right of access to information
tra-1	DECLARATION (No 17) on the right of access to information
src-2	Conclusiones del Consejo sobre el comercio electrónico y los impuestos indirectos.
ref-2	Council conclusions on electronic commerce and indirect taxation.
tra-2	Council conclusions on e-commerce and indirect taxation.
src-3	participación de los países candidatos en los programas comunitarios.
ref-3	participation of the applicant countries in Community programmes.
tra-3	countries' involvement in Community programmes.

Example 1: Examples of initial fully automatically generated sentences classified as correct by the CMs.

Figure 1 shows WSR (WSR IMT-CM) and BLEU (BLEU IMT-CM) scores obtained varying τ_s for the MEAN CM. Additionally, we also show the BLEU score (BLEU SMT) obtained by a fully automatic SMT system as translation quality baseline, and the WSR score (WSR IMT) obtained by a conventional IMT system as user effort baseline. This figure shows a continuous transition between the fully automatic SMT system and the conventional IMT system. This transition occurs when ranging τ_s between 0.0 and 0.6. This is an undesired effect, since for almost a half of the possible values for τ_s there is no change in the behaviour of our proposed IMT system.

The RATIO CM confidence values depend on a word classification threshold τ_w . We have carried out experimentation ranging τ_w between 0.0 and 1.0 and found that this value can be used to solve the above mentioned undesired effect for the MEAN CM. Specifically, varying the value of τ_w we can stretch the interval in which the transition between the fully automatic SMT system and the conventional IMT system is produced, allowing us to obtain smoother transitions. Figure 2 shows WSR and BLEU scores for different values of the sentence classification threshold τ_s using $\tau_w = 0.4$. We show results only for this value of τ_w due to paper space limitations and because $\tau_w = 0.4$ produced the smoothest transition. According to Figure 2, using a sentence classification threshold value of 0.6 we obtain a WSR reduction of 20% relative and an almost perfect translation quality of 87 BLEU points.

It is worth of notice that the final translations are compared with only one reference, therefore, the reported translation quality scores are clearly pessimistic. Better results are expected using a multi-reference corpus. Example 1 shows the source sentence (src), the reference translation

(ref) and the final translation (tra) for three of the initial fully automatically generated translations that were classified as correct by our CMs, and thus, were not interactively translated by the user. The first translation (tra-1) is identical to the corresponding reference translation (ref-1). The second translation (tra-2) corresponds to a correct translation of the source sentence (src-2) that is different from the corresponding reference (ref-2). Finally, the third translation (tra-3) is an example of a slightly incorrect translation.

4 Concluding Remarks

In this paper, we have presented a novel proposal that introduces sentence CMs into an IMT system to reduce user effort. Our proposal entails a modification of the user-machine interaction protocol that allows to achieve a balance between the user effort and the final translation error.

We have carried out experimentation using two different sentence CMs. Varying the value of the sentence classification threshold, we can range from a fully automatic SMT system to a conventional IMT system. Empirical results show that our proposal allows to obtain almost perfect translations while significantly reducing user effort.

Future research aims at the investigation of improved CMs to be integrated in our IMT system.

Acknowledgments

Work supported by the EC (FEDER/FSE) and the Spanish MEC/MICINN under the MIPRCV “Consolider Ingenio 2010” program (CSD2007-00018), the iTransDoc (TIN2006-15694-CO2-01) and iTrans2 (TIN2009-14511) projects and the FPU scholarship AP2006-00691. Also supported by the Spanish MITYC under the erudito.com (TSI-020110-2009-439) project and by the Generalitat Valenciana under grant Prometeo/2009/014.

References

- S. Barrachina, O. Bender, F. Casacuberta, J. Civera, E. Cubel, S. Khadivi, A. Lagarda, H. Ney, J. Tomás, and E. Vidal. 2009. Statistical approaches to computer-assisted translation. *Computational Linguistics*, 35(1):3–28.
- J. Blatz, E. Fitzgerald, G. Foster, S. Gandrabur, C. Goutte, A. Kulesza, A. Sanchis, and N. Ueffing. 2003. Confidence estimation for machine translation.
- J. Blatz, E. Fitzgerald, G. Foster, S. Gandrabur, C. Goutte, A. Kulesza, A. Sanchis, and N. Ueffing. 2004. Confidence estimation for machine translation. In *Proc. COLING*, page 315.
- P. F. Brown, S. A. Della Pietra, V. J. Della Pietra, and R. L. Mercer. 1993. The Mathematics of Statistical Machine Translation: Parameter Estimation. *Computational Linguistics*, 19(2):263–311.
- J. Esteban, J. Lorenzo, A. Valderrábanos, and G. Lapalme. 2004. Transtype2: an innovative computer-assisted translation system. In *Proc. ACL*, page 1.
- G. Foster, P. Isabelle, and P. Plamondon. 1997. Target-text mediated interactive machine translation. *Machine Translation*, 12:12–175.
- S. Gandrabur and G. Foster. 2003. Confidence estimation for text prediction. In *Proc. CoNLL*, pages 315–321.
- P. Koehn, H. Hoang, A. Birch, C. Callison-Burch, M. Federico, N. Bertoldi, B. Cowan, W. Shen, C. Moran, R. Zens, C. Dyer, O. Bojar, A. Constantin, and E. Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In *Proc. ACL*, pages 177–180.
- P. Langlais, G. Lapalme, and M. Loranger. 2002. Transtype: Development-evaluation cycles to boost translator’s productivity. *Machine Translation*, 15(4):77–98.
- F. J. Och. 2003. Minimum error rate training in statistical machine translation. In *Proc. ACL*, pages 160–167.
- D. Ortiz, I. García-Varea, and F. Casacuberta. 2005. Thot: a toolkit to train phrase-based statistical translation models. In *Proc. MT Summit*, pages 141–148.
- K. Papineni, S. Roukos, T. Ward, and W. Zhu. 2002. BLEU: a method for automatic evaluation of MT. In *Proc. ACL*, pages 311–318.
- G. Sanchis-Trilles, D. Ortiz-Martínez, J. Civera, F. Casacuberta, E. Vidal, and H. Hoang. 2008. Improving interactive machine translation via mouse actions. In *Proc. EMNLP*, pages 25–27.
- N. Ueffing and H. Ney. 2005. Application of word-level confidence measures in interactive statistical machine translation. In *Proc. EAMT*, pages 262–270.
- N. Ueffing and H. Ney. 2007. Word-level confidence estimation for machine translation. *Comput. Linguist.*, 33(1):9–40.
- N. Ueffing, F.J. Och, and H. Ney. 2002. Generation of word graphs in statistical machine translation. In *Proc. EMNLP*, pages 156–163.

Improving Arabic-to-English Statistical Machine Translation by Reordering Post-verbal Subjects for Alignment

Marine Carpuat Yuval Marton Nizar Habash

Columbia University

Center for Computational Learning Systems

475 Riverside Drive, New York, NY 10115

{marine,ymarton,habash}@ccls.columbia.edu

Abstract

We study the challenges raised by Arabic verb and subject detection and reordering in Statistical Machine Translation (SMT). We show that post-verbal subject (VS) constructions are hard to translate because they have highly ambiguous reordering patterns when translated to English. In addition, implementing reordering is difficult because the boundaries of VS constructions are hard to detect accurately, even with a state-of-the-art Arabic dependency parser. We therefore propose to reorder VS constructions into SV order for SMT word alignment only. This strategy significantly improves BLEU and TER scores, even on a strong large-scale baseline and despite noisy parses.

1 Introduction

Modern Standard Arabic (MSA) is a morpho-syntactically complex language, with different phenomena from English, a fact that raises many interesting issues for natural language processing and Arabic-to-English statistical machine translation (SMT). While comprehensive Arabic preprocessing schemes have been widely adopted for handling Arabic morphology in SMT (e.g., Sadat and Habash (2006), Zollmann et al. (2006), Lee (2004)), syntactic issues have not received as much attention by comparison (Green et al. (2009), Crego and Habash (2008), Habash (2007)). Arabic verbal constructions are particularly challenging since subjects can occur in pre-verbal (SV), post-verbal (VS) or pro-dropped (“null subject”) constructions. As a result, training data for learning verbal construction translations is split between the different constructions and their patterns; and complex reordering schemas are needed in order to translate them into primarily

pre-verbal subject languages (SVO) such as English.

These issues are particularly problematic in phrase-based SMT (Koehn et al., 2003). Standard phrase-based SMT systems memorize phrasal translation of verb and subject constructions as observed in the training bitext. They do not capture any generalizations between occurrences in VS and SV orders, even for the same verbs. In addition, their distance-based reordering models are not well suited to handling complex reordering operations which can include long distance dependencies, and may vary by context. Despite these limitations, phrase-based SMT systems have achieved competitive results in Arabic-to-English benchmark evaluations.¹ However, error analysis shows that verbs are still often dropped or incorrectly translated, and subjects are split or garbled in translation. This suggests that better syntactic modeling should further improve SMT.

We attempt to get a better understanding of translation patterns for Arabic verb constructions, particularly VS constructions, by studying their occurrence and reordering patterns in a hand-aligned Arabic-English parallel treebank. Our analysis shows that VS reordering rules are not straightforward and that SMT should therefore benefit from direct modeling of Arabic verb subject translation. In order to detect VS constructions, we use our state-of-the-art Arabic dependency parser, which is essentially the CATIBEX baseline in our subsequent parsing work in Marton et al. (2010), and is further described there. We show that VS subjects and their exact boundaries are hard to identify accurately. Given the noise in VS detection, existing strategies for source-side reordering (e.g., Xia and McCord (2004), Collins et al. (2005), Wang et al. (2007)) or using de-

¹<http://www.itl.nist.gov/iad/mig/tests/mt/2009/ResultsRelease/currentArabic.html>

Table 1: How are Arabic SV and VS translated in the manually word-aligned Arabic-English parallel treebank? We check whether V and S are translated in a “monotone” or “inverted” order for all VS and SV constructions. “Overlap” represents instances where translations of the Arabic verb and subject have some English words in common, and are not monotone nor inverted.

	gold reordering	all verbs	%
SV	monotone	2588	98.2
SV	inverted	15	0.5
SV	overlap	35	1.3
SV	total	2638	100
VS	monotone	1700	27.3
VS	inverted	4033	64.7
VS	overlap	502	8.0
VS	total	6235	100

pendency parses as cohesion constraints in decoding (e.g., Cherry (2008); Bach et al. (2009)) are not effective at this stage. While these approaches have been successful for language pairs such as German-English for which syntactic parsers are more developed and relevant reordering patterns might be less ambiguous, their impact potential on Arabic-English translation is still unclear.

In this work, we focus on VS constructions only, and propose a new strategy in order to benefit from their noisy detection: for the word alignment stage only, we reorder phrases detected as VS constructions into an SV order. Then, for phrase extraction, weight optimization and decoding, we use the original (non-reordered) text. This approach significantly improves both BLEU and TER on top of strong medium and large-scale phrase-based SMT baselines.

2 VS reordering in gold Arabic-English translation

We use the manually word-aligned parallel Arabic-English Treebank (LDC2009E82) to study how Arabic VS constructions are translated into English by humans. Given the gold Arabic syntactic parses and the manual Arabic-English word alignments, we can determine the gold reorderings for SV and VS constructions. We extract VS representations from the gold constituent parses by deterministic conversion to a simplified dependency structure, CATiB (Habash and Roth, 2009)

(see Section 3). We then check whether the English translations of the Arabic verb and the Arabic subject occur in the same order as in Arabic (monotone) or not (inverted). Table 1 summarizes the reordering patterns for each category. As expected, 98% of Arabic SV are translated in a monotone order in English. For VS constructions, the picture is surprisingly more complex. The monotone VS translations are mostly explained by changes to passive voice or to non-verbal constructions (such as nominalization) in the English translation.

In addition, Table 1 shows that verb subjects occur more frequently in VS order (70%) than in SV order (30%). These numbers do not include pro-dropped (“null subject”) constructions.

3 Arabic VS construction detection

Even if the SMT system had perfect knowledge of VS reordering, it has to accurately detect VS constructions and their spans in order to apply the reordering correctly. For that purpose, we use our state-of-the-art parsing model, which is essentially the CATIBEX baseline model in Marton et al. (2010), and whose details we summarize next. We train a syntactic dependency parser, MaltParser v1.3 with the Nivre “eager” algorithm (Nivre, 2003; Nivre et al., 2006; Nivre, 2008) on the training portion of the Penn Arabic Treebank part 3 v3.1, hereafter PATB3 (Maamouri et al., 2008; Maamouri et al., 2009). The training / development split is the same as in Zitouni et al. (2006). We convert the PATB3 representation into the succinct CATiB format, with 8 dependency relations and 6 POS tags, which we then extend to a set of 44 tags using regular expressions of the basic POS and the normalized surface word form, similarly to Marton et al. (2010), following Habash and Roth (2009). We normalize Alif Maq-sura to Ya, and Hamzated Alifs to bare Alif, as is commonly done in Arabic SMT.

For analysis purposes, we evaluate our subject and verb detection on the development part of PATB3 using gold POS tags. There are various ways to go about it. We argue that combined detection statistics of constructions of verbs and their subjects (VATS), for which we achieve an F-score of 74%, are more telling for the task at hand.²

²We divert from the CATiB representation in that a non-matrix subject of a pseudo verb (*An and her sisters*) is treated as a subject of the verb that is under the same pseudo verb. This treatment of said subjects is comparable to the PATB’s.

These scores take into account the spans of both the subject and the specific verb it belongs to, and potentially reorder with. We also provide statistics of VS detection separately (F-score 63%), since we only handle VS here. This low score can be explained by the difficulty in detecting the post-verbal subject’s end boundary, and the correct verb the subject belongs to. The SV construction scores are higher, presumably since the pre-verbal subject’s end is bounded by the verb it belongs to. See Table 2.

Although not directly comparable, our VS scores are similar to those of Green et al. (2009). Their VS detection technique with conditional random fields (CRF) is different from ours in bypassing full syntactic parsing, and in only detecting maximal (non-nested) subjects of verb-initial clauses. Additionally, they use a different training / test split of the PATB data (parts 1, 2 and 3). They report 65.9% precision and 61.3% F-score. Note that a closer score comparison should take into account their reported verb detection accuracy of 98.1%.

Table 2: Precision, Recall and F-scores for constructions of Arabic verbs and their subjects, evaluated on our development part of PATB3.

construction	P	R	F
VATS (verbs & their subj.)	73.84	74.37	74.11
VS	66.62	59.41	62.81
SV	86.75	61.07	71.68
VNS (verbs w/ null subj.)	76.32	92.04	83.45
verbal subj. exc. null subj.	72.46	60.18	65.75
verbal subj. inc. null subj.	73.97	74.50	74.23
verbs with non-null subj.	91.94	76.17	83.31
SV or VS	72.19	59.95	65.50

4 Reordering Arabic VS for SMT word alignment

Based on these analyses, we propose a new method to help phrase-based SMT systems deal with Arabic-English word order differences due to VS constructions. As in related work on syntactic reordering by preprocessing, our method attempts to make Arabic and English word order closer to each other by reordering Arabic VS constructions into SV. However, unlike in previous work, the reordered Arabic sentences are used only for word alignment. Phrase translation extraction and de-

coding are performed on the original Arabic word order. Preliminary experiments on an earlier version of the large-scale SMT system described in Section 6 showed that forcing reordering of all VS constructions at training and test time does not have a consistent impact on translation quality: for instance, on the NIST MT08-NW test set, TER slightly improved from 44.34 to 44.03, while BLEU score decreased from 49.21 to 49.09.

Limiting reordering to alignment allows the system to be more robust and recover from incorrect changes introduced either by incorrect VS detection, or by incorrect reordering of a correctly detected VS. Given a parallel sentence (a, e) , we proceed as follows:

1. automatically tag VS constructions in a
2. generate new sentence $a' = reorder(a)$ by reordering Arabic VS into SV
3. get word alignment wa' on new sentence pair (a', e)
4. using mapping from a to a' , get corresponding word alignment $wa = unreorder(wa')$ for the original sentence pair (a, e)

5 Experiment set-up

We use the open-source Moses toolkit (Koehn et al., 2007) to build two phrase-based SMT systems trained on two different data conditions:

- **medium-scale** the bitext consists of 12M words on the Arabic side (LDC2007E103). The language model is trained on the English side of the large bitext.
- **large-scale** the bitext consists of several newswire LDC corpora, and has 64M words on the Arabic side. The language model is trained on the English side of the bitext augmented with Gigaword data.

Except from this difference in training data, the two systems are identical. They use a standard phrase-based architecture. The parallel corpus is word-aligned using the GIZA++ (Och and Ney, 2003), which sequentially learns word alignments for the IBM1, HMM, IBM3 and IBM4 models. The resulting alignments in both translation directions are intersected and augmented using the grow-diag-final-and heuristic (Koehn et al., 2007). Phrase translations of up to 10 words are extracted in the Moses phrase-table. We apply statistical significance tests to prune unreliable phrase-pairs

and score remaining phrase-table entries (Chen et al., 2009). We use a 5-gram language model with modified Kneser-Ney smoothing. Feature weights are tuned to maximize BLEU on the NIST MT06 test set.

For all systems, the English data is tokenized using simple punctuation-based rules. The Arabic side is segmented according to the Arabic Treebank (PATB3) tokenization scheme (Maamouri et al., 2009) using the MADA+TOKAN morphological analyzer and tokenizer (Habash and Rambow, 2005). MADA-produced Arabic lemmas are used for word alignment.

6 Results

We evaluate translation quality using both BLEU (Papineni et al., 2002) and TER (Snover et al., 2006) scores on three standard evaluation test sets from the NIST evaluations, which yield more than 4400 test sentences with 4 reference translations. On this large data set, our VS reordering method remarkably yields statistically significant improvements in BLEU and TER on the medium and large SMT systems at the 99% confidence level (Table 3).

Results per test set are reported in Table 4. TER scores are improved in all 10 test configurations, and BLEU scores are improved in 8 out of the 10 configurations. Results on the MT08 test set show that improvements are obtained both on newswire and on web text as measured by TER (but not BLEU score on the web section.) It is worth noting that consistent improvements are obtained even on the large-scale system, and that both baselines are full-fledged systems, which include lexicalized reordering and large 5-gram language models.

Analysis shows that our VS reordering technique improves word alignment coverage (yielding 48k and 330k additional links on the medium and large scale systems respectively). This results in larger phrase-tables which improve translation quality.

7 Related work

To the best of our knowledge, the only other approach to detecting and using Arabic verb-subject constructions for SMT is that of Green et al. (2009) (see Section 3), which failed to improve Arabic-English SMT. In contrast with our reordering approach, they integrate subject span information as a log-linear model feature which encour-

Table 3: Evaluation on all test sets: on the total of 4432 test sentences, improvements are statistically significant at the 99% level using bootstrap resampling (Koehn, 2004)

system	BLEU r4n4 (%)	TER (%)
medium baseline	44.35	48.34
+ VS reordering	44.65 (+0.30)	47.78 (-0.56)
large baseline	51.45	42.45
+ VS reordering	51.70 (+0.25)	42.21 (-0.24)

ages a phrase-based SMT decoder to use phrasal translations that do not break subject boundaries.

Syntactically motivated reordering for phrase-based SMT has been more successful on language pairs other than Arabic-English, perhaps due to more accurate parsers and less ambiguous reordering patterns than for Arabic VS. For instance, Collins et al. (2005) apply six manually defined transformations to German parse trees which improve German-English translation by 0.4 BLEU on the Europarl task. Xia and McCord (2004) learn reordering rules for French to English translations, which arguably presents less syntactic distortion than Arabic-English. Zhang et al. (2007) limit reordering to decoding for Chinese-English SMT using a lattice representation. Cherry (2008) uses dependency parses as cohesion constraints in decoding for French-English SMT.

For Arabic-English phrase-based SMT, the impact of syntactic reordering as preprocessing is less clear. Habash (2007) proposes to learn syntactic reordering rules targeting Arabic-English word order differences and integrates them as deterministic preprocessing. He reports improvements in BLEU compared to phrase-based SMT limited to monotonic decoding, but these improvements do not hold with distortion. Instead of applying reordering rules deterministically, Crego and Habash (2008) use a lattice input to represent alternate word orders which improves a ngram-based SMT system. But they do not model VS constructions explicitly.

Most previous syntax-aware word alignment models were specifically designed for syntax-based SMT systems. These models are often bootstrapped from existing word alignments, and could therefore benefit from our VS reordering approach. For instance, Fossum et al. (2008) report improvements ranging from 0.1 to 0.5 BLEU on Arabic translation by learning to delete alignment

Table 4: VS reordering improves BLEU and TER scores in almost all test conditions on 5 test sets, 2 metrics, and 2 MT systems

BLEU r4n4 (%)					
test set	MT03	MT04	MT05	MT08nw	MT08wb
medium baseline	45.95	44.94	48.05	44.86	32.05
+ VS reordering	46.33 (+0.38)	45.03 (+0.09)	48.69 (+0.64)	45.06 (+0.20)	31.96 (-0.09)
large baseline	52.3	52.45	54.66	52.60	39.22
+ VS reordering	52.63 (+0.33)	52.34 (-0.11)	55.29 (+0.63)	52.85 (+0.25)	39.87 (+0.65)
TER (%)					
test set	MT03	MT04	MT05	MT08nw	MT08wb
medium baseline	48.77	46.45	45.00	47.74	58.02
+ VS reordering	48.31 (-0.46)	46.10 (-0.35)	44.29 (-0.71)	47.11 (-0.63)	57.30 (-0.72)
large baseline	43.33	40.42	39.15	41.81	52.05
+ VS reordering	42.95 (-0.38)	40.40 (-0.02)	38.75 (-0.40)	41.51 (-0.30)	51.86 (-0.19)

links if they degrade their syntax-based translation system. Departing from commonly-used alignment models, Hermjakob (2009) aligns Arabic and English content words using pointwise mutual information, and in this process indirectly uses English sentences reordered into VS order to collect cooccurrence counts. The approach outperforms GIZA++ on a small-scale translation task, but the impact of reordering alone is not evaluated.

8 Conclusion and future work

We presented a novel method for improving overall SMT quality using a noisy syntactic parser: we use these parses to reorder VS constructions into SV for word alignment only. This approach increases word alignment coverage and significantly improves BLEU and TER scores on two strong SMT baselines.

In subsequent work, we show that matrix (main-clause) VS constructions are reordered much more frequently than non-matrix VS, and that limiting reordering to matrix VS constructions for word alignment further improves translation quality (Carpuat et al., 2010). In the future, we plan to improve robustness to parsing errors by using not just one, but multiple subject boundary hypotheses. We will also investigate the integration of VS reordering in SMT decoding.

Acknowledgements

The authors would like to thank Mona Diab, Owen Rambow, Ryan Roth, Kristen Parton and Joakim Nivre for helpful discussions and assistance. This material is based upon work supported by the Defense Advanced Research Projects Agency (DARPA) under GALE Contract No HR0011-08-C-

0110. Any opinions, findings and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of DARPA.

References

- Nguyen Bach, Stephan Vogel, and Colin Cherry. 2009. Cohesive constraints in a beam search phrase-based decoder. In *Proceedings of the 10th Meeting of the North American Chapter of the Association for Computational Linguistics, Companion Volume: Short Papers*, pages 1–4.
- Marine Carpuat, Yuval Marton, and Nizar Habash. 2010. Reordering matrix post-verbal subjects for arabic-to-english smt. In *Proceedings of the Conference Traitement Automatique des Langues Naturelles (TALN)*.
- Boxing Chen, George Foster, and Roland Kuhn. 2009. Phrase translation model enhanced with association based features. In *Proceedings of MT-Summit XII*, Ottawa, Ontario, September.
- Colin Cherry. 2008. Cohesive phrase-based decoding for statistical machine translation. In *Proceedings of the 46th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 72–80, Columbus, Ohio, June.
- Michael Collins, Philipp Koehn, and Ivona Kucerova. 2005. Clause restructuring for statistical machine translation. In *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 531–540, Ann Arbor, MI, June.
- Josep M. Crego and Nizar Habash. 2008. Using shallow syntax information to improve word alignment and reordering for SMT. In *Proceedings of the Third Workshop on Statistical Machine Translation*, pages 53–61, June.
- Victoria Fossum, Kevin Knight, and Steven Abney. 2008. Using syntax to improve word alignment precision for syntax-based machine translation. In *Proceedings of the Third Workshop on Statistical Machine Translation*, pages 44–52.
- Spence Green, Conal Sathi, and Christopher D. Manning. 2009. NP subject detection in verb-initial Arabic clauses.

- In *Proceedings of the Third Workshop on Computational Approaches to Arabic Script-based Languages (CAASL3)*.
- Nizar Habash and Owen Rambow. 2005. Arabic Tokenization, Part-of-Speech Tagging and Morphological Disambiguation in One Fell Swoop. In *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL'05)*, pages 573–580, Ann Arbor, Michigan, June.
- Nizar Habash and Ryan Roth. 2009. CATiB: The Columbia Arabic treebank. In *Proceedings of the ACL-IJCNLP 2009 Conference Short Papers*, pages 221–224, Suntec, Singapore, August. Association for Computational Linguistics.
- Nizar Habash. 2007. Syntactic preprocessing for statistical machine translation. In *Proceedings of the Machine Translation Summit (MT-Summit)*, Copenhagen.
- Ulf Hermjakob. 2009. Improved word alignment with statistics and linguistic heuristics. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing*, pages 229–237, Singapore, August.
- Philipp Koehn, Franz Och, and Daniel Marcu. 2003. Statistical phrase-based translation. In *Proceedings of HLT/NAACL-2003*, Edmonton, Canada, May.
- Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondrej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In *Annual Meeting of the Association for Computational Linguistics (ACL), demonstration session*, Prague, Czech Republic, June.
- Philipp Koehn. 2004. Statistical significance tests for machine translation evaluation. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing (EMNLP-2004)*, Barcelona, Spain, July.
- Young-Suk Lee. 2004. Morphological analysis for statistical machine translation. In *Proceedings of the Human Language Technology Conference of the NAACL*, pages 57–60, Boston, MA.
- Mohamed Maamouri, Ann Bies, and Seth Kulick. 2008. Enhancing the arabic treebank: a collaborative effort toward new annotation guidelines. In *Proceedings of the Sixth International Language Resources and Evaluation (LREC'08)*, Marrakech, Morocco.
- Mohamed Maamouri, Ann Bies, Seth Kulick, Fatma Gaddeche, Wigdan Mekki, Sondos Krouna, and Basma Bouziri. 2009. The penn arabic treebank part 3 version 3.1. Linguistic Data Consortium LDC2008E22.
- Yuval Marton, Nizar Habash, and Owen Rambow. 2010. Improving arabic dependency parsing with lexical and inflectional morphological features. In *Proceedings of the 11th Meeting of the North American Chapter of the Association for Computational Linguistics (NAACL) workshop on Statistical Parsing of Morphologically Rich Languages (SPMRL)*, Los Angeles.
- Joakim Nivre, Johan Hall, and Jens Nilsson. 2006. Malt-Parser: A Data-Driven Parser-Generator for Dependency Parsing. In *Proceedings of the Conference on Language Resources and Evaluation (LREC)*.
- Joakim Nivre. 2003. An efficient algorithm for projective dependency parsing. In *Proceedings of the 8th International Conference on Parsing Technologies (IWPT)*, pages 149–160, Nancy, France.
- Joakim Nivre. 2008. Algorithms for Deterministic Incremental Dependency Parsing. *Computational Linguistics*, 34(4).
- Franz Josef Och and Hermann Ney. 2003. A systematic comparison of various statistical alignment models. *Computational Linguistics*, 29(1):19–52.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*.
- Fatiha Sadat and Nizar Habash. 2006. Combination of arabic preprocessing schemes for statistical machine translation. In *Proceedings of the 21st International Conference on Computational Linguistics and the 44th annual meeting of the Association for Computational Linguistics*, pages 1–8, Morristown, NJ, USA.
- Matthew Snover, Bonnie Dorr, Richard Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of AMTA*, pages 223–231, Boston, MA.
- Chao Wang, Michael Collins, and Philipp Koehn. 2007. Chinese syntactic reordering for statistical machine translation. In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 737–745.
- Fei Xia and Michael McCord. 2004. Improving a statistical mt system with automatically learned rewrite patterns. In *Proceedings of COLING 2004*, pages 508–514, Geneva, Switzerland, August.
- Yuqi Zhang, Richard Zens, and Hermann Ney. 2007. Chunk-level reordering of source language sentences with automatically learned rules for statistical machine translation. In *Human Language Technology Conf. / North American Chapter of the Assoc. for Computational Linguistics Annual Meeting*, Rochester, NY, April.
- Imed Zitouni, Jeffrey S. Sorensen, and Ruhi Sarikaya. 2006. Maximum Entropy Based Restoration of Arabic Diacritics. In *Proceedings of COLING-ACL, the joint conference of the International Committee on Computational Linguistics and the Association for Computational Linguistics*, pages 577–584, Sydney, Australia.
- Andreas Zollmann, Ashish Venugopal, and Stephan Vogel. 2006. Bridging the inflection morphology gap for arabic statistical machine translation. In *Proceedings of the Human Language Technology Conference of the NAACL, Companion Volume: Short Papers*, pages 201–204, New York City, USA.

Learning Common Grammar from Multilingual Corpus

Tomoharu Iwata

Daichi Mochihashi

Hiroshi Sawada

NTT Communication Science Laboratories

2-4 Hikaridai, Seika-cho, Soraku-gun, Kyoto, Japan

{iwata, daichi, sawada}@cslab.kecl.ntt.co.jp

Abstract

We propose a corpus-based probabilistic framework to extract hidden common syntax across languages from non-parallel multilingual corpora in an unsupervised fashion. For this purpose, we assume a generative model for multilingual corpora, where each sentence is generated from a language dependent probabilistic context-free grammar (PCFG), and these PCFGs are generated from a prior grammar that is common across languages. We also develop a variational method for efficient inference. Experiments on a non-parallel multilingual corpus of eleven languages demonstrate the feasibility of the proposed method.

1 Introduction

Languages share certain common properties (Pinker, 1994). For example, the word order in most European languages is subject-verb-object (SVO), and some words with similar forms are used with similar meanings in different languages. The reasons for these common properties can be attributed to: 1) a common ancestor language, 2) borrowing from nearby languages, and 3) the innate abilities of humans (Chomsky, 1965).

We assume hidden commonalities in syntax across languages, and try to extract a common grammar from non-parallel multilingual corpora. For this purpose, we propose a generative model for multilingual grammars that is learned in an unsupervised fashion. There are some computational models for capturing commonalities at the phoneme and word level (Oakes, 2000; Bouchard-Côté et al., 2008), but, as far as we know, no attempt has been made to extract commonalities in syntax level from non-parallel and non-annotated multilingual corpora.

In our scenario, we use probabilistic context-free grammars (PCFGs) as our monolingual grammar model. We assume that a PCFG for each language is generated from a general model that are common across languages, and each sentence in multilingual corpora is generated from the language dependent PCFG. The inference of the general model as well as the multilingual PCFGs can be performed by using a variational method for efficiency. Our approach is based on a Bayesian multitask learning framework (Yu et al., 2005; Daumé III, 2009). Hierarchical Bayesian modeling provides a natural way of obtaining a joint regularization for individual models by assuming that the model parameters are drawn from a common prior distribution (Yu et al., 2005).

2 Related work

The unsupervised grammar induction task has been extensively studied (Carroll and Charniak, 1992; Stolcke and Omohundro, 1994; Klein and Manning, 2002; Klein and Manning, 2004; Liang et al., 2007). Recently, models have been proposed that outperform PCFG in the grammar induction task (Klein and Manning, 2002; Klein and Manning, 2004). We used PCFG as a first step for capturing commonalities in syntax across languages because of its simplicity. The proposed framework can be used for probabilistic grammar models other than PCFG.

Grammar induction using bilingual parallel corpora has been studied mainly in machine translation research (Wu, 1997; Melamed, 2003; Eisner, 2003; Chiang, 2005; Blunsom et al., 2009; Snyder et al., 2009). These methods require sentence-aligned parallel data, which can be costly to obtain and difficult to scale to many languages. On the other hand, our model does not require sentences to be aligned. Moreover, since the complexity of our model increases linearly with the number of languages, our model is easily applicable to cor-

pora of more than two languages, as we will show in the experiments. To our knowledge, the only grammar induction work on non-parallel corpora is (Cohen and Smith, 2009), but their method does not model a common grammar, and requires prior information such as part-of-speech tags. In contrast, our method does not require any such prior information.

3 Proposed Method

3.1 Model

Let $\mathbf{X} = \{\mathbf{X}_l\}_{l \in \mathbf{L}}$ be a non-parallel and non-annotated multilingual corpus, where \mathbf{X}_l is a set of sentences in language l , and \mathbf{L} is a set of languages. The task is to learn multilingual PCFGs $\mathbf{G} = \{\mathbf{G}_l\}_{l \in \mathbf{L}}$ and a common grammar that generates these PCFGs. Here, $\mathbf{G}_l = (\mathbf{K}, \mathbf{W}_l, \Phi_l)$ represents a PCFG of language l , where \mathbf{K} is a set of nonterminals, \mathbf{W}_l is a set of terminals, and Φ_l is a set of rule probabilities. Note that a set of nonterminals \mathbf{K} is shared among languages, but a set of terminals \mathbf{W}_l and rule probabilities Φ_l are specific to the language. For simplicity, we consider Chomsky normal form grammars, which have two types of rules: emissions rewrite a nonterminal as a terminal $A \rightarrow w$, and binary productions rewrite a nonterminal as two nonterminals $A \rightarrow BC$, where $A, B, C \in \mathbf{K}$ and $w \in \mathbf{W}_l$.

The rule probabilities for each nonterminal A of PCFG G_l in language l consist of: 1) $\theta_{lA} = \{\theta_{lAt}\}_{t \in \{0,1\}}$, where θ_{lA0} and θ_{lA1} represent probabilities of choosing the emission rule and the binary production rule, respectively, 2) $\phi_{lA} = \{\phi_{lABC}\}_{B,C \in \mathbf{K}}$, where ϕ_{lABC} represents the probability of nonterminal production $A \rightarrow BC$, and 3) $\psi_{lA} = \{\psi_{lAw}\}_{w \in \mathbf{W}_l}$, where ψ_{lAw} represents the probability of terminal emission $A \rightarrow w$. Note that $\theta_{lA0} + \theta_{lA1} = 1$, $\theta_{lAt} \geq 0$, $\sum_{B,C} \phi_{lABC} = 1$, $\phi_{lABC} \geq 0$, $\sum_w \psi_{lAw} = 1$, and $\psi_{lAw} \geq 0$. In the proposed model, multinomial parameters θ_{lA} and ϕ_{lA} are generated from Dirichlet distributions that are common across languages: $\theta_{lA} \sim \text{Dir}(\alpha_A^\theta)$ and $\phi_{lA} \sim \text{Dir}(\alpha_A^\phi)$, since we assume that languages share a common syntax structure. α_A^θ and α_A^ϕ represent the parameters of a common grammar. We use the Dirichlet prior because it is the conjugate prior for the multinomial distribution. In summary, the proposed model assumes the following generative process for a multilingual corpus,

1. For each nonterminal $A \in \mathbf{K}$:

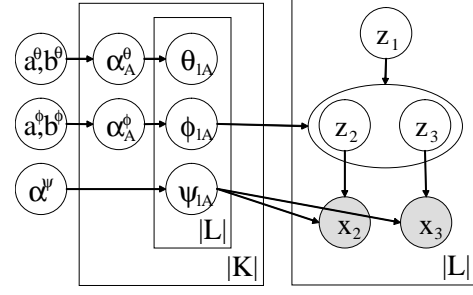


Figure 1: Graphical model.

- (a) For each rule type $t \in \{0, 1\}$:
 - i. Draw common rule type parameters $\alpha_{At}^\theta \sim \text{Gam}(a^\theta, b^\theta)$
 - (b) For each nonterminal pair (B, C) :
 - i. Draw common production parameters $\alpha_{ABC}^\phi \sim \text{Gam}(a^\phi, b^\phi)$
2. For each language $l \in \mathbf{L}$:
 - (a) For each nonterminal $A \in \mathbf{K}$:
 - i. Draw rule type parameters $\theta_{lA} \sim \text{Dir}(\alpha_A^\theta)$
 - ii. Draw binary production parameters $\phi_{lA} \sim \text{Dir}(\alpha_A^\phi)$
 - iii. Draw emission parameters $\psi_{lA} \sim \text{Dir}(\alpha^\psi)$
 - (b) For each node i in the parse tree:
 - i. Choose rule type $t_{li} \sim \text{Mult}(\theta_{lz_i})$
 - ii. If $t_{li} = 0$:
 - A. Emit terminal $x_{li} \sim \text{Mult}(\psi_{lz_i})$
 - iii. Otherwise:
 - A. Generate children nonterminals $(z_{lL(i)}, z_{lR(i)}) \sim \text{Mult}(\phi_{lz_i})$,

where $L(i)$ and $R(i)$ represent the left and right children of node i . Figure 1 shows a graphical model representation of the proposed model, where the shaded and unshaded nodes indicate observed and latent variables, respectively.

3.2 Inference

The inference of the proposed model can be efficiently computed using a variational Bayesian method. We extend the variational method to the monolingual PCFG learning of Kurihara and Sato (2004) for multilingual corpora. The goal is to estimate posterior $p(\mathbf{Z}, \Phi, \alpha | \mathbf{X})$, where \mathbf{Z} is a set of parse trees, $\Phi = \{\Phi_l\}_{l \in \mathbf{L}}$ is a set of language dependent parameters, $\Phi_l = \{\theta_{lA}, \phi_{lA}, \psi_{lA}\}_{A \in \mathbf{K}}$, and $\alpha = \{\alpha_A^\theta, \alpha_A^\phi\}_{A \in \mathbf{K}}$ is a set of common parameters. In the variational method, posterior $p(\mathbf{Z}, \Phi, \alpha | \mathbf{X})$ is approximated by a tractable variational distribution $q(\mathbf{Z}, \Phi, \alpha)$.

We use the following variational distribution,

$$q(\mathbf{Z}, \Phi, \alpha) = \prod_A q(\alpha_A^\theta) q(\alpha_A^\phi) \prod_{l,d} q(\mathbf{z}_{ld}) \\ \times \prod_{l,A} q(\theta_{lA}) q(\phi_{lA}) q(\psi_{lA}), \quad (1)$$

where we assume that hyperparameters $q(\alpha_A^\theta)$ and $q(\alpha_A^\phi)$ are degenerated, or $q(\alpha) = \delta_{\alpha^*}(\alpha)$, and infer them by point estimation instead of distribution estimation. We find an approximate posterior distribution that minimizes the Kullback-Leibler divergence from the true posterior. The variational distribution of the parse tree of the d th sentence in language l is obtained as follows,

$$q(\mathbf{z}_{ld}) \propto \prod_{A \rightarrow BC} \left(\pi_{lA1}^\theta \pi_{lABC}^\phi \right)^{C(A \rightarrow BC; \mathbf{z}_{ld}, l, d)} \\ \times \prod_{A \rightarrow w} \left(\pi_{lA0}^\theta \pi_{lAw}^\psi \right)^{C(A \rightarrow w; \mathbf{z}_{ld}, l, d)}, \quad (2)$$

where $C(r; \mathbf{z}, l, d)$ is the count of rule r that occurs in the d th sentence of language l with parse tree \mathbf{z} . The multinomial weights are calculated as follows,

$$\pi_{lAt}^\theta = \exp(\mathbb{E}_{q(\theta_{lA})} [\log \theta_{lAt}]), \quad (3)$$

$$\pi_{lABC}^\phi = \exp(\mathbb{E}_{q(\phi_{lA})} [\log \phi_{lABC}]), \quad (4)$$

$$\pi_{lAw}^\psi = \exp(\mathbb{E}_{q(\psi_{lA})} [\log \psi_{lAw}]). \quad (5)$$

The variational Dirichlet parameters for $q(\theta_{lA}) = \text{Dir}(\gamma_{lA}^\theta)$, $q(\phi_{lA}) = \text{Dir}(\gamma_{lA}^\phi)$, and $q(\psi_{lA}) = \text{Dir}(\gamma_{lA}^\psi)$, are obtained as follows,

$$\gamma_{lAt}^\theta = \alpha_{At}^\theta + \sum_{d, \mathbf{z}_{ld}} q(\mathbf{z}_{ld}) C(A, t; \mathbf{z}_{ld}, l, d), \quad (6)$$

$$\gamma_{lABC}^\phi = \alpha_{ABC}^\phi + \sum_{d, \mathbf{z}_{ld}} q(\mathbf{z}_{ld}) C(A \rightarrow BC; \mathbf{z}_{ld}, l, d), \quad (7)$$

$$\gamma_{lAw}^\psi = \alpha_{Aw}^\psi + \sum_{d, \mathbf{z}_{ld}} q(\mathbf{z}_{ld}) C(A \rightarrow w; \mathbf{z}_{ld}, l, d), \quad (8)$$

where $C(A, t; \mathbf{z}, l, d)$ is the count of rule type t that is selected in nonterminal A in the d th sentence of language l with parse tree \mathbf{z} .

The common rule type parameter α_{At}^θ that minimizes the KL divergence between the true posterior and the approximate posterior can be obtained by using the fixed-point iteration method

described in (Minka, 2000). The update rule is as follows,

$$\alpha_{At}^{\theta(\text{new})} \leftarrow \frac{a^\theta - 1 + \alpha_{At}^\theta L(\Psi(\sum_{t'} \alpha_{At'}^\theta) - \Psi(\alpha_{At}^\theta))}{b^\theta + \sum_l (\Psi(\sum_{t'} \gamma_{lAt'}^\theta) - \Psi(\gamma_{lAt}^\theta))}, \quad (9)$$

where L is the number of languages, and $\Psi(x) = \frac{\partial \log \Gamma(x)}{\partial x}$ is the digamma function. Similarly, the common production parameter α_{ABC}^ϕ can be updated as follows,

$$\alpha_{ABC}^{\phi(\text{new})} \leftarrow \frac{a^\phi - 1 + \alpha_{ABC}^\phi L J_{ABC}}{b^\phi + \sum_l J'_{lABC}}, \quad (10)$$

where $J_{ABC} = \Psi(\sum_{B', C'} \alpha_{AB'C'}^\phi) - \Psi(\alpha_{ABC}^\phi)$, and $J'_{lABC} = \Psi(\sum_{B', C'} \gamma_{lAB'C'}^\phi) - \Psi(\gamma_{lABC}^\phi)$.

Since factored variational distributions depend on each other, an optimal approximated posterior can be obtained by updating parameters by (2) - (10) alternatively until convergence. The updating of language dependent distributions by (2) - (8) is also described in (Kurihara and Sato, 2004; Liang et al., 2007) while the updating of common grammar parameters by (9) and (10) is new. The inference can be carried out efficiently using the inside-outside algorithm based on dynamic programming (Lari and Young, 1990).

After the inference, the probability of a common grammar rule $A \rightarrow BC$ is calculated by $\hat{\phi}_{A \rightarrow BC} = \hat{\theta}_1 \hat{\phi}_{ABC}$, where $\hat{\theta}_1 = \alpha_1^\theta / (\alpha_0^\theta + \alpha_1^\theta)$ and $\hat{\phi}_{ABC} = \alpha_{ABC}^\phi / \sum_{B', C'} \alpha_{AB'C'}^\phi$ represent the mean values of θ_{l0} and ϕ_{lABC} , respectively.

4 Experimental results

We evaluated our method by employing the EuroParl corpus (Koehn, 2005). The corpus consists of the proceedings of the European Parliament in eleven western European languages: Danish (da), German (de), Greek (el), English (en), Spanish (es), Finnish (fi), French (fr), Italian (it), Dutch (nl), Portuguese (pt), and Swedish (sv), and it contains roughly 1,500,000 sentences in each language. We set the number of nonterminals at $|\mathbf{K}| = 20$, and omitted sentences with more than ten words for tractability. We randomly sampled 100,000 sentences for each language, and analyzed them using our method. It should be noted that our random samples are not sentence-aligned.

Figure 2 shows the most probable terminals of emission for each language and nonterminal with a high probability of selecting the emission rule.

2: verb and auxiliary verb (V)

[da] det jeg vi der de derfor dette forhandlingen hvad
[de] ist sind haben wird hat müssen möchte meine werden kann
[el] , είναι για να πρέπει δεν και αυτό με ότι
[en] is are , will have has must was should you
[es] es hay gracias mañana qué tiene tenemos trata lugar son
[fi] on ei ovat olisi oli toimitetaan eivät voi koskee voidaan
[fr] ' ne est n' en a lieu aura vous avons
[it] , è che non si discussione svolgerà presidente ' sono
[nl] is zijn moeten heeft hebben moet kan zal wil wordt
[pt] , que é não senhor parlamento de aprova amanhã com
[sv] det jag vi detta vad därför de debatten ni den

5: noun (N)

[da] det dem her dette sig dag os noget betænkningen over
[de] abstimmung aussprache kommission bericht frage parlament
[el] θέμα έκθεση τροπολογία κ. Επιτροπή ψήφισμα πρόταση
[en] vote debate president commission much like council minutes
[es] debate presidente informe parlamento ello comisario sr.
[fi] unionin " yhteisön hetkellä enemmän uudelleen kerran heidän
[fr] vote débat parlement rapport commission question président
[it] votazione parlamento commissione relazione risoluzione
[nl] debat stemming commissie parlement verslag voorzitter
[pt] votação comissão questão relatório situação sessão proposta
[sv] det rum damer här detta oss sig frågan första dem

7: subject (SBJ)

[da] er har vil skal kan må var finder bør ville
[de] ich wir das , es sie - dies vielen (
[el] , θα δεν ψηφοφορία ότι Σώμα πολύ αυτό Πρόεδρε Το
[en] we i that it this there what (thank they
[es] no se ¿ esto por lo me pero muchas tendrá
[fi] , että puhemies kiitoksia mietintö joka parlamentin mitä
[fr] nous je il c' cela j' ce l mais vous
[it] non (la e questo ma si vorrei signor mi
[nl] dat ik wij het er we dit u daar wat
[pt] o , que encerrado presidente terá obrigado lugar isso não
[sv] är har måste kommer kan vill skall finns skulle var

9: preposition (PR)

[da] , af for i og til på med om fra
[de] und , in für auf von zu mit an auch
[el] και / (σε , που από είναι των για
[en] to of in , for not and on with take
[es] , de que a en por con y para sobre
[fi] ja euroopan , on kuin / : tai ovat eikä
[fr] , de à que pour sur dans d' et par
[it] di e della del in a (dell' dei da
[nl] van in , voor en op met aan over maar
[pt] de da do e para em dos / é com
[sv] i och för av till på om med som :

11: punctuation (.)

[da] . ?) ! : f.eks. ... sessionen vedtoges bl.a.
[de] . ! ?) : protokolls ... sitzungsperiode " ...
[el] .) ! ; : κ. προκτικών ... π.μ. κ.κ.
[en] . ?) ! a.m. : p.m. ... ' ;
[es] . ?) ! : ... ; » " anterior
[fi] . ?) ! : ... " ...; ääntä
[fr] . ?) ! la : ... » ; ...
[it] . ?) ! : sessione ... " ; precedente
[nl] . ?) ! : zitting ... gesloten " onderbroken
[pt] . ?) ! : ... " anterior ; urgentes
[sv] . !) ? : protokoll ... sessionen t.ex. "

13: determiner (DT)

[da] ikke at en den også de et gerne det være
[de] die der eine den das ein diese im des dieser
[el] να το την η τη τις είναι τα στην μεα
[en] the a be mr very an been not no in
[es] el ha este señor un se hemos debemos han debe
[fi] ole myös kuitenkin vielä hyvin siis erittäin nyt jo vain
[fr] le la les l' une cette un ce ces des
[it] la il l' un una le in i a gli
[nl] de het een deze dit geen die onze mijn mijnheer
[pt] a o uma um os de as esta este em
[sv] att inte en mycket också ett för vara om äga

Figure 2: Probable terminals of emission for each language and nonterminal.

0 → 16 11	(R → S .)	0.11
16 → 7 6	(S → SBJ VP)	0.06
6 → 2 12	(VP → V NP)	0.04
12 → 13 5	(NP → DT N)	0.19
15 → 17 19	(NP → NP N)	0.07
17 → 5 9	(NP → N PR)	0.07
15 → 13 5	(NP → DT N)	0.06

Figure 3: Examples of inferred common grammar rules in eleven languages, and their probabilities. Hand-provided annotations have the following meanings, R: root, S: sentence, NP: noun phrase, VP: verb phrase, and others appear in Figure 2.

We named nonterminals by using grammatical categories after the inference. We can see that words in the same grammatical category clustered across languages as well as within a language. Figure 3 shows examples of inferred common grammar rules with high probabilities. Grammar rules that seem to be common to European languages have been extracted.

5 Discussion

We have proposed a Bayesian hierarchical PCFG model for capturing commonalities at the syntax level for non-parallel multilingual corpora. Although our results have been encouraging, a number of directions remain in which we must extend our approach. First, we need to evaluate our model quantitatively using corpora with a greater diversity of languages. Measurement examples include the perplexity, and machine translation score. Second, we need to improve our model. For example, we can infer the number of nonterminals with a nonparametric Bayesian model (Liang et al., 2007), infer the model more robustly based on a Markov chain Monte Carlo inference (Johnson et al., 2007), and use probabilistic grammar models other than PCFGs. In our model, all the multilingual grammars are generated from a general model. We can extend it hierarchically using the coalescent (Kingman, 1982). That model may help to infer an evolutionary tree of languages in terms of grammatical structure without the etymological information that is generally used (Gray and Atkinson, 2003). Finally, the proposed approach may help to indicate the presence of a universal grammar (Chomsky, 1965), or to find it.

References

- Phil Blunsom, Trevor Cohn, and Miles Osborne. 2009. Bayesian synchronous grammar induction. In D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou, editors, *Advances in Neural Information Processing Systems 21*, pages 161–168.
- Alexandre Bouchard-Côté, Percy Liang, Thomas Griffiths, and Dan Klein. 2008. A probabilistic approach to language change. In J.C. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems 20*, pages 169–176, Cambridge, MA. MIT Press.
- Glenn Carroll and Eugene Charniak. 1992. Two experiments on learning probabilistic dependency grammars from corpora. In *Working Notes of the Workshop Statistically-Based NLP Techniques*, pages 1–13. AAAI.
- David Chiang. 2005. A hierarchical phrase-based model for statistical machine translation. In *ACL '05: Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics*, pages 263–270, Morristown, NJ, USA. Association for Computational Linguistics.
- Norm Chomsky. 1965. *Aspects of the Theory of Syntax*. MIT Press.
- Shay B. Cohen and Noah A. Smith. 2009. Shared logistic normal distributions for soft parameter tying in unsupervised grammar induction. In *NAACL '09: Proceedings of the Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 74–82, Morristown, NJ, USA. Association for Computational Linguistics.
- Hal Daumé III. 2009. Bayesian multitask learning with latent hierarchies. In *Proceedings of the Twenty-Fifth Annual Conference on Uncertainty in Artificial Intelligence (UAI-09)*, pages 135–142, Corvallis, Oregon. AUAI Press.
- Jason Eisner. 2003. Learning non-isomorphic tree mappings for machine translation. In *ACL '03: Proceedings of the 41st Annual Meeting on Association for Computational Linguistics*, pages 205–208, Morristown, NJ, USA. Association for Computational Linguistics.
- Russell D. Gray and Quentin D. Atkinson. 2003. Language-tree divergence times support the Anatolian theory of Indo-European origin. *Nature*, 426(6965):435–439, November.
- Mark Johnson, Thomas Griffiths, and Sharon Goldwater. 2007. Bayesian inference for PCFGs via Markov chain Monte Carlo. In *Human Language Technologies 2007: The Conference of the North American Chapter of the Association for Computational Linguistics; Proceedings of the Main Conference*, pages 139–146, Rochester, New York, April. Association for Computational Linguistics.
- J. F. C. Kingman. 1982. The coalescent. *Stochastic Processes and their Applications*, 13:235–248.
- Dan Klein and Christopher D. Manning. 2002. A generative constituent-context model for improved grammar induction. In *ACL '02: Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, pages 128–135, Morristown, NJ, USA. Association for Computational Linguistics.
- Dan Klein and Christopher D. Manning. 2004. Corpus-based induction of syntactic structure: models of dependency and constituency. In *ACL '04: Proceedings of the 42nd Annual Meeting on Association for Computational Linguistics*, page 478, Morristown, NJ, USA. Association for Computational Linguistics.
- Philipp Koehn. 2005. Europarl: A parallel corpus for statistical machine translation. In *Proceedings of the 10th Machine Translation Summit*, pages 79–86.
- Kenichi Kurihara and Taisuke Sato. 2004. An application of the variational Bayesian approach to probabilistic context-free grammars. In *International Joint Conference on Natural Language Processing Workshop Beyond Shallow Analysis*.
- K. Lari and S.J. Young. 1990. The estimation of stochastic context-free grammars using the inside-outside algorithm. *Computer Speech and Language*, 4:35–56.
- Percy Liang, Slav Petrov, Michael I. Jordan, and Dan Klein. 2007. The infinite PCFG using hierarchical dirichlet processes. In *EMNLP '07: Proceedings of the Empirical Methods on Natural Language Processing*, pages 688–697.
- I. Dan Melamed. 2003. Multitext grammars and synchronous parsers. In *NAACL '03: Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology*, pages 79–86, Morristown, NJ, USA. Association for Computational Linguistics.
- Thomas Minka. 2000. Estimating a Dirichlet distribution. Technical report, M.I.T.
- Michael P. Oakes. 2000. Computer estimation of vocabulary in a protolanguage from word lists in four daughter languages. *Journal of Quantitative Linguistics*, 7(3):233–243.
- Steven Pinker. 1994. *The Language Instinct: How the Mind Creates Language*. HarperCollins, New York.
- Benjamin Snyder, Tahira Naseem, and Regina Barzilay. 2009. Unsupervised multilingual grammar induction. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, pages 73–81, Suntec, Singapore, August. Association for Computational Linguistics.
- Andreas Stolcke and Stephen M. Omohundro. 1994. Inducing probabilistic grammars by Bayesian model merging. In *ICGI '94: Proceedings of the Second International Colloquium on Grammatical Inference and Applications*, pages 106–118, London, UK. Springer-Verlag.
- Dekai Wu. 1997. Stochastic inversion transduction grammars and bilingual parsing of parallel corpora. *Comput. Linguist.*, 23(3):377–403.
- Kai Yu, Volker Tresp, and Anton Schwaighofer. 2005. Learning gaussian processes from multiple tasks. In *ICML '05: Proceedings of the 22nd International Conference on Machine Learning*, pages 1012–1019, New York, NY, USA. ACM.

Tree-Based Deterministic Dependency Parsing — An Application to Nivre’s Method —

Kotaro Kitagawa Kumiko Tanaka-Ishii

Graduate School of Information Science and Technology,

The University of Tokyo

kitagawa@cl.ci.i.u-tokyo.ac.jp kumiko@i.u-tokyo.ac.jp

Abstract

Nivre’s method was improved by enhancing deterministic dependency parsing through application of a tree-based model. The model considers all words necessary for selection of parsing actions by including words in the form of trees. It chooses the most probable head candidate from among the trees and uses this candidate to select a parsing action.

In an evaluation experiment using the Penn Treebank (WSJ section), the proposed model achieved higher accuracy than did previous deterministic models. Although the proposed model’s worst-case time complexity is $O(n^2)$, the experimental results demonstrated an average parsing time not much slower than $O(n)$.

1 Introduction

Deterministic parsing methods achieve both effective time complexity and accuracy not far from those of the most accurate methods. One such deterministic method is Nivre’s method, an incremental parsing method whose time complexity is linear in the number of words (Nivre, 2003). Still, deterministic methods can be improved. As a specific example, Nivre’s model greedily decides the parsing action only from two words and their locally relational words, which can lead to errors.

In the field of Japanese dependency parsing, Iwatate et al. (2008) proposed a tournament model that takes all head candidates into account in judging dependency relations. This method assumes backward parsing because the Japanese dependency structure has a head-final constraint, so that any word’s head is located to its right.

Here, we propose a tree-based model, applicable to any projective language, which can be considered as a kind of generalization of Iwatate’s

idea. Instead of selecting a parsing action for two words, as in Nivre’s model, our tree-based model first chooses the most probable head candidate from among the trees through a tournament and then decides the parsing action between two trees.

Global-optimization parsing methods are another common approach (Eisner, 1996; McDonald et al., 2005). Koo et al. (2008) studied semi-supervised learning with this approach. Hybrid systems have improved parsing by integrating outputs obtained from different parsing models (Zhang and Clark, 2008).

Our proposal can be situated among global-optimization parsing methods as follows. The proposed tree-based model is deterministic but takes a step towards global optimization by widening the search space to include all necessary words connected by previously judged head-dependent relations, thus achieving a higher accuracy yet largely retaining the speed of deterministic parsing.

2 Deterministic Dependency Parsing

2.1 Dependency Parsing

A dependency parser receives an input sentence $x = w_1, w_2, \dots, w_n$ and computes a dependency graph $G = (W, A)$. The set of nodes $W = \{w_0, w_1, \dots, w_n\}$ corresponds to the words of a sentence, and the node w_0 is the root of G . A is the set of arcs (w_i, w_j) , each of which represents a dependency relation where w_i is the *head* and w_j is the *dependent*.

In this paper, we assume that the resulting dependency graph for a sentence is well-formed and projective (Nivre, 2008). G is well-formed if and only if it satisfies the following three conditions of being **single-headed**, **acyclic**, and **rooted**.

2.2 Nivre’s Method

An incremental dependency parsing algorithm was first proposed by (Covington, 2001). After

Table 1: Transitions for Nivre’s method and the proposed method.

		Transition	Precondition
Nivre’s Method	Left-Arc	$(\sigma w_i, w_j \beta, A) \Rightarrow (\sigma, w_j \beta, A \cup \{(w_j, w_i)\})$	$i \neq 0 \wedge \neg \exists w_k (w_k, w_i) \in A$ $\exists w_k (w_k, w_i) \in A$
	Right-Arc	$(\sigma w_i, w_j \beta, A) \Rightarrow (\sigma w_i w_j, \beta, A \cup \{(w_i, w_j)\})$	
	Reduce	$(\sigma w_i, \beta, A) \Rightarrow (\sigma, \beta, A)$	
	Shift	$(\sigma, w_j \beta, A) \Rightarrow (\sigma w_j, \beta, A)$	
Proposed Method	Left-Arc	$(\sigma t_i, t_j \beta, A) \Rightarrow (\sigma, t_j \beta, A \cup \{(w_j, w_i)\})$	$i \neq 0$
	Right-Arc	$(\sigma t_i, t_j \beta, A) \Rightarrow (\sigma t_i, \beta, A \cup \{(mphc(t_i, t_j), w_j)\})$	
	Shift	$(\sigma, t_j \beta, A) \Rightarrow (\sigma t_j, \beta, A)$	

studies taking data-driven approaches, by (Kudo and Matsumoto, 2002), (Yamada and Matsumoto, 2003), and (Nivre, 2003), the deterministic incremental parser was generalized to a state transition system in (Nivre, 2008).

Nivre’s method applying an arc-eager algorithm works by using a stack of words denoted as σ , for a buffer β initially containing the sentence x . Parsing is formulated as a quadruple (S, T_s, s_{init}, S_t) , where each component is defined as follows:

- S is a set of states, each of which is denoted as $(\sigma, \beta, A) \in S$.
- T_s is a set of transitions, and each element of T_s is a function $t_s : S \rightarrow S$.
- $s_{init} = ([w_0], [w_1, \dots, w_n], \phi)$ is the initial state.
- S_t is a set of terminal states.

Syntactic analysis generates a sequence of optimal transitions t_s provided by an oracle $o : S \rightarrow T_s$, applied to a target consisting of the stack’s top element w_i and the first element w_j in the buffer. The oracle is constructed as a classifier trained on tree-bank data. Each transition is defined in the upper block of Table 1 and explained as follows:

Left-Arc Make w_j the head of w_i and pop w_i , where w_i is located at the stack top (denoted as $\sigma|w_i$), when the buffer head is w_j (denoted as $w_j|\beta$).

Right-Arc Make w_i the head of w_j , and push w_j .

Reduce Pop w_i , located at the stack top.

Shift Push the word w_j , located at the buffer head, onto the stack top.

The method explained thus far has the following drawbacks.

Locality of Parsing Action Selection

The dependency relations are greedily determined, so when the transition Right-Arc adds a dependency arc (w_i, w_j) , a more probable head of w_j located in the stack is disregarded as a candidate.

Features Used for Selecting Reduce

The features used in (Nivre and Scholz, 2004) to define a state transition are basically obtained from the two target words w_i and w_j , and their related words. These words are not sufficient to select Reduce, because this action means that w_j has no dependency relation with any word in the stack.

Preconditions

When the classifier selects a transition, the resulting graph satisfies well-formedness and projectivity only under the preconditions listed in Table 1. Even though the parsing seems to be formulated as a four-class classifier problem, it is in fact formed of two types of three-class classifiers.

Solving these problems and selecting a more suitable dependency relation requires a parser that considers more global dependency relations.

3 Tree-Based Parsing Applied to Nivre’s Method

3.1 Overall Procedure

Tree-based parsing uses trees as the procedural elements instead of words. This allows enhancement of previously proposed deterministic models such as (Covington, 2001; Yamada and Matsumoto, 2003). In this paper, we show the application of tree-based parsing to Nivre’s method. The parser is formulated as a state transition system (S, T_s, s_{init}, S_t) , similarly to Nivre’s parser, but σ and β for a state $s = (\sigma, \beta, A) \in S$ denote a stack of trees and a buffer of trees, respectively. A tree $t_i \in T$ is defined as the tree rooted by the word w_i , and the initial state is $s_{init} = ([t_0], [t_1, \dots, t_n], \phi)$, which is formed from the input sentence x .

The state transitions T_s are decided through the following two steps.

1. **Select the most probable head candidate (MPHC):** For the tree t_i located at the stack top, search for and select the MPHC for w_j , which is the root word of t_j located at the buffer head. This procedure is denoted as a

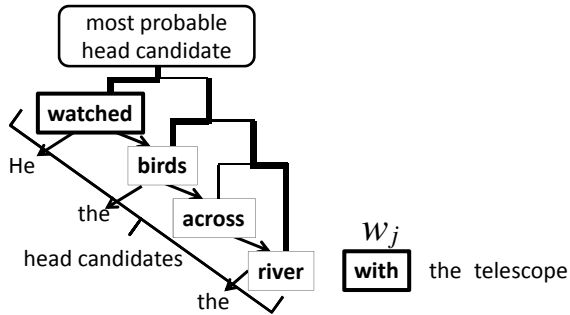


Figure 1: Example of a tournament.

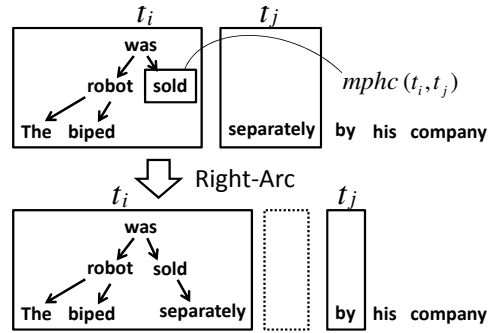


Figure 2: Example of the transition Right.

function $mphc(t_i, t_j)$, and its details are explained in §3.2.

2. **Select a transition:** Choose a transition, by using an oracle, from among the following three possibilities (explained in detail in §3.3):

Left-Arc Make w_j the head of w_i and pop t_i , where t_i is at the stack top (denoted as $\sigma|t_i$, with the tail being σ), when the buffer head is t_j (denoted as $t_j|\beta$).

Right-Arc Make the MPHC the head of w_j , and pop the MPHC.

Shift Push the tree t_j located at the buffer head onto the stack top.

These transitions correspond to three possibilities for the relation between t_i and t_j : (1) a word of t_i is a dependent of a word of t_j ; (2) a word of t_j is a dependent of a word of t_i ; or (3) the two trees are not related.

The formulations of these transitions in the lower block of Table 1 correspond to Nivre’s transitions of the same name, except that here a transition is applied to a tree. This enhancement from words to trees allows removal of both the Reduce transition and certain preconditions.

3.2 Selection of Most Probable Head Candidate

By using $mphc(t_i, t_j)$, a word located far from w_j (the head of t_j) can be selected as the head candidate in t_i . This selection process decreases the number of errors resulting from greedy decision considering only a few candidates.

Various procedures can be considered for implementing $mphc(t_i, t_j)$. One way is to apply the tournament procedure to the words in t_i . The tournament procedure was originally introduced for parsing methods in Japanese by (Iwatate et al.,

2008). Since the Japanese language has the head-final property, the tournament model itself constitutes parsing, whereas for parsing a general projective language, the tournament model can only be used as part of a parsing algorithm.

Figure 1 shows a tournament for the example of “with,” where the word “watched” finally wins. Although only the words on the left-hand side of tree t_j are searched, this does not mean that the tree-based method considers only one side of a dependency relation. For example, when we apply the tree-based parsing to Yamada’s method, the search problems on both sides are solved.

To implement $mphc(t_i, t_j)$, a binary classifier is built to judge which of two given words is more appropriate as the head for another input word. This classifier concerns three words, namely, the two words l (left) and r (right) in t_i , whose appropriateness as the head is compared for the dependent w_j . All word pairs of l and r in t_i are compared repeatedly in a “tournament,” and the survivor is regarded as the MPHC of w_j .

The classifier is generated through learning of training examples for all t_i and w_j pairs, each of which generates examples comparing the true head and other (inappropriate) heads in t_i . Table 2 lists the features used in the classifier. Here, $\text{lex}(X)$ and $\text{pos}(X)$ mean the surface form and part of speech of X , respectively. X^{left} means the dependents of X located on the left-hand side of X , while X^{right} means those on the right. Also, X^{head} means the head of X . The feature design concerns three additional words occurring after w_j , as well, denoted as $w_{j+1}, w_{j+2}, w_{j+3}$.

3.3 Transition Selection

A transition is selected by a three-class classifier after deciding the MPHC, as explained in §3.1. Table 1 lists the three transitions and one precon-

Table 2: Features used for a tournament.

pos(l), lex(l) pos(l^{head}), pos(l^{left}), pos(l^{right})
pos(r), lex(r) pos(r^{head}), pos(r^{left}), pos(r^{right})
pos(w_j), lex(w_j), pos(w_j^{left})
pos(w_{j+1}), lex(w_{j+1}), pos(w_{j+2}), lex(w_{j+2}) pos(w_{j+3}), lex(w_{j+3})

Table 3: Features used for a state transition.

pos(w_i), lex(w_i) pos(w_i^{left}), pos(w_i^{right}), lex(w_i^{left}), lex(w_i^{right})
pos(MPHC), lex(MPHC) pos(MPHC ^{head}), pos(MPHC ^{left}), pos(MPHC ^{right}) lex(MPHC ^{head}), lex(MPHC ^{left}), lex(MPHC ^{right})
pos(w_j), lex(w_j), pos(w_j^{left}), lex(w_j^{left})
pos(w_{j+1}), lex(w_{j+1}), pos(w_{j+2}), lex(w_{j+2}), pos(w_{j+3}), lex(w_{j+3})

dition. The transition Shift indicates that the target trees t_i and t_j have no dependency relations. The transition Right-Arc indicates generation of the dependent-head relation between w_j and the result of $mphc(t_i, t_j)$, i.e., the MPHC for w_j . Figure 2 shows an example of this transition. The transition Left-Arc indicates generation of the dependency relation in which w_j is the head of w_i . While Right-Arc requires searching for the MPHC in t_i , this is not the case for Left-Arc¹.

The key to obtaining an accurate tree-based parsing model is to extend the search space while at the same time providing ways to narrow down the space and find important information, such as the MPHC, for proper judgment of transitions.

The three-class classifier is constructed as follows. The dependency relation between the target trees is represented by the three words w_i , MPHC, and w_j . Therefore, the features are designed to incorporate these words, their relational words, and the three words next to w_j . Table 3 lists the exact set of features used in this work. Since this transition selection procedure presumes selection of the MPHC, the result of $mphc(t_i, t_j)$ is also incorporated among the features.

4 Evaluation

4.1 Data and Experimental Setting

In our experimental evaluation, we used Yamada’s head rule to extract unlabeled dependencies from the Wall Street Journal section of a Penn Treebank. Sections 2-21 were used as the training data, and section 23 was used as the test data. This test data

¹The head word of w_i can only be w_j without searching within t_j , because the relations between the other words in t_j and w_i have already been inferred from the decisions made within previous transitions. If t_j has a child w_k that could become the head of w_i under projectivity, this w_k must be located between w_i and w_j . The fact that w_k ’s head is w_j means that there were two phases before t_i and t_j (i.e., w_i and w_j) became the target:

- t_i and t_k became the target, and Shift was selected.
- t_k and t_j became the target, and Left-Arc was selected.

The first phase precisely indicates that w_i and w_k are unrelated.

was used in several other previous works, enabling mutual comparison with the methods reported in those works.

The SVM^{light} package² was used to build the support vector machine classifiers. The binary classifier for MPHC selection and the three-class classifier for transition selection were built using a cubic polynomial kernel. The parsing speed was evaluated on a Core2Duo (2.53 GHz) machine.

4.2 Parsing Accuracy

We measured the ratio of words assigned correct heads to all words (accuracy), and the ratio of sentences with completely correct dependency graphs to all sentences (complete match). In the evaluation, we consistently excluded punctuation marks.

Table 4 compares our results for the proposed method with those reported in some previous works using equivalent training and test data. The first column lists the four previous methods and our method, while the second through fourth columns list the accuracy, complete match accuracy, and time complexity, respectively, for each method. Here, we obtained the scores for the previous works from the corresponding articles listed in the first column. Note that every method used different features, which depend on the method.

The proposed method achieved higher accuracy than did the previous deterministic models. Although the accuracy of our method did not reach that of (McDonald and Pereira, 2006), the scores were competitive even though our method is deterministic. These results show the capability of the tree-based approach in effectively extending the search space.

4.3 Parsing Time

Such extension of the search space also concerns the speed of the method. Here, we compare its computational time with that of Nivre’s method. We re-implemented Nivre’s method to use SVMs with cubic polynomial kernel, similarly to our

²<http://svmlight.joachims.org/>

Table 4: Dependency parsing performance.

	Accuracy	Complete match	Time complexity	Global vs. deterministic	Learning method
McDonald & Pereira (2006)	91.5	42.1	$O(n^3)$	global	MIRA
McDonald et al. (2005)	90.9	37.5	$O(n^3)$	global	MIRA
Yamada & Matsumoto (2003)	90.4	38.4	$O(n^2)$	deterministic	support vector machine
Goldberg & Elhadad (2010)	89.7	37.5	$O(n \log n)$	deterministic	structured perceptron
Nivre (2004)	87.1	30.4	$O(n)$	deterministic	memory based learning
Proposed method	91.3	41.7	$O(n^2)$	deterministic	support vector machine

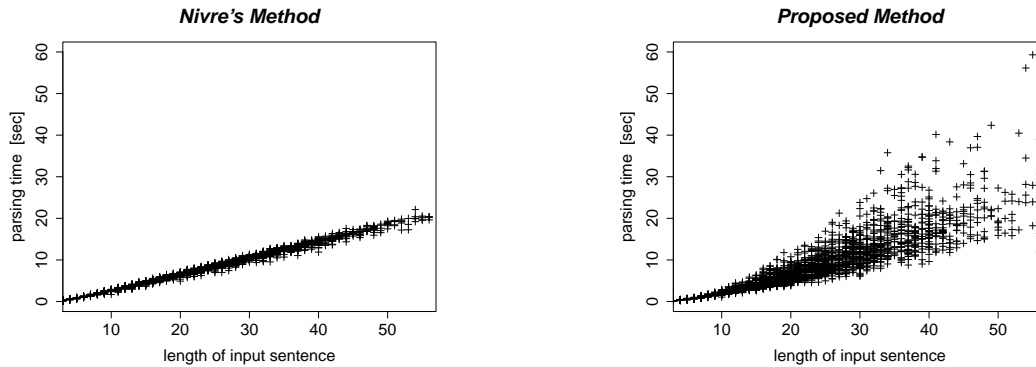


Figure 3: Parsing time for sentences.

method. Figure 3 shows plots of the parsing times for all sentences in the test data. The average parsing time for our method was 8.9 sec, whereas that for Nivre’s method was 7.9 sec.

Although the worst-case time complexity for Nivre’s method is $O(n)$ and that for our method is $O(n^2)$, worst-case situations (e.g., all words having heads on their left) did not appear frequently. This can be seen from the sparse appearance of the upper bound in the second figure.

5 Conclusion

We have proposed a tree-based model that decides head-dependency relations between trees instead of between words. This extends the search space to obtain the best head for a word within a deterministic model. The tree-based idea is potentially applicable to various previous parsing methods; in this paper, we have applied it to enhance Nivre’s method.

Our tree-based model outperformed various deterministic parsing methods reported previously. Although the worst-case time complexity of our method is $O(n^2)$, the average parsing time is not much slower than $O(n)$.

References

Xavier Carreras. 2007. Experiments with a higher-order projective dependency parse. *Proceedings of the CoNLL Shared Task Session of EMNLP-CoNLL*, pp. 957-961.

Michael A. Covington. 2001. A fundamental algorithm for dependency parsing. *Proceedings of ACM*, pp. 95-102.

Jason M. Eisner. 1996. Three new probabilistic models for dependency parsing: An exploration. *Proceedings of COLING*, pp. 340-345.

Yoav Goldberg and Michael Elhadad. 2010. An Efficient Algorithm for Easy-First Non-Directional Dependency Parsing. *Proceedings of NAACL*.

Masakazu Iwatate, Masayuki Asahara, and Yuji Matsumoto. 2008. Japanese dependency parsing using a tournament model. *Proceedings of COLING*, pp. 361-368.

Terry Koo, Xavier Carreras, and Michael Collins. 2008. Simple semi-supervised dependency parsing. *Proceedings of ACL*, pp. 595-603.

Taku Kudo and Yuji Matsumoto. 2002. Japanese dependency analysis using cascaded chunking. *Proceedings of CoNLL*, pp. 63-69.

Ryan McDonald, Koby Crammer, and Fernando Pereira. 2005. Online large-margin training of dependency parsers. *Proceedings of ACL*, pp. 91-98.

Ryan McDonald and Fernando Pereira. 2006. Online learning of approximate dependency parsing algorithms. *Proceedings of the EACL*, pp. 81-88.

Joakim Nivre. 2003. An efficient algorithm for projective dependency parsing. *Proceedings of IWPT*, pp. 149-160.

Joakim Nivre. 2008. Algorithms for deterministic incremental dependency parsing. *Computational Linguistics*, vol. 34, num. 4, pp. 513-553.

Joakim Nivre and Mario Scholz. 2004. Deterministic dependency parsing of English text. *Proceedings of COLING*, pp. 64-70.

Hiroyasu Yamada and Yuji Matsumoto. 2003. Statistical dependency analysis with support vector machines. *Proceedings of IWPT*, pp. 195-206.

Yue Zhang and Stephen Clark. 2008. A tale of two parsers: investigating and combining graph-based and transition-based dependency parsing using beamsearch. *Proceedings of EMNLP*, pp. 562-571.

Sparsity in Dependency Grammar Induction

Jennifer Gillenwater and **Kuzman Ganchev**

University of Pennsylvania
Philadelphia, PA, USA

{jengi,kuzman}@cis.upenn.edu

João Graça

L²F INESC-ID
Lisboa, Portugal

joao.graca@l2f.inesc-id.pt

Fernando Pereira

Google Inc.
Mountain View, CA, USA
pereira@google.com

Ben Taskar

University of Pennsylvania
Philadelphia, PA, USA
taskar@cis.upenn.edu

Abstract

A strong inductive bias is essential in unsupervised grammar induction. We explore a particular sparsity bias in dependency grammars that encourages a small number of unique dependency types. Specifically, we investigate sparsity-inducing penalties on the posterior distributions of parent-child POS tag pairs in the posterior regularization (PR) framework of Graça et al. (2007). In experiments with 12 languages, we achieve substantial gains over the standard expectation maximization (EM) baseline, with average improvement in attachment accuracy of 6.3%. Further, our method outperforms models based on a standard Bayesian sparsity-inducing prior by an average of 4.9%. On English in particular, we show that our approach improves on several other state-of-the-art techniques.

1 Introduction

We investigate an unsupervised learning method for dependency parsing models that imposes sparsity biases on the dependency types. We assume a corpus annotated with POS tags, where the task is to induce a dependency model from the tags for corpus sentences. In this setting, the *type* of a dependency is defined as a pair: tag of the dependent (also known as the child), and tag of the head (also known as the parent). Given that POS tags are designed to convey information about grammatical relations, it is reasonable to assume that only some of the possible dependency types will be realized

for a given language. For instance, in English it is ungrammatical for nouns to dominate verbs, adjectives to dominate adverbs, and determiners to dominate almost any part of speech. Thus, the realized dependency types should be a sparse subset of all possible types.

Previous work in unsupervised grammar induction has tried to achieve sparsity through priors. Liang et al. (2007), Finkel et al. (2007) and Johnson et al. (2007) proposed hierarchical Dirichlet process priors. Cohen et al. (2008) experimented with a discounting Dirichlet prior, which encourages a standard dependency parsing model (see Section 2) to limit the number of dependent types for each head type.

Our experiments show a more effective sparsity pattern is one that limits the total number of unique head-dependent tag pairs. This kind of sparsity bias avoids inducing competition between dependent types for each head type. We can achieve the desired bias with a constraint on model posteriors during learning, using the posterior regularization (PR) framework (Graça et al., 2007). Specifically, to implement PR we augment the maximum marginal likelihood objective of the dependency model with a term that penalizes head-dependent tag distributions that are too permissive.

Although not focused on sparsity, several other studies use soft parameter sharing to couple different types of dependencies. To this end, Cohen et al. (2008) and Cohen and Smith (2009) investigated logistic normal priors, and Headden III et al. (2009) used a backoff scheme. We compare to their results in Section 5.

The remainder of this paper is organized as fol-

lows. Section 2 and 3 review the models and several previous approaches for learning them. Section 4 describes learning with PR. Section 5 describes experiments across 12 languages and Section 6 analyzes the results. For additional details on this work see Gillenwater et al. (2010).

2 Parsing Model

The models we use are based on the generative dependency model with valence (DMV) (Klein and Manning, 2004). For a sentence with tags \mathbf{x} , the root POS $r(\mathbf{x})$ is generated first. Then the model decides whether to generate a right dependent conditioned on the POS of the root and whether other right dependents have already been generated for this head. Upon deciding to generate a right dependent, the POS of the dependent is selected by conditioning on the head POS and the directionality. After stopping on the right, the root generates left dependents using the mirror reversal of this process. Once the root has generated all its dependents, the dependents generate their own dependents in the same manner.

2.1 Model Extensions

For better comparison with previous work we implemented three model extensions, borrowed from Headden III et al. (2009). The first extension alters the stopping probability by conditioning it not only on whether there are *any* dependents in a particular direction already, but also on *how many* such dependents there are. When we talk about models with maximum stop valency $V_s = S$, this means it distinguishes S different cases: $0, 1, \dots, S-2$, and $\geq S-1$ dependents in a given direction. The basic DMV has $V_s = 2$.

The second model extension we implement is analogous to the first, but applies to dependent tag probabilities instead of stop probabilities. Again, we expand the conditioning such that the model considers how many other dependents were already generated in the same direction. When we talk about a model with maximum child valency $V_c = C$, this means we distinguish C different cases. The basic DMV has $V_c = 1$. Since this extension to the dependent probabilities dramatically increases model complexity, the third model extension we implement is to add a backoff for the dependent probabilities that does not condition on the identity of the parent POS (see Equation 2).

More formally, under the extended DMV the

probability of a sentence with POS tags \mathbf{x} and dependency tree \mathbf{y} is given by:

$$p_{\theta}(\mathbf{x}, \mathbf{y}) = p_{root}(r(\mathbf{x})) \times \prod_{y \in \mathbf{y}} p_{stop}(false | y_p, y_d, y_{v_s}) p_{child}(y_c | y_p, y_d, y_{v_c}) \times \prod_{x \in \mathbf{x}} p_{stop}(true | x, left, x_{v_l}) p_{stop}(true | x, right, x_{v_r}) \quad (1)$$

where y is the dependency of y_c on head y_p in direction y_d , and $y_{v_c}, y_{v_s}, x_{v_r}$, and x_{v_l} indicate valence. For the third model extension, the backoff to a probability not dependent on parent POS can be formally expressed as:

$$\lambda p_{child}(y_c | y_p, y_d, y_{v_c}) + (1 - \lambda) p_{child}(y_c | y_d, y_{v_c}) \quad (2)$$

for $\lambda \in [0, 1]$. We fix $\lambda = 1/3$, which is a crude approximation to the value learned by Headden III et al. (2009).

3 Previous Learning Approaches

In our experiments, we compare PR learning to standard expectation maximization (EM) and to Bayesian learning with a sparsity-inducing prior. The EM algorithm optimizes marginal likelihood $\mathcal{L}(\theta) = \log \sum_{\mathbf{Y}} p_{\theta}(\mathbf{X}, \mathbf{Y})$, where $\mathbf{X} = \{\mathbf{x}^1, \dots, \mathbf{x}^n\}$ denotes the entire unlabeled corpus and $\mathbf{Y} = \{\mathbf{y}^1, \dots, \mathbf{y}^n\}$ denotes a set of corresponding parses for each sentence. Neal and Hinton (1998) view EM as block coordinate ascent on a function that lower-bounds $\mathcal{L}(\theta)$. Starting from an initial parameter estimate θ^0 , the algorithm iterates two steps:

$$\mathbf{E} : q^{t+1} = \arg \min_q \mathbf{KL}(q(\mathbf{Y}) \| p_{\theta^t}(\mathbf{Y} | \mathbf{X})) \quad (3)$$

$$\mathbf{M} : \theta^{t+1} = \arg \max_{\theta} \mathbf{E}_{q^{t+1}} [\log p_{\theta}(\mathbf{X}, \mathbf{Y})] \quad (4)$$

Note that the E-step just sets $q^{t+1}(\mathbf{Y}) = p_{\theta^t}(\mathbf{Y} | \mathbf{X})$, since it is an unconstrained minimization of a KL-divergence. The PR method we present modifies the E-step by adding constraints.

Besides EM, we also compare to learning with several Bayesian priors that have been applied to the DMV. One such prior is the Dirichlet, whose hyperparameter we will denote by α . For $\alpha < 0.5$, this prior encourages parameter sparsity. Cohen et al. (2008) use this method with $\alpha = 0.25$ for training the DMV and achieve improvements over basic EM. In this paper we will refer to our own implementation of the Dirichlet prior as the ‘‘discounting Dirichlet’’ (DD) method. In addition to

the Dirichlet, other types of priors have been applied, in particular logistic normal priors (LN) and shared logistic normal priors (SLN) (Cohen et al., 2008; Cohen and Smith, 2009). LN and SLN aim to tie parameters together. Essentially, this has a similar goal to sparsity-inducing methods in that it posits a more concise explanation for the grammar of a language. Headden III et al. (2009) also implement a sort of parameter tying for the E-DMV through a learning a backoff distribution on child probabilities. We compare against results from all these methods.

4 Learning with Sparse Posteriors

We would like to penalize models that predict a large number of distinct dependency types. To enforce this penalty, we use the posterior regularization (PR) framework (Graça et al., 2007). PR is closely related to generalized expectation constraints (Mann and McCallum, 2007; Mann and McCallum, 2008; Bellare et al., 2009), and is also indirectly related to a Bayesian view of learning with constraints on posteriors (Liang et al., 2009). The PR framework uses constraints on posterior expectations to guide parameter estimation. Here, PR allows a natural and tractable representation of sparsity constraints based on edge type counts that cannot easily be encoded in model parameters. We use a version of PR where the desired bias is a penalty on the log likelihood (see Ganchev et al. (2010) for more details). For a distribution p_θ , we define a penalty as the (generic) β -norm of expectations of some features ϕ :

$$\|\mathbf{E}_{p_\theta}[\phi(\mathbf{X}, \mathbf{Y})]\|_\beta \quad (5)$$

For computational tractability, rather than penalizing the model’s posteriors directly, we use an auxiliary distribution q , and penalize the marginal log-likelihood of a model by the KL-divergence of p_θ from q , plus the penalty term with respect to q . For a fixed set of model parameters θ the full PR penalty term is:

$$\min_q \text{KL}(q(\mathbf{Y}) \parallel p_\theta(\mathbf{Y}|\mathbf{X})) + \sigma \|\mathbf{E}_q[\phi(\mathbf{X}, \mathbf{Y})]\|_\beta \quad (6)$$

where σ is the strength of the regularization. PR seeks to maximize $\mathcal{L}(\theta)$ minus this penalty term. The resulting objective can be optimized by a variant of the EM (Dempster et al., 1977) algorithm used to optimize $\mathcal{L}(\theta)$.

4.1 ℓ_1/ℓ_∞ Regularization

We now define precisely how to count dependency types. For each child tag c , let i range over an enumeration of all occurrences of c in the corpus, and let p be another tag. Let the indicator $\phi_{cpi}(\mathbf{X}, \mathbf{Y})$ have value 1 if p is the parent tag of the i th occurrence of c , and value 0 otherwise. The number of unique dependency types is then:

$$\sum_{cp} \max_i \phi_{cpi}(\mathbf{X}, \mathbf{Y}) \quad (7)$$

Note there is an asymmetry in this count: occurrences of child type c are enumerated with i , but all occurrences of parent type p are or-ed in ϕ_{cpi} . That is, $\phi_{cpi} = 1$ if *any* occurrence of p is the parent of the i th occurrence of c . We will refer to PR training with this constraint as PR-AS. Instead of counting pairs of a child token and a parent type, we can alternatively count pairs of a child token and a parent token by letting p range over all *tokens* rather than *types*. Then each potential dependency corresponds to a different indicator ϕ_{cpij} , and the penalty is symmetric with respect to parents and children. We will refer to PR training with this constraint as PR-S. Both approaches perform very well, so we report results for both.

Equation 7 can be viewed as a mixed-norm penalty on the features ϕ_{cpi} or ϕ_{cpij} : the sum corresponds to an ℓ_1 norm and the max to an ℓ_∞ norm. Thus, the quantity we want to minimize fits precisely into the PR penalty framework. Formally, to optimize the PR objective, we complete the following E-step:

$$\arg \min_q \text{KL}(q(\mathbf{Y}) \parallel p_\theta(\mathbf{Y}|\mathbf{X})) + \sigma \sum_{cp} \max_i \mathbf{E}_q[\phi(\mathbf{X}, \mathbf{Y})], \quad (8)$$

which can equivalently be written as:

$$\begin{aligned} \min_{q(\mathbf{Y}), \xi_{cp}} \quad & \text{KL}(q(\mathbf{Y}) \parallel p_\theta(\mathbf{Y}|\mathbf{X})) + \sigma \sum_{cp} \xi_{cp} \\ \text{s. t.} \quad & \xi_{cp} \leq \mathbf{E}_q[\phi(\mathbf{X}, \mathbf{Y})] \end{aligned} \quad (9)$$

where ξ_{cp} corresponds to the maximum expectation of ϕ over all instances of c and p . Note that the projection problem can be solved efficiently in the dual (Ganchev et al., 2010).

5 Experiments

We evaluate on 12 languages. Following the example of Smith and Eisner (2006), we strip punctuation from the sentences and keep only sentences of length ≤ 10 . For simplicity, for all models we use the “harmonic” initializer from Klein

Model	EM	PR	Type	σ
DMV	45.8	62.1	PR-S	140
2-1	45.1	62.7	PR-S	100
2-2	54.4	62.9	PR-S	80
3-3	55.3	64.3	PR-S	140
4-4	55.1	64.4	PR-AS	140

Table 1: Attachment accuracy results. **Column 1:** V_c - V_s used for the E-DMV models. **Column 3:** Best PR result for each model, which is chosen by applying each of the two types of constraints (PR-S and PR-AS) and trying $\sigma \in \{80, 100, 120, 140, 160, 180\}$. **Columns 4 & 5:** Constraint type and σ that produced the values in column 3.

and Manning (2004), which we refer to as K&M. We always train for 100 iterations and evaluate on the test set using Viterbi parses. Before evaluating, we smooth the resulting models by adding e^{-10} to each learned parameter, merely to remove the chance of zero probabilities for unseen events. (We did not tune this as it should make very little difference for final parses.) We score models by their attachment accuracy — the fraction of words assigned the correct parent.

5.1 Results on English

We start by comparing English performance for EM, PR, and DD. To find α for DD we searched over five values: $\{0.01, 0.1, 0.25, 1\}$. We found 0.25 to be the best setting for the DMV, the same as found by Cohen et al. (2008). DD achieves accuracy 46.4% with this α . For the E-DMV we tested four model complexities with valencies V_c - V_s of 2-1, 2-2, 3-3, and 4-4. DD’s best accuracy was 53.6% with the 4-4 model at $\alpha = 0.1$. A comparison between EM and PR is shown in Table 1. PR-S generally performs better than the PR-AS for English. Comparing PR-S to EM, we also found PR-S is always better, independent of the particular σ , with improvements ranging from 2% to 17%. Note that in this work we do not perform the PR projection at test time; we found it detrimental, probably due to a need to set the (corpus-size-dependent) σ differently for the test set. We also note that development likelihood and the best setting for σ are not well-correlated, which unfortunately makes it hard to pick these parameters without some supervision.

5.2 Comparison with Previous Work

In this section we compare to previously published unsupervised dependency parsing results for English. It might be argued that the comparison is unfair since we do supervised selection of model

Learning Method	Accuracy		
	≤ 10	≤ 20	all
PR-S ($\sigma = 140$)	62.1	53.8	49.1
LN families	59.3	45.1	39.0
SLN TieV & N	61.3	47.4	41.4
PR-AS ($\sigma = 140$)	64.4	55.2	50.5
DD ($\alpha = 1, \lambda$ learned)	65.0 (± 5.7)		

Table 2: Comparison with previous published results. Rows 2 and 3 are taken from Cohen et al. (2008) and Cohen and Smith (2009), and row 5 from Headen III et al. (2009).

complexity and regularization strength. However, we feel the comparison is not so unfair as we perform only a very limited search of the model- σ space. Specifically, the only values of σ we search over are $\{80, 100, 120, 140, 160, 180\}$.

First, we consider the top three entries in Table 2, which are for the basic DMV. The first entry was generated using our implementation of PR-S. The second two entries are logistic normal and shared logistic normal parameter tying results (Cohen et al., 2008; Cohen and Smith, 2009). The PR-S result is the clear winner, especially as length of test sentences increases. For the bottom two entries in the table, which are for the E-DMV, the last entry is best, corresponding to using a DD prior with $\alpha = 1$ (non-sparsifying), but with a special “random pools” initialization and a learned weight λ for the child backoff probability. The result for PR-AS is well within the variance range of this last entry, and thus we conjecture that combining PR-AS with random pools initialization and learned λ would likely produce the best-performing model of all.

5.3 Results on Other Languages

Here we describe experiments on 11 additional languages. For each we set σ and model complexity (DMV versus one of the four E-DMV experimented with previously) based on the best configuration found for English. This likely will not result in the ideal parameters for all languages, but provides a realistic test setting: a user has available a labeled corpus in one language, and would like to induce grammars for many other languages. Table 3 shows the performance for all models and training procedures. We see that the sparsifying methods tend to improve over EM most of the time. For the basic DMV, average improvements are 1.6% for DD, 6.0% for PR-S, and 7.5% for PR-AS. PR-AS beats PR-S in 8 out of 12 cases,

	Bg	Cz	De	Dk	En	Es	Jp	Nl	Pt	Se	Si	Tr
	DMV Model											
EM	37.8	29.6	35.7	47.2	45.8	40.3	52.8	37.1	35.7	39.4	42.3	46.8
DD 0.25	39.3	30.0	38.6	43.1	46.4	47.5	57.8	35.1	38.7	40.2	48.8	43.8
PR-S 140	53.7	31.5	39.6	44.0	62.1	61.1	58.8	31.0	47.0	42.2	39.9	51.4
PR-AS 140	54.0	32.0	39.6	42.4	61.9	62.4	60.2	37.9	47.8	38.7	50.3	53.4
	Extended Model											
EM (3,3)	41.7	48.9	40.1	46.4	55.3	44.3	48.5	47.5	35.9	48.6	47.5	46.2
DD 0.1 (4,4)	47.6	48.5	42.0	44.4	53.6	48.9	57.6	45.2	48.3	47.6	35.6	48.9
PR-S 140 (3,3)	59.0	54.7	47.4	45.8	64.3	57.9	60.8	33.9	54.3	45.6	49.1	56.3
PR-AS 140 (4,4)	59.8	54.6	45.7	46.6	64.4	57.9	59.4	38.8	49.5	41.4	51.2	56.9

Table 3: Attachment accuracy results. The parameters used are the best settings found for English. Values for hyperparameters (α or σ) are given after the method name. For the extended model (V_c, V_s) are indicated in parentheses. En is the English Penn Treebank (Marcus et al., 1993) and the other 11 languages are from the CoNLL X shared task: Bulgarian [Bg] (Simov et al., 2002), Czech [Cz] (Bohomovà et al., 2001), German [De] (Brants et al., 2002), Danish [Dk] (Kromann et al., 2003), Spanish [Es] (Civit and Martí, 2004), Japanese [Jp] (Kawata and Bartels, 2000), Dutch [Nl] (Van der Beek et al., 2002), Portuguese [Pt] (Afonso et al., 2002), Swedish [Se] (Nilsson et al., 2005), Slovene [Sl] (Džeroski et al., 2006), and Turkish [Tr] (Ofłazer et al., 2003).

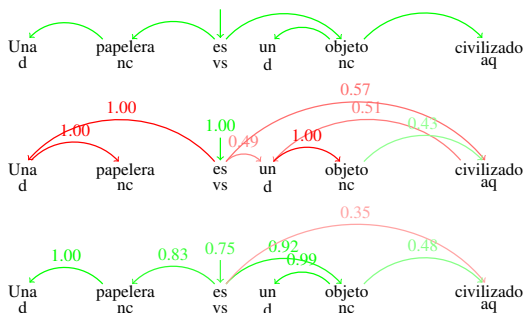


Figure 1: Posterior edge probabilities for an example sentence from the Spanish test corpus. At the top are the gold dependencies, the middle are EM posteriors, and bottom are PR posteriors. Green indicates correct dependencies and red indicates incorrect dependencies. The numbers on the edges are the values of the posterior probabilities.

though the average increase is only 1.5%. PR-S is also better than DD for 10 out of 12 languages. If we instead consider these methods for the E-DMV, DD performs worse, just 1.4% better than the E-DMV EM, while both PR-S and PR-AS continue to show substantial average improvements over EM, 6.5% and 6.3%, respectively.

6 Analysis

One common EM error that PR fixes in many languages is the directionality of the noun-determiner relation. Figure 1 shows an example of a Spanish sentence where PR significantly outperforms EM because of this. Sentences such as “Lleva tiempo entenderlos” which has tags “main-verb common-noun main-verb” (no determiner tag) provide an explanation for PR’s improvement—when PR sees that sometimes nouns can appear without determiners but that the opposite situation

does not occur, it shifts the model parameters to make nouns the parent of determiners instead of the reverse. Then it does not have to pay the cost of assigning a parent with a new tag to cover each noun that doesn’t come with a determiner.

7 Conclusion

In this paper we presented a new method for unsupervised learning of dependency parsers. In contrast to previous approaches that constrain model parameters, we constrain model posteriors. Our approach consistently outperforms the standard EM algorithm and a discounting Dirichlet prior.

We have several ideas for further improving our constraints, such as: taking into account the directionality of the edges, using different regularization strengths for the root probabilities than for the child probabilities, and working directly on word types rather than on POS tags. In the future, we would also like to try applying similar constraints to the more complex task of joint induction of POS tags and dependency parses.

Acknowledgments

J. Gillenwater was supported by NSF-IGERT 0504487. K. Ganchev was supported by ARO MURI SUBTLE W911NF-07-1-0216. J. Graça was supported by FCT fellowship SFRH/BD/27528/2006 and by FCT project CMU-PT/HuMach/0039/2008. B. Taskar was partly supported by DARPA CSSG and ONR Young Investigator Award N000141010746.

References

- S. Afonso, E. Bick, R. Haber, and D. Santos. 2002. Floresta Sinta(c)tica: a treebank for Portuguese. In *Proc. LREC*.
- K. Bellare, G. Druck, and A. McCallum. 2009. Alternating projections for learning with expectation constraints. In *Proc. UAI*.
- A. Bohomová, J. Hajic, E. Hajicova, and B. Hladka. 2001. The prague dependency treebank: Three-level annotation scenario. In Anne Abeillé, editor, *Treebanks: Building and Using Syntactically Annotated Corpora*.
- S. Brants, S. Dipper, S. Hansen, W. Lezius, and G. Smith. 2002. The TIGER treebank. In *Proc. Workshop on Treebanks and Linguistic Theories*.
- M. Civit and M.A. Martí. 2004. Building cast3lb: A Spanish Treebank. *Research on Language & Computation*.
- S.B. Cohen and N.A. Smith. 2009. The shared logistic normal distribution for grammar induction. In *Proc. NAACL*.
- S.B. Cohen, K. Gimpel, and N.A. Smith. 2008. Logistic normal priors for unsupervised probabilistic grammar induction. In *Proc. NIPS*.
- A.P. Dempster, N.M. Laird, and D.B. Rubin. 1977. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society*, 39(1):1–38.
- S. Džeroski, T. Erjavec, N. Ledinek, P. Pajas, Z. Žabokrtsky, and A. Žele. 2006. Towards a Slovene dependency treebank. In *Proc. LREC*.
- J. Finkel, T. Grenager, and C. Manning. 2007. The infinite tree. In *Proc. ACL*.
- K. Ganchev, J. Graça, J. Gillenwater, and B. Taskar. 2010. Posterior regularization for structured latent variable models. *Journal of Machine Learning Research*.
- J. Gillenwater, K. Ganchev, J. Graça, F. Pereira, and B. Taskar. 2010. Posterior sparsity in unsupervised dependency parsing. Technical report, MS-CIS-10-19, University of Pennsylvania.
- J. Graça, K. Ganchev, and B. Taskar. 2007. Expectation maximization and posterior constraints. In *Proc. NIPS*.
- W.P. Headen III, M. Johnson, and D. McClosky. 2009. Improving unsupervised dependency parsing with richer contexts and smoothing. In *Proc. NAACL*.
- M. Johnson, T.L. Griffiths, and S. Goldwater. 2007. Adaptor grammars: A framework for specifying compositional nonparametric Bayesian models. In *Proc. NIPS*.
- Y. Kawata and J. Bartels. 2000. Stylebook for the Japanese Treebank in VERBMOBIL. Technical report, Eberhard-Karls-Universität Tübingen.
- D. Klein and C. Manning. 2004. Corpus-based induction of syntactic structure: Models of dependency and constituency. In *Proc. ACL*.
- M.T. Kromann, L. Mikkelsen, and S.K. Lynge. 2003. Danish Dependency Treebank. In *Proc. TLT*.
- P. Liang, S. Petrov, M.I. Jordan, and D. Klein. 2007. The infinite PCFG using hierarchical Dirichlet processes. In *Proc. EMNLP*.
- P. Liang, M.I. Jordan, and D. Klein. 2009. Learning from measurements in exponential families. In *Proc. ICML*.
- G. Mann and A. McCallum. 2007. Simple, robust, scalable semi-supervised learning via expectation regularization. In *Proc. ICML*.
- G. Mann and A. McCallum. 2008. Generalized expectation criteria for semi-supervised learning of conditional random fields. In *Proc. ACL*.
- M. Marcus, M. Marcinkiewicz, and B. Santorini. 1993. Building a large annotated corpus of English: The Penn Treebank. *Computational Linguistics*, 19(2):313–330.
- R. Neal and G. Hinton. 1998. A new view of the EM algorithm that justifies incremental, sparse and other variants. In M. I. Jordan, editor, *Learning in Graphical Models*, pages 355–368. MIT Press.
- J. Nilsson, J. Hall, and J. Nivre. 2005. MAMBA meets TIGER: Reconstructing a Swedish treebank from antiquity. *NODALIDA Special Session on Treebanks*.
- K. Oflazer, B. Say, D.Z. Hakkani-Tür, and G. Tür. 2003. Building a Turkish treebank. *Treebanks: Building and Using Parsed Corpora*.
- K. Simov, P. Osenova, M. Slavcheva, S. Kolkovska, E. Balabanova, D. Doikoff, K. Ivanova, A. Simov, E. Simov, and M. Kouylekov. 2002. Building a linguistically interpreted corpus of bulgarian: the bul-treebank. In *Proc. LREC*.
- N. Smith and J. Eisner. 2006. Annealing structural bias in multilingual weighted grammar induction. In *Proc. ACL*.
- L. Van der Beek, G. Bouma, R. Malouf, and G. Van Noord. 2002. The Alpino dependency treebank. *Language and Computers*.

Top-Down K-Best A* Parsing

Adam Pauls and Dan Klein

Computer Science Division
University of California at Berkeley
{adpauls, klein}@cs.berkeley.edu

Chris Quirk

Microsoft Research
Redmond, WA, 98052
chrisq@microsoft.com

Abstract

We propose a top-down algorithm for extracting k -best lists from a parser. Our algorithm, TKA* is a variant of the k -best A* (KA*) algorithm of Pauls and Klein (2009). In contrast to KA*, which performs an inside and outside pass before performing k -best extraction bottom up, TKA* performs only the inside pass before extracting k -best lists top down. TKA* maintains the same optimality and efficiency guarantees of KA*, but is simpler to both specify and implement.

1 Introduction

Many situations call for a parser to return a k -best list of parses instead of a single best hypothesis.¹ Currently, there are two efficient approaches known in the literature. The k -best algorithm of Jiménez and Marzal (2000) and Huang and Chiang (2005), referred to hereafter as LAZY, operates by first performing an exhaustive Viterbi inside pass and then lazily extracting k -best lists in top-down manner. The k -best A* algorithm of Pauls and Klein (2009), hereafter KA*, computes Viterbi inside *and* outside scores before extracting k -best lists bottom up.

Because these additional passes are only partial, KA* can be significantly faster than LAZY, especially when a heuristic is used (Pauls and Klein, 2009). In this paper, we propose TKA*, a top-down variant of KA* that, like LAZY, performs only an inside pass before extracting k -best lists top-down, but maintains the same optimality and efficiency guarantees as KA*. This algorithm can be seen as a generalization of the lattice k -best algorithm of Soong and Huang (1991) to parsing. Because TKA* eliminates the outside pass from KA*, TKA* is simpler both in implementation and specification.

¹See Huang and Chiang (2005) for a review.

2 Review

Because our algorithm is very similar to KA*, which is in turn an extension of the (1-best) A* parsing algorithm of Klein and Manning (2003), we first introduce notation and review those two algorithms before presenting our new algorithm.

2.1 Notation

Assume we have a PCFG² \mathcal{G} and an input sentence $s_0 \dots s_{n-1}$ of length n . The grammar \mathcal{G} has a set of symbols denoted by capital letters, including a distinguished goal (root) symbol G . Without loss of generality, we assume Chomsky normal form: each non-terminal rule r in \mathcal{G} has the form $r = A \rightarrow B C$ with weight w_r . Edges are labeled spans $e = (A, i, j)$. *Inside derivations* of an edge (A, i, j) are trees with root non-terminal A , spanning $s_i \dots s_{j-1}$. The weight (negative log-probability) of the best (minimum) inside derivation for an edge e is called the Viterbi inside score $\beta(e)$, and the weight of the best derivation of $G \rightarrow s_0 \dots s_{i-1} A s_j \dots s_{n-1}$ is called the Viterbi outside score $\alpha(e)$. The goal of a k -best parsing algorithm is to compute the k best (minimum weight) inside derivations of the edge $(G, 0, n)$.

We formulate the algorithms in this paper in terms of prioritized weighted deduction rules (Shieber et al., 1995; Nederhof, 2003). A *prioritized weighted deduction rule* has the form

$$\phi_1 : w_1, \dots, \phi_n : w_n \xrightarrow{p(w_1, \dots, w_n)} \phi_0 : g(w_1, \dots, w_n)$$

where ϕ_1, \dots, ϕ_n are the *antecedent items* of the deduction rule and ϕ_0 is the *conclusion item*. A deduction rule states that, given the antecedents ϕ_1, \dots, ϕ_n with weights w_1, \dots, w_n , the conclusion ϕ_0 can be formed with weight $g(w_1, \dots, w_n)$ and priority $p(w_1, \dots, w_n)$.

²While we present the algorithm specialized to parsing with a PCFG, this algorithm generalizes to a wide range of

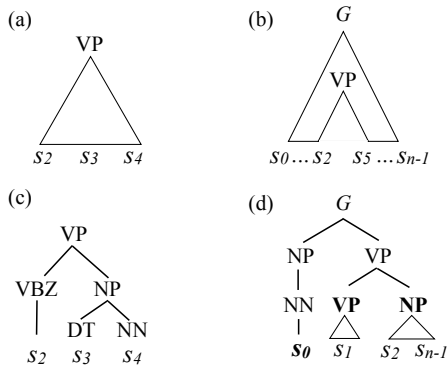


Figure 1: Representations of the different types of items used in parsing. (a) An inside edge item $I(\text{VP}, 2, 5)$. (b) An outside edge item $O(\text{VP}, 2, 5)$. (c) An inside derivation item: $D(T^{\text{VP}}, 2, 5)$. (d) An outside derivation item: $Q(T_{\text{VP}}^G, 1, 2, \{(NP, 2, n)\})$. The edges in boldface are frontier edges.

These deduction rules are “executed” within a generic agenda-driven algorithm, which constructs items in a prioritized fashion. The algorithm maintains an *agenda* (a priority queue of items), as well as a *chart* of items already processed. The fundamental operation of the algorithm is to pop the highest priority item ϕ from the agenda, put it into the chart with its current weight, and apply deduction rules to form any items which can be built by combining ϕ with items already in the chart. When the resulting items are either new or have a weight smaller than an item’s best score so far, they are put on the agenda with priority given by $p(\cdot)$. Because all antecedents must be constructed before a deduction rule is executed, we sometimes refer to particular conclusion item as “waiting” on another item before it can be built.

2.2 A*

A* parsing (Klein and Manning, 2003) is an algorithm for computing the 1-best parse of a sentence. A* operates on items called *inside edge items* $I(A, i, j)$, which represent the many possible inside derivations of an edge (A, i, j) . Inside edge items are constructed according to the IN deduction rule of Table 1. This deduction rule constructs inside edge items in a bottom-up fashion, combining items representing smaller edges $I(B, i, k)$ and $I(C, k, j)$ with a grammar rule $r = A \rightarrow B C$ to form a larger item $I(A, i, j)$. The weight of a newly constructed item is given by the sum of the weights of the antecedent items and the grammar rule r , and its priority is given by

hypergraph search problems as shown in Klein and Manning (2001).

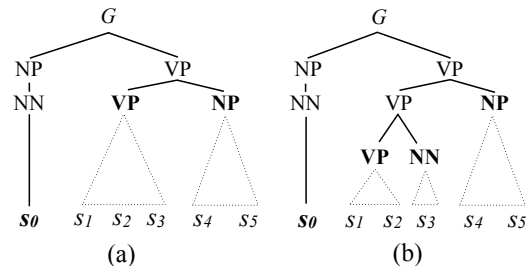


Figure 2: (a) An outside derivation item before expansion at the edge $(\text{VP}, 1, 4)$. (b) A possible expansion of the item in (a) using the rule $\text{VP} \rightarrow \text{VP NN}$. Frontier edges are marked in boldface.

its weight plus a *heuristic* $h(A, i, j)$. For consistent and admissible heuristics $h(\cdot)$, this deduction rule guarantees that when an inside edge item is removed from the agenda, its current weight is its true Viterbi inside score.

The heuristic h controls the speed of the algorithm. It can be shown that an edge e satisfying $\beta(e) + h(A, i, j) > \beta(G, 0, n)$ will never be removed from the agenda, allowing some edges to be safely pruned during parsing. The more closely $h(e)$ approximates the Viterbi outside cost $\alpha(e)$, the more items are pruned.

2.3 KA*

The use of inside edge items in A* exploits the optimal substructure property of derivations – since a best derivation of a larger edge is always composed of best derivations of smaller edges, it is only necessary to compute the best way of building a particular inside edge item. When finding k -best lists, this is no longer possible, since we are interested in suboptimal derivations.

Thus, KA^* , the k -best extension of A*, must search not in the space of inside edge items, but rather in the space of *inside derivation items* $D(T^A, i, j)$, which represent specific derivations of the edge (A, i, j) using tree T^A . However, the number of inside derivation items is exponential in the length of the input sentence, and even with a very accurate heuristic, running A* directly in this space is not feasible.

Fortunately, Pauls and Klein (2009) show that with a *perfect* heuristic, that is, $h(e) = \alpha(e) \forall e$, A* search on inside derivation items will only remove items from the agenda that participate in the true k -best lists (up to ties). In order to compute this perfect heuristic, KA^* makes use of *outside edge items* $O(A, i, j)$ which represent the many possible derivations of $G \rightarrow$

IN^{*†}:		$I(B, i, l) : w_1$	$I(C, l, j) : w_2$	$\xrightarrow{w_1+w_2+w_r+h(A,i,j)}$	$I(A, i, j) : w_1 + w_2 + w_r$
IN-D[†]:	$O(A, i, j) : w_1$	$D(T^B, i, l) : w_2$	$D(T^C, l, j) : w_3$	$\xrightarrow{w_2+w_3+w_r+w_1}$	$D(T^A, i, j) : w_2 + w_3 + w_r$
OUT-L[†]:	$O(A, i, j) : w_1$	$I(B, i, l) : w_2$	$I(C, l, j) : w_3$	$\xrightarrow{w_1+w_3+w_r+w_2}$	$O(B, i, l) : w_1 + w_3 + w_r$
OUT-R[†]:	$O(A, i, j) : w_1$	$I(B, i, l) : w_2$	$I(C, l, j) : w_3$	$\xrightarrow{w_1+w_2+w_r+w_3}$	$O(C, l, j) : w_1 + w_2 + w_r$
OUT-D[*]:	$Q(T_A^G, i, j, \mathcal{F}) : w_1$	$I(B, i, l) : w_2$	$I(C, l, j) : w_3$	$\xrightarrow{w_1+w_r+w_2+w_3+\beta(\mathcal{F})}$	$Q(T_B^G, i, l, \mathcal{F}_C) : w_1 + w_r$

Table 1: The deduction rules used in this paper. Here, r is the rule $A \rightarrow B C$. A superscript $*$ indicates that the rule is used in TKA^{*}, and a superscript \dagger indicates that the rule is used in KA^{*}. In IN-D, the tree T_A is rooted at (A, i, j) and has children T^B and T^C . In OUT-D, the tree T_B^G is the tree T_A^G extended at (A, i, j) with rule r , \mathcal{F}_C is the list \mathcal{F} with (C, l, j) prepended, and $\beta(\mathcal{F})$ is $\sum_{e \in \mathcal{F}} \beta(e)$. Whenever the left child $I(B, i, l)$ of an application of OUT-D represents a terminal, the next edge is removed from \mathcal{F} and is used as the new point of expansion.

$s_1 \dots s_i A s_{j+1} \dots s_n$ (see Figure 1(b)).

Outside items are built using the OUT-L and OUT-R deduction rules shown in Table 1. OUT-L and OUT-R combine, in a top-down fashion, an outside edge over a larger span and inside edge over a smaller span to form a new outside edge over a smaller span. Because these rules make reference to inside edge items $I(A, i, j)$, these items must also be built using the IN deduction rules from 1-best A^{*}. Outside edge items must thus wait until the necessary inside edge items have been built. The outside pass is initialized with the item $O(G, 0, n)$ when the inside edge item $I(G, 0, n)$ is popped from the agenda.

Once we have started populating outside scores using the outside deductions, we can initiate a search on inside derivation items.³ These items are built bottom-up using the IN-D deduction rule. The crucial element of this rule is that derivation items for a particular edge wait until the exact outside score of that edge has been computed. The algorithm terminates when k derivation items rooted at $(G, 0, n)$ have been popped from the agenda.

3 TKA^{*}

KA^{*} efficiently explores the space of inside derivation items because it waits for the exact Viterbi outside cost before building each derivation item. However, these outside costs and associated deduction items are only auxiliary quantities used to guide the exploration of inside derivations: they allow KA^{*} to prioritize currently constructed inside derivation items (i.e., constructed derivations of the goal) by their optimal completion costs. Outside costs are thus only necessary because we construct partial derivations bottom-up; if we constructed partial derivations in a top-down fashion, all we would need to compute opti-

³We stress that the order of computation is entirely specified by the deduction rules – we only speak about e.g. “initiating a search” as an appeal to intuition.

mal completion costs are Viterbi inside scores, and we could forget the outside pass.

TKA^{*} does exactly that. Inside edge items are constructed in the same way as KA^{*}, but once the inside edge item $I(G, 0, n)$ has been discovered, TKA^{*} begins building partial derivations from the goal outwards. We replace the inside derivation items of KA^{*} with *outside derivation items*, which represent trees rooted at the goal and expanding downwards. These items bottom out in a list of edges called the *frontier* edges. See Figure 1(d) for a graphical representation. When a frontier edge represents a single word in the input, i.e. is of the form $(s_i, i, i + 1)$, we say that edge is *complete*. An outside derivation can be expanded by applying a rule to one of its incomplete frontier edges; see Figure 2. In the same way that inside derivation items wait on exact outside scores before being built, outside derivation items wait on the inside edge items of all frontier edges before they can be constructed.

Although building derivations top-down obviates the need for a 1-best outside pass, it raises a new issue. When building derivations bottom-up, the only way to expand a particular partial inside derivation is to combine it with another partial inside derivation to build a bigger tree. In contrast, an outside derivation item can be expanded anywhere along its frontier. Naively building derivations top-down would lead to a prohibitively large number of expansion choices.

We solve this issue by always expanding the left-most incomplete frontier edge of an outside derivation item. We show the deduction rule OUT-D which performs this deduction in Figure 1(d). We denote an outside derivation item as $Q(T_A^G, i, j, \mathcal{F})$, where T_A^G is a tree rooted at the goal with left-most incomplete edge (A, i, j) , and \mathcal{F} is the list of incomplete frontier edges excluding (A, i, j) , ordered from left to right. Whenever the application of this rule “completes” the left-

most edge, the next edge is removed from \mathcal{F} and is used as the new point of expansion. Once all frontier edges are complete, the item represents a correctly scored derivation of the goal, explored in a pre-order traversal.

3.1 Correctness

It should be clear that expanding the left-most incomplete frontier edge first eventually explores the same set of derivations as expanding all frontier edges simultaneously. The only worry in fixing this canonical order is that we will somehow explore the Q items in an incorrect order, possibly building some complete derivation Q'_C before a more optimal complete derivation Q_C . However, note that all items Q along the left-most construction of Q_C have priority equal to or better than any less optimal complete derivation Q'_C . Therefore, when Q'_C is enqueued, it will have lower priority than all Q ; Q'_C will therefore not be dequeued until all Q – and hence Q_C – have been built.

Furthermore, it can be shown that the top-down expansion strategy maintains the same efficiency and optimality guarantees as KA^* for all item types: for consistent heuristics h , the first k entirely complete outside derivation items are the true k -best derivations (modulo ties), and that only derivation items which participate in those k -best derivations will be removed from the queue (up to ties).

3.2 Implementation Details

Building derivations bottom-up is convenient from an indexing point of view: since larger derivations are built from smaller ones, it is not necessary to construct the larger derivation from scratch. Instead, one can simply construct a new tree whose children point to the old trees, saving both memory and CPU time.

In order keep the same efficiency when building trees top-down, a slightly different data structure is necessary. We represent top-down derivations as a lazy list of expansions. The top node T_G^G is an empty list, and whenever we expand an outside derivation item $Q(T_A^G, i, j, \mathcal{F})$ with a rule $r = A \rightarrow B C$ and split point l , the resulting derivation T_B^G is a new list item with (r, l) as the head data, and T_A^G as its tail. The tree can be reconstructed later by recursively reconstructing the parent, and adding the edges (B, i, l) and (C, l, j) as children of (A, i, j) .

3.3 Advantages

Although our algorithm eliminates the 1-best outside pass of KA^* , in practice, even for $k = 10^4$, the 1-best inside pass remains the overwhelming bottleneck (Pauls and Klein, 2009), and our modifications leave that pass unchanged.

However, we argue that our implementation is simpler to specify and implement. In terms of deduction rules, our algorithm eliminates the 2 outside deduction rules and replaces the IN-D rule with the OUT-D rule, bringing the total number of rules from four to two.

The ease of specification translates directly into ease of implementation. In particular, if high-quality heuristics are not available, it is often more efficient to implement the 1-best inside pass as an exhaustive dynamic program, as in Huang and Chiang (2005). In this case, one would only need to implement a single, agenda-based k -best extraction phase, instead of the 2 needed for KA^* .

3.4 Performance

The contribution of this paper is theoretical, not empirical. We have argued that TKA^* is simpler than KA^* , but we do not expect it to do any more or less work than KA^* , modulo grammar specific optimizations. Therefore, we simply verify, like KA^* , that the additional work of extracting k -best lists with TKA^* is negligible compared to the time spent building 1-best inside edges.

We examined the time spent building 100-best lists for the same experimental setup as Pauls and Klein (2009).⁴ On 100 sentences, our implementation of TKA^* constructed 3.46 billion items, of which about 2% were outside derivation items. Our implementation of KA^* constructed 3.41 billion edges, of which about 0.1% were outside edge items or inside derivation items. In other words, the cost of k -best extraction is dwarfed by the the 1-best inside edge computation in both cases. The reason for the slight performance advantage of KA^* is that our implementation of KA^* uses lazy optimizations discussed in Pauls and Klein (2009), and while such optimizations could easily be incorporated in TKA^* , we have not yet done so in our implementation.

⁴This setup used 3- and 6-round state-split grammars from Petrov et al. (2006), the former used to compute a heuristic for the latter, tested on sentences of length up to 25.

4 Conclusion

We have presented TKA*, a simplification to the KA* algorithm. Our algorithm collapses the 1-best outside and bottom-up derivation passes of KA* into a single, top-down pass without sacrificing efficiency or optimality. This reduces the number of non base-case deduction rules, making TKA* easier both to specify and implement.

Acknowledgements

This project is funded in part by the NSF under grant 0643742 and an NSERC Postgraduate Fellowship.

References

- Liang Huang and David Chiang. 2005. Better k-best parsing. In *Proceedings of the International Workshop on Parsing Technologies (IWPT)*, pages 53–64.
- Víctor M. Jiménez and Andrés Marzal. 2000. Computation of the n best parse trees for weighted and stochastic context-free grammars. In *Proceedings of the Joint IAPR International Workshops on Advances in Pattern Recognition*, pages 183–192, London, UK. Springer-Verlag.
- Dan Klein and Christopher D. Manning. 2001. Parsing and hypergraphs. In *Proceedings of the International Workshop on Parsing Technologies (IWPT)*, pages 123–134.
- Dan Klein and Christopher D. Manning. 2003. A* parsing: Fast exact Viterbi parse selection. In *Proceedings of the Human Language Technology Conference and the North American Association for Computational Linguistics (HLT-NAACL)*, pages 119–126.
- Mark-Jan Nederhof. 2003. Weighted deductive parsing and Knuth’s algorithm. *Computational Linguistics*, 29(1):135–143.
- Adam Pauls and Dan Klein. 2009. K-best A* parsing. In *Proceedings of the Association for Computational Linguistics (ACL)*.
- Slav Petrov, Leon Barrett, Romain Thibaux, and Dan Klein. 2006. Learning accurate, compact, and interpretable tree annotation. In *Proceedings of the Association for Computational Linguistics (ACL)*.
- Stuart M. Shieber, Yves Schabes, and Fernando C. N. Pereira. 1995. Principles and implementation of deductive parsing. *Journal of Logic Programming*, 24:3–36.
- Frank K. Soong and Eng-Fong Huang. 1991. A tree-trellis based fast search for finding the n best sentence hypotheses in continuous speech recognition. In *Proceedings of the Workshop on Speech and Natural Language*.

Simple semi-supervised training of part-of-speech taggers

Anders Søgaard

Center for Language Technology
University of Copenhagen
soegaard@hum.ku.dk

Abstract

Most attempts to train part-of-speech taggers on a mixture of labeled and unlabeled data have failed. In this work stacked learning is used to reduce tagging to a classification task. This simplifies semi-supervised training considerably. Our preferred semi-supervised method combines tri-training (Li and Zhou, 2005) and disagreement-based co-training. On the Wall Street Journal, we obtain an error reduction of 4.2% with SVMTool (Gimenez and Marquez, 2004).

1 Introduction

Semi-supervised part-of-speech (POS) tagging is relatively rare, and the main reason seems to be that results have mostly been negative. Meritaldo (1994), in a now famous negative result, attempted to improve HMM POS tagging by expectation maximization with unlabeled data. Clark et al. (2003) reported positive results with little labeled training data but negative results when the amount of labeled training data increased; the same seems to be the case in Wang et al. (2007) who use co-training of two diverse POS taggers. Huang et al. (2009) present positive results for self-training a simple bigram POS tagger, but results are considerably below state-of-the-art.

Recently researchers have explored alternative methods. Suzuki and Isozaki (2008) introduce a semi-supervised extension of conditional random fields that combines supervised and unsupervised probability models by so-called MDF parameter estimation, which reduces error on Wall Street Journal (WSJ) standard splits by about 7% relative to their supervised baseline. Spoustova et al. (2009) use a new pool of unlabeled data tagged by an ensemble of state-of-the-art taggers in every training step of an averaged perceptron

POS tagger with 4–5% error reduction. Finally, Søgaard (2009) stacks a POS tagger on an unsupervised clustering algorithm trained on large amounts of unlabeled data with mixed results.

This work combines a new semi-supervised learning method to POS tagging, namely tri-training (Li and Zhou, 2005), with stacking on unsupervised clustering. It is shown that this method can be used to improve a state-of-the-art POS tagger, SVMTool (Gimenez and Marquez, 2004). Finally, we introduce a variant of tri-training called tri-training with disagreement, which seems to perform equally well, but which imports much less unlabeled data and is therefore more efficient.

2 Tagging as classification

This section describes our dataset and our input tagger. We also describe how stacking is used to reduce POS tagging to a classification task. Finally, we introduce the supervised learning algorithms used in our experiments.

2.1 Data

We use the POS-tagged WSJ from the Penn Treebank Release 3 (Marcus et al., 1993) with the standard split: Sect. 0–18 is used for training, Sect. 19–21 for development, and Sect. 22–24 for testing. Since we need to train our classifiers on material distinct from the training material for our input POS tagger, we save Sect. 19 for training our classifiers. Finally, we use the (untagged) Brown corpus as our unlabeled data. The number of tokens we use for training, developing and testing the classifiers, and the amount of unlabeled data available to it, are thus:

	tokens
train	44,472
development	103,686
test	129,281
unlabeled	1,170,811

The amount of unlabeled data available to our classifiers is thus a bit more than 25 times the amount of labeled data.

2.2 Input tagger

In our experiments we use SVMTool (Gimenez and Marquez, 2004) with model type 4 run incrementally in both directions. SVMTool has an accuracy of 97.15% on WSJ Sect. 22-24 with this parameter setting. Gimenez and Marquez (2004) report that SVMTool has an accuracy of 97.16% with an optimized parameter setting.

2.3 Classifier input

The way classifiers are constructed in our experiments is very simple. We train SVMTool and an unsupervised tagger, Unsupos (Biemann, 2006), on our training sections and apply them to the development, test and unlabeled sections. The results are combined in tables that will be the input of our classifiers. Here is an excerpt:¹

Gold standard	SVMTool	Unsupos
DT	DT	17
NNP	NNP	27
NNP	NNS	17*
NNP	NNP	17
VBD	VBD	26

Each row represents a word and lists the gold standard POS tag, the predicted POS tag and the word cluster selected by Unsupos. For example, the first word is labeled 'DT', which SVMTool correctly predicts, and it belongs to cluster 17 of about 500 word clusters. The first column is blank in the table for the unlabeled section.

Generally, the idea is that a classifier will learn to trust SVMTool in some cases, but that it may also learn that if SVMTool predicts a certain tag for some word cluster the correct label is another tag. This way of combining taggers into a single end classifier can be seen as a form of stacking (Wolpert, 1992). It has the advantage that it reduces POS tagging to a classification task. This may simplify semi-supervised learning considerably.

2.4 Learning algorithms

We assume some knowledge of supervised learning algorithms. Most of our experiments are implementations of wrapper methods that call off-

¹The numbers provided by Unsupos refer to clusters; "*" marks out-of-vocabulary words.

the-shelf implementations of supervised learning algorithms. Specifically we have experimented with support vector machines (SVMs), decision trees, bagging and random forests. Tri-training, explained below, is a semi-supervised learning method which requires large amounts of data. Consequently, we only used very fast learning algorithms in the context of tri-training. On the development section, decision trees performed better than bagging and random forests. The decision tree algorithm is the C4.5 algorithm first introduced in Quinlan (1993). We used SVMs with polynomial kernels of degree 2 to provide a stronger stacking-only baseline.

3 Tri-training

This section first presents the tri-training algorithm originally proposed by Li and Zhou (2005) and then considers a novel variant: tri-training with disagreement.

Let L denote the labeled data and U the unlabeled data. Assume that three classifiers c_1, c_2, c_3 (same learning algorithm) have been trained on three bootstrap samples of L . In tri-training, an unlabeled datapoint in U is now labeled for a classifier, say c_1 , if the other two classifiers agree on its label, i.e. c_2 and c_3 . Two classifiers inform the third. If the two classifiers agree on a labeling, there is a good chance that they are right. The algorithm stops when the classifiers no longer change. The three classifiers are combined by majority voting. Li and Zhou (2005) show that under certain conditions the increase in classification noise rate is compensated by the amount of newly labeled data points.

The most important condition is that the three classifiers are diverse. If the three classifiers are identical, tri-training degenerates to self-training. Diversity is obtained in Li and Zhou (2005) by training classifiers on bootstrap samples. In their experiments, they consider classifiers based on the C4.5 algorithm, BP neural networks and naive Bayes classifiers. The algorithm is sketched in a simplified form in Figure 1; see Li and Zhou (2005) for all the details.

Tri-training has to the best of our knowledge not been applied to POS tagging before, but it has been applied to other NLP classification tasks, incl. Chinese chunking (Chen et al., 2006) and question classification (Nguyen et al., 2008).

```

1: for  $i \in \{1..3\}$  do
2:    $S_i \leftarrow \text{bootstrap\_sample}(L)$ 
3:    $c_i \leftarrow \text{train\_classifier}(S_i)$ 
4: end for
5: repeat
6:   for  $i \in \{1..3\}$  do
7:     for  $x \in U$  do
8:        $L_i \leftarrow \emptyset$ 
9:       if  $c_j(x) = c_k(x) (j, k \neq i)$  then
10:         $L_i \leftarrow L_i \cup \{(x, c_j(x))\}$ 
11:       end if
12:     end for
13:      $c_i \leftarrow \text{train\_classifier}(L \cup L_i)$ 
14:   end for
15: until none of  $c_i$  changes
16: apply majority vote over  $c_i$ 

```

Figure 1: Tri-training (Li and Zhou, 2005).

3.1 Tri-training with disagreement

We introduce a possible improvement of the tri-training algorithm: If we change lines 9–10 in the algorithm in Figure 1 with the lines:

```

if  $c_j(x) = c_k(x) \neq c_i(x) (j, k \neq i)$  then
   $L_i \leftarrow L_i \cup \{(x, c_j(x))\}$ 
end if

```

two classifiers, say c_1 and c_2 , only label a data-point for the third classifier, c_3 , if c_1 and c_2 agree on its label, but c_3 *disagrees*. The intuition is that we only want to strengthen a classifier in its weak points, and we want to avoid skewing our labeled data by easy data points. Finally, since tri-training with disagreement imports less unlabeled data, it is much more efficient than tri-training. No one has to the best of our knowledge applied tri-training with disagreement to real-life classification tasks before.

4 Results

Our results are presented in Figure 2. The stacking result was obtained by training a SVM on top of the predictions of SVMTool and the word clusters of Unsupos. SVMs performed better than decision trees, bagging and random forests on our development section, but improvements on test data were modest. Tri-training refers to the original algorithm sketched in Figure 1 with C4.5 as learning algorithm. Since tri-training degenerates to

self-training if the three classifiers are trained on the same sample, we used our implementation of tri-training to obtain self-training results and validated our results by a simpler implementation. We varied poolsize to optimize self-training. Finally, we list results for a technique called co-forests (Li and Zhou, 2007), which is a recent alternative to tri-training presented by the same authors, and for tri-training with disagreement (tri-disagr). The p -values are computed using 10,000 stratified shuffles.

Tri-training and tri-training with disagreement gave the best results. Note that since tri-training leads to much better results than stacking alone, it is unlabeled data that gives us most of the improvement, not the stacking itself. The difference between tri-training and self-training is near-significant ($p < 0.0150$). It seems that tri-training with disagreement is a competitive technique in terms of accuracy. The main advantage of tri-training with disagreement compared to ordinary tri-training, however, is that it is very efficient. This is reflected by the average number of tokens in L_i over the three learners in the worst round of learning:

	av. tokens in L_i
tri-training	1,170,811
tri-disagr	173

Note also that self-training gave very good results. Self-training was, again, much slower than tri-training with disagreement since we had to train on a large pool of unlabeled data (but only once). Of course this is not a standard self-training set-up, but self-training informed by unsupervised word clusters.

4.1 Follow-up experiments

SVMTool is one of the most accurate POS taggers available. This means that the predictions that are added to the labeled data are of very high quality. To test if our semi-supervised learning methods were sensitive to the quality of the input taggers we repeated the self-training and tri-training experiments with a less competitive POS tagger, namely the maximum entropy-based POS tagger first described in (Ratnaparkhi, 1998) that comes with the maximum entropy library in (Zhang, 2004). Results are presented as the second line in Figure 2. Note that error reduction is much lower in this case.

	BL	stacking	tri-tr.	self-tr.	co-forests	tri-disagr	error red.	<i>p</i> -value
SVMTool	97.15%	97.19%	97.27%	97.26%	97.13%	97.27%	4.21%	<0.0001
MaxEnt	96.31%	-	96.36%	96.36%	96.28%	96.36%	1.36%	<0.0001

Figure 2: Results on Wall Street Journal Sect. 22-24 with different semi-supervised methods.

5 Conclusion

This paper first shows how stacking can be used to reduce POS tagging to a classification task. This reduction seems to enable robust semi-supervised learning. The technique was used to improve the accuracy of a state-of-the-art POS tagger, namely SVMTool. Four semi-supervised learning methods were tested, incl. self-training, tri-training, co-forests and tri-training with disagreement. All methods increased the accuracy of SVMTool significantly. Error reduction on Wall Street Journal Sect. 22-24 was 4.2%, which is comparable to related work in the literature, e.g. Suzuki and Isozaki (2008) (7%) and Spoustova et al. (2009) (4–5%).

References

- Chris Biemann. 2006. Unsupervised part-of-speech tagging employing efficient graph clustering. In *COLING-ACL Student Session*, Sydney, Australia.
- Wenliang Chen, Yujie Zhang, and Hitoshi Isahara. 2006. Chinese chunking with tri-training learning. In *Computer processing of oriental languages*, pages 466–473. Springer, Berlin, Germany.
- Stephen Clark, James Curran, and Mike Osborne. 2003. Bootstrapping POS taggers using unlabeled data. In *CONLL*, Edmonton, Canada.
- Jesus Gimenez and Lluís Marquez. 2004. SVMTool: a general POS tagger generator based on support vector machines. In *LREC*, Lisbon, Portugal.
- Zhongqiang Huang, Vladimir Eidelman, and Mary Harper. 2009. Improving a simple bigram HMM part-of-speech tagger by latent annotation and self-training. In *NAACL-HLT*, Boulder, CO.
- Ming Li and Zhi-Hua Zhou. 2005. Tri-training: exploiting unlabeled data using three classifiers. *IEEE Transactions on Knowledge and Data Engineering*, 17(11):1529–1541.
- Ming Li and Zhi-Hua Zhou. 2007. Improve computer-aided diagnosis with machine learning techniques using undiagnosed samples. *IEEE Transactions on Systems, Man and Cybernetics*, 37(6):1088–1098.
- Mitchell Marcus, Mary Marcinkiewicz, and Beatrice Santorini. 1993. Building a large annotated corpus of English: the Penn Treebank. *Computational Linguistics*, 19(2):313–330.
- Bernard Merialdo. 1994. Tagging English text with a probabilistic model. *Computational Linguistics*, 20(2):155–171.
- Tri Nguyen, Le Nguyen, and Akira Shimazu. 2008. Using semi-supervised learning for question classification. *Journal of Natural Language Processing*, 15:3–21.
- Ross Quinlan. 1993. *Programs for machine learning*. Morgan Kaufmann.
- Adwait Ratnaparkhi. 1998. *Maximum entropy models for natural language ambiguity resolution*. Ph.D. thesis, University of Pennsylvania.
- Anders Søgaard. 2009. Ensemble-based POS tagging of Italian. In *IAAI-EVALITA*, Reggio Emilia, Italy.
- Drahomira Spoustova, Jan Hajic, Jan Raab, and Miroslav Spousta. 2009. Semi-supervised training for the averaged perceptron POS tagger. In *EACL*, Athens, Greece.
- Jun Suzuki and Hideki Isozaki. 2008. Semi-supervised sequential labeling and segmentation using gigaword scale unlabeled data. In *ACL*, pages 665–673, Columbus, Ohio.
- Wen Wang, Zhongqiang Huang, and Mary Harper. 2007. Semi-supervised learning for part-of-speech tagging of Mandarin transcribed speech. In *ICASSP*, Hawaii.
- David Wolpert. 1992. Stacked generalization. *Neural Networks*, 5:241–259.
- Le Zhang. 2004. Maximum entropy modeling toolkit for Python and C++. University of Edinburgh.

Efficient Optimization of an MDL-Inspired Objective Function for Unsupervised Part-of-Speech Tagging

Ashish Vaswani¹

Adam Pauls²

David Chiang¹

¹Information Sciences Institute
University of Southern California
4676 Admiralty Way, Suite 1001
Marina del Rey, CA 90292
{avaswani, chiang}@isi.edu

²Computer Science Division
University of California at Berkeley
Soda Hall
Berkeley, CA 94720
adpauls@eecs.berkeley.edu

Abstract

The Minimum Description Length (MDL) principle is a method for model selection that trades off between the explanation of the data by the model and the complexity of the model itself. Inspired by the MDL principle, we develop an objective function for generative models that captures the description of the data by the model (log-likelihood) and the description of the model (model size). We also develop an efficient general search algorithm based on the MAP-EM framework to optimize this function. Since recent work has shown that minimizing the model size in a Hidden Markov Model for part-of-speech (POS) tagging leads to higher accuracies, we test our approach by applying it to this problem. The search algorithm involves a simple change to EM and achieves high POS tagging accuracies on both English and Italian data sets.

1 Introduction

The Minimum Description Length (MDL) principle is a method for model selection that provides a generic solution to the overfitting problem (Barron et al., 1998). A formalization of Ockham’s Razor, it says that the parameters are to be chosen that minimize the description length of the data given the model plus the description length of the model itself.

It has been successfully shown that minimizing the model size in a Hidden Markov Model (HMM) for part-of-speech (POS) tagging leads to higher accuracies than simply running the Expectation-Maximization (EM) algorithm (Dempster et al., 1977). Goldwater and Griffiths (2007) employ a Bayesian approach to POS tagging and use sparse Dirichlet priors to minimize model size. More re-

cently, Ravi and Knight (2009) alternately minimize the model using an integer linear program and maximize likelihood using EM to achieve the highest accuracies on the task so far. However, in the latter approach, because there is no single objective function to optimize, it is not entirely clear how to generalize this technique to other problems. In this paper, inspired by the MDL principle, we develop an objective function for generative models that captures both the description of the data by the model (log-likelihood) and the description of the model (model size). By using a simple prior that encourages sparsity, we cast our problem as a search for the maximum *a posteriori* (MAP) hypothesis and present a variant of EM to approximately search for the minimum-description-length model. Applying our approach to the POS tagging problem, we obtain higher accuracies than both EM and Bayesian inference as reported by Goldwater and Griffiths (2007). On an Italian POS tagging task, we obtain even larger improvements. We find that our objective function correlates well with accuracy, suggesting that this technique might be useful for other problems.

2 MAP EM with Sparse Priors

2.1 Objective function

In the unsupervised POS tagging task, we are given a word sequence $\mathbf{w} = w_1, \dots, w_N$ and want to find the best tagging $\mathbf{t} = t_1, \dots, t_N$, where $t_i \in \mathcal{T}$, the tag vocabulary. We adopt the problem formulation of Merialdo (1994), in which we are given a dictionary of possible tags for each word type.

We define a bigram HMM

$$P(\mathbf{w}, \mathbf{t} | \theta) = \prod_{i=1}^N P(\mathbf{w}, \mathbf{t} | \theta) \cdot P(t_i | t_{i-1}) \quad (1)$$

In maximum likelihood estimation, the goal is to

find parameter estimates

$$\hat{\theta} = \arg \max_{\theta} \log P(\mathbf{w} | \theta) \quad (2)$$

$$= \arg \max_{\theta} \log \sum_{\mathbf{t}} P(\mathbf{w}, \mathbf{t} | \theta) \quad (3)$$

The EM algorithm can be used to find a solution. However, we would like to maximize likelihood and minimize the size of the model simultaneously. We define the size of a model as the number of non-zero probabilities in its parameter vector. Let $\theta_1, \dots, \theta_n$ be the components of θ . We would like to find

$$\hat{\theta} = \arg \min_{\theta} (-\log P(\mathbf{w} | \theta) + \alpha \|\theta\|_0) \quad (4)$$

where $\|\theta\|_0$, called the L0 norm of θ , simply counts the number of non-zero parameters in θ . The hyperparameter α controls the tradeoff between likelihood maximization and model minimization. Note the similarity of this objective function with MDL's, where α would be the space (measured in nats) needed to describe one parameter of the model.

Unfortunately, minimization of the L0 norm is known to be NP-hard (Hyder and Mahata, 2009). It is not smooth, making it unamenable to gradient-based optimization algorithms. Therefore, we use a smoothed approximation,

$$\|\theta\|_0 \approx \sum_i \left(1 - e^{-\frac{\theta_i}{\beta}}\right) \quad (5)$$

where $0 < \beta \leq 1$ (Mohimani et al., 2007). For smaller values of β , this closely approximates the desired function (Figure 1). Inverting signs and ignoring constant terms, our objective function is now:

$$\hat{\theta} = \arg \max_{\theta} \left(\log P(\mathbf{w} | \theta) + \alpha \sum_i e^{-\frac{\theta_i}{\beta}} \right) \quad (6)$$

We can think of the approximate model size as a kind of prior:

$$P(\theta) = \frac{\exp \alpha \sum_i e^{-\frac{\theta_i}{\beta}}}{Z} \quad (7)$$

$$\log P(\theta) = \alpha \cdot \sum_i e^{-\frac{\theta_i}{\beta}} - \log Z \quad (8)$$

where $Z = \int_{d\theta} \exp \alpha \sum_i e^{-\frac{\theta_i}{\beta}}$ is a normalization constant. Then our goal is to find the maximum

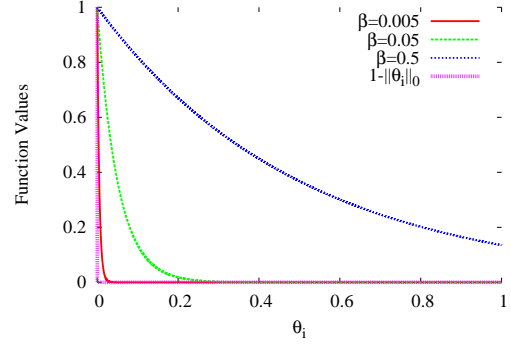


Figure 1: Ideal model-size term and its approximations.

a posteriori parameter estimate, which we find using MAP-EM (Bishop, 2006):

$$\hat{\theta} = \arg \max_{\theta} \log P(\mathbf{w}, \theta) \quad (9)$$

$$= \arg \max_{\theta} (\log P(\mathbf{w} | \theta) + \log P(\theta)) \quad (10)$$

Substituting (8) into (10) and ignoring the constant term $\log Z$, we get our objective function (6) again.

We can exercise finer control over the sparsity of the tag-bigram and channel probability distributions by using a different α for each:

$$\arg \max_{\theta} \left(\log P(\mathbf{w} | \theta) + \alpha_c \sum_{w,t} e^{-\frac{P(w|t)}{\beta}} + \alpha_t \sum_{t,t'} e^{-\frac{P(t'|t)}{\beta}} \right) \quad (11)$$

In our experiments, we set $\alpha_c = 0$ since previous work has shown that minimizing the number of tag n -gram parameters is more important (Ravi and Knight, 2009; Goldwater and Griffiths, 2007).

A common method for preferring smaller models is minimizing the L1 norm, $\sum_i |\theta_i|$. However, for a model which is a product of multinomial distributions, the L1 norm is a constant.

$$\begin{aligned} \sum_i |\theta_i| &= \sum_i \theta_i \\ &= \sum_t \left(\sum_w P(w | t) + \sum_{t'} P(t' | t) \right) \\ &= 2|\mathcal{T}| \end{aligned}$$

Therefore, we cannot use the L1 norm as part of the size term as the result will be the same as the EM algorithm.

2.2 Parameter optimization

To optimize (11), we use MAP EM, which is an iterative search procedure. The E step is the same as in standard EM, which is to calculate $P(\mathbf{t} \mid \mathbf{w}, \theta^t)$, where the θ^t are the parameters in the current iteration t . The M step in iteration $(t + 1)$ looks like

$$\theta^{t+1} = \arg \max_{\theta} \left(E_{P(\mathbf{t} \mid \mathbf{w}, \theta^t)} [\log P(\mathbf{w}, \mathbf{t} \mid \theta)] + \alpha_t \sum_{t, t'} e^{-\frac{P(t'|t)}{\beta}} \right) \quad (12)$$

Let $C(t, w; \mathbf{t}, \mathbf{w})$ count the number of times the word w is tagged as t in \mathbf{t} , and $C(t, t'; \mathbf{t})$ the number of times the tag bigram (t, t') appears in \mathbf{t} . We can rewrite the M step as

$$\theta^{t+1} = \arg \max_{\theta} \left(\sum_t \sum_w E[C(t, w)] \log P(w \mid t) + \sum_t \sum_{t'} \left(E[C(t, t')] \log P(t' \mid t) + \alpha_t e^{-\frac{P(t'|t)}{\beta}} \right) \right) \quad (13)$$

subject to the constraints $\sum_w P(w \mid t) = 1$ and $\sum_{t'} P(t' \mid t) = 1$. Note that we can optimize each term of both summations over t separately. For each t , the term

$$\sum_w E[C(t, w)] \log P(w \mid t) \quad (14)$$

is easily optimized as in EM: just let $P(w \mid t) \propto E[C(t, w)]$. But the term

$$\sum_{t'} \left(E[C(t, t')] \log P(t' \mid t) + \alpha_t e^{-\frac{P(t'|t)}{\beta}} \right) \quad (15)$$

is trickier. This is a non-convex optimization problem for which we invoke a publicly available constrained optimization tool, ALGENCAN (Andreani et al., 2007). To carry out its optimization, ALGENCAN requires computation of the following in every iteration:

- **Objective function**, defined in equation (15). This is calculated in polynomial time using dynamic programming.
- **Constraints**: $g_t = \sum_{t'} P(t' \mid t) - 1 = 0$ for each tag $t \in \mathcal{T}$. Also, we constrain $P(t' \mid t)$ to the interval $[\epsilon, 1]$.¹

¹We must have $\epsilon > 0$ because of the $\log P(t' \mid t)$ term in equation (15). It seems reasonable to set $\epsilon \ll \frac{1}{N}$; in our experiments, we set $\epsilon = 10^{-7}$.

- **Gradient of objective function**:

$$\frac{\partial F}{\partial P(t' \mid t)} = \frac{E[C(t, t')]}{P(t' \mid t)} - \frac{\alpha_t}{\beta} e^{-\frac{P(t'|t)}{\beta}} \quad (16)$$

- **Gradient of equality constraints**:

$$\frac{\partial g_t}{\partial P(t'' \mid t')} = \begin{cases} 1 & \text{if } t = t' \\ 0 & \text{otherwise} \end{cases} \quad (17)$$

- **Hessian of objective function**, which is not required but greatly speeds up the optimization:

$$\frac{\partial^2 F}{\partial P(t' \mid t) \partial P(t' \mid t)} = -\frac{E[C(t, t')]}{P(t' \mid t)^2} + \alpha_t \frac{e^{-\frac{P(t'|t)}{\beta}}}{\beta^2} \quad (18)$$

The other second-order partial derivatives are all zero, as are those of the equality constraints.

We perform this optimization for each instance of (15). These optimizations could easily be performed in parallel for greater scalability.

3 Experiments

We carried out POS tagging experiments on English and Italian.

3.1 English POS tagging

To set the hyperparameters α_t and β , we prepared three held-out sets H_1, H_2 , and H_3 from the Penn Treebank. Each H_i comprised about 24,000 words annotated with POS tags. We ran MAP-EM for 100 iterations, with uniform probability initialization, for a suite of hyperparameters and averaged their tagging accuracies over the three held-out sets. The results are presented in Table 2. We then picked the hyperparameter setting with the highest average accuracy. These were $\alpha_t = 80, \beta = 0.05$. We then ran MAP-EM again on the test data with these hyperparameters and achieved a tagging accuracy of 87.4% (see Table 1). This is higher than the 85.2% that Goldwater and Griffiths (2007) obtain using Bayesian methods for inferring both POS tags and hyperparameters. It is much higher than the 82.4% that standard EM achieves on the test set when run for 100 iterations.

Using $\alpha_t = 80, \beta = 0.05$, we ran multiple random restarts on the test set (see Figure 2). We find that the objective function correlates well with accuracy, and picking the point with the highest objective function value achieves 87.1% accuracy.

α_t	β								
	0.75	0.5	0.25	0.075	0.05	0.025	0.0075	0.005	0.0025
10	82.81	82.78	83.10	83.50	83.76	83.70	84.07	83.95	83.75
20	82.78	82.82	83.26	83.60	83.89	84.88	83.74	84.12	83.46
30	82.78	83.06	83.26	83.29	84.50	84.82	84.54	83.93	83.47
40	82.81	83.13	83.50	83.98	84.23	85.31	85.05	83.84	83.46
50	82.84	83.24	83.15	84.08	82.53	84.90	84.73	83.69	82.70
60	83.05	83.14	83.26	83.30	82.08	85.23	85.06	83.26	82.96
70	83.09	83.10	82.97	82.37	83.30	86.32	83.98	83.55	82.97
80	83.13	83.15	82.71	83.00	86.47	86.24	83.94	83.26	82.93
90	83.20	83.18	82.53	84.20	86.32	84.87	83.49	83.62	82.03
100	83.19	83.51	82.84	84.60	86.13	85.94	83.26	83.67	82.06
110	83.18	83.53	83.29	84.40	86.19	85.18	80.76	83.32	82.05
120	83.08	83.65	83.71	84.11	86.03	85.39	80.66	82.98	82.20
130	83.10	83.19	83.52	84.02	85.79	85.65	80.08	82.04	81.76
140	83.11	83.17	83.34	85.26	85.86	85.84	79.09	82.51	81.64
150	83.14	83.20	83.40	85.33	85.54	85.18	78.90	81.99	81.88

Table 2: Average accuracies over three held-out sets for English.

system	accuracy (%)
Standard EM	82.4
+ random restarts	84.5
(Goldwater and Griffiths, 2007)	85.2
our approach	87.4
+ random restarts	87.1

Table 1: MAP-EM with a L0 norm achieves higher tagging accuracy on English than (2007) and much higher than standard EM.

system	zero parameters	bigram types
maximum possible	1389	–
EM, 100 iterations	444	924
MAP-EM, 100 iterations	695	648

Table 3: MAP-EM with a smoothed L0 norm yields much smaller models than standard EM.

We also carried out the same experiment with standard EM (Figure 3), where picking the point with the highest corpus probability achieves 84.5% accuracy.

We also measured the minimization effect of the sparse prior against that of standard EM. Since our method lower-bounds all the parameters by ϵ , we consider a parameter θ_i as a zero if $\theta_i \leq \epsilon$. We also measured the number of unique tag bigram types in the Viterbi tagging of the word sequence. Table 3 shows that our method produces much smaller models than EM, and produces Viterbi taggings with many fewer tag-bigram types.

3.2 Italian POS tagging

We also carried out POS tagging experiments on an Italian corpus from the Italian Turin Univer-

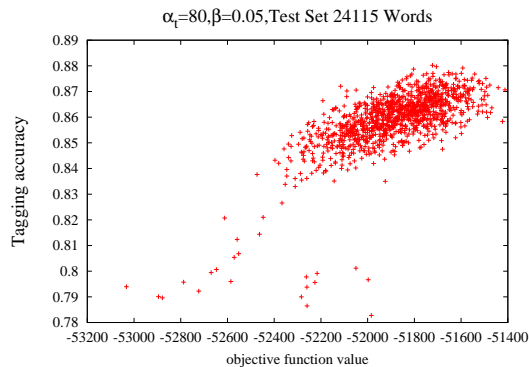


Figure 2: Tagging accuracy vs. objective function for 1152 random restarts of MAP-EM with smoothed L0 norm.

sity Treebank (Bos et al., 2009). This test set comprises 21, 878 words annotated with POS tags and a dictionary for each word type. Since this is all the available data, we could not tune the hyperparameters on a held-out data set. Using the hyperparameters tuned on English ($\alpha_t = 80, \beta = 0.05$), we obtained 89.7% tagging accuracy (see Table 4), which was a large improvement over 81.2% that standard EM achieved. When we tuned the hyperparameters on the test set, the best setting ($\alpha_t = 120, \beta = 0.05$) gave an accuracy of 90.28%.

4 Conclusion

A variety of other techniques in the literature have been applied to this unsupervised POS tagging task. Smith and Eisner (2005) use conditional random fields with contrastive estimation to achieve

α_t	β								
	0.75	0.5	0.25	0.075	0.05	0.025	0.0075	0.005	0.0025
10	81.62	81.67	81.63	82.47	82.70	84.64	84.82	84.96	84.90
20	81.67	81.63	81.76	82.75	84.28	84.79	85.85	88.49	85.30
30	81.66	81.63	82.29	83.43	85.08	88.10	86.16	88.70	88.34
40	81.64	81.79	82.30	85.00	86.10	88.86	89.28	88.76	88.80
50	81.71	81.71	78.86	85.93	86.16	88.98	88.98	89.11	88.01
60	81.65	82.22	78.95	86.11	87.16	89.35	88.97	88.59	88.00
70	81.69	82.25	79.55	86.32	89.79	89.37	88.91	85.63	87.89
80	81.74	82.23	80.78	86.34	89.70	89.58	88.87	88.32	88.56
90	81.70	81.85	81.00	86.35	90.08	89.40	89.09	88.09	88.50
100	81.70	82.27	82.24	86.53	90.07	88.93	89.09	88.30	88.72
110	82.19	82.49	82.22	86.77	90.12	89.22	88.87	88.48	87.91
120	82.23	78.60	82.76	86.77	90.28	89.05	88.75	88.83	88.53
130	82.20	78.60	83.33	87.48	90.12	89.15	89.30	87.81	88.66
140	82.24	78.64	83.34	87.48	90.12	89.01	88.87	88.99	88.85
150	82.28	78.69	83.32	87.75	90.25	87.81	88.50	89.07	88.41

Table 4: Accuracies on test set for Italian.

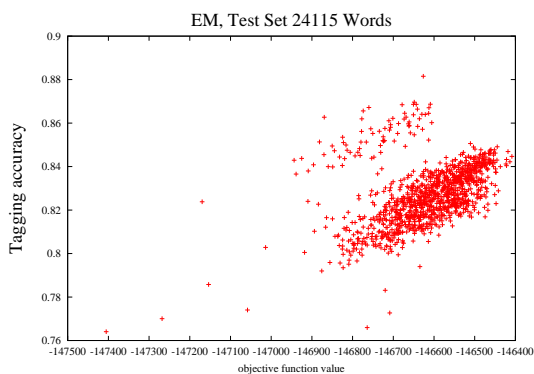


Figure 3: Tagging accuracy vs. likelihood for 1152 random restarts of standard EM.

88.6% accuracy. Goldberg et al. (2008) provide a linguistically-informed starting point for EM to achieve 91.4% accuracy. More recently, Chiang et al. (2010) use Gibbs sampling for Bayesian inference along with automatic run selection and achieve 90.7%.

In this paper, our goal has been to investigate whether EM can be extended in a generic way to use an MDL-like objective function that simultaneously maximizes likelihood and minimizes model size. We have presented an efficient search procedure that optimizes this function for generative models and demonstrated that maximizing this function leads to improvement in tagging accuracy over standard EM. We infer the hyperparameters of our model using held out data and achieve better accuracies than (Goldwater and Griffiths, 2007). We have also shown that the objective function correlates well with tagging accu-

racy supporting the MDL principle. Our approach performs quite well on POS tagging for both English and Italian. We believe that, like EM, our method can benefit from more unlabeled data, and there is reason to hope that the success of these experiments will carry over to other tasks as well.

Acknowledgements

We would like to thank Sujith Ravi, Kevin Knight and Steve DeNeefe for their valuable input, and Jason Baldridge for directing us to the Italian POS data. This research was supported in part by DARPA contract HR0011-06-C-0022 under sub-contract to BBN Technologies and DARPA contract HR0011-09-1-0028.

References

- R. Andreati, E. G. Birgin, J. M. Martinez, and M. L. Schuverdt. 2007. On Augmented Lagrangian methods with general lower-level constraints. *SIAM Journal on Optimization*, 18:1286–1309.
- A. Barron, J. Rissanen, and B. Yu. 1998. The minimum description length principle in coding and modeling. *IEEE Transactions on Information Theory*, 44(6):2743–2760.
- C. Bishop. 2006. *Pattern Recognition and Machine Learning*. Springer.
- J. Bos, C. Bosco, and A. Mazzei. 2009. Converting a dependency treebank to a categorical grammar treebank for italian. In *Eighth International Workshop on Treebanks and Linguistic Theories (TLT8)*.
- D. Chiang, J. Graehl, K. Knight, A. Pauls, and S. Ravi. 2010. Bayesian inference for Finite-State transducers. In *Proceedings of the North American Association of Computational Linguistics*.

- A. P. Dempster, N. M. Laird, and D. B. Rubin. 1977. Maximum likelihood from incomplete data via the EM algorithm. *Computational Linguistics*, 39(4):1–38.
- Y. Goldberg, M. Adler, and M. Elhadad. 2008. EM can find pretty good HMM POS-taggers (when given a good start). In *Proceedings of the ACL*.
- S. Goldwater and T. L. Griffiths. 2007. A fully Bayesian approach to unsupervised part-of-speech tagging. In *Proceedings of the ACL*.
- M. Hyder and K. Mahata. 2009. An approximate L0 norm minimization algorithm for compressed sensing. In *Proceedings of the 2009 IEEE International Conference on Acoustics, Speech and Signal Processing*.
- B. Merialdo. 1994. Tagging English text with a probabilistic model. *Computational Linguistics*, 20(2):155–171.
- H. Mohimani, M. Babaie-Zadeh, and C. Jutten. 2007. Fast sparse representation based on smoothed L0 norm. In *Proceedings of the 7th International Conference on Independent Component Analysis and Signal Separation (ICA2007)*.
- S. Ravi and K. Knight. 2009. Minimized models for unsupervised part-of-speech tagging. In *Proceedings of ACL-IJCNLP*.
- N. Smith, and J. Eisner. 2005. Contrastive estimation: Training log-linear models on unlabeled data. In *Proceedings of the ACL*.

SVD and Clustering for Unsupervised POS Tagging

Michael Lamar*

Division of Applied Mathematics
Brown University
Providence, RI, USA
mlamar@dam.brown.edu

Yariv Maron*

Gonda Brain Research Center
Bar-Ilan University
Ramat-Gan, Israel
syarivm@yahoo.com

Mark Johnson

Department of Computing
Faculty of Science
Macquarie University
Sydney, Australia
mjohnson@science.mq.edu.au

Elie Bienenstock

Division of Applied Mathematics
and Department of Neuroscience
Brown University
Providence, RI, USA
elie@brown.edu

Abstract

We revisit the algorithm of Schütze (1995) for unsupervised part-of-speech tagging. The algorithm uses reduced-rank singular value decomposition followed by clustering to extract latent features from context distributions. As implemented here, it achieves state-of-the-art tagging accuracy at considerably less cost than more recent methods. It can also produce a range of finer-grained taggings, with potential applications to various tasks.

1 Introduction

While supervised approaches are able to solve the part-of-speech (POS) tagging problem with over 97% accuracy (Collins 2002; Toutanova et al. 2003), unsupervised algorithms perform considerably less well. These models attempt to tag text without resources such as an annotated corpus, a dictionary, etc. The use of singular value decomposition (SVD) for this problem was introduced in Schütze (1995). Subsequently, a number of methods for POS tagging without a dictionary were examined, e.g., by Clark (2000), Clark (2003), Haghighi and Klein (2006), Johnson (2007), Goldwater and Griffiths (2007), Gao and Johnson (2008), and Graça et al. (2009). The latter two, using Hidden Markov Models (HMMs), exhibit the highest performances to

date for fully unsupervised POS tagging.

The revisited SVD-based approach presented here, which we call “two-step SVD” or SVD2, has four important characteristics. First, it achieves state-of-the-art tagging accuracy. Second, it requires drastically less computational effort than the best currently available models. Third, it demonstrates that state-of-the-art accuracy can be realized without disambiguation, *i.e.*, without attempting to assign different tags to different tokens of the same type. Finally, with no significant increase in computational cost, SVD2 can create much finer-grained labelings than typically produced by other algorithms. When combined with some minimal supervision in post-processing, this makes the approach useful for tagging languages that lack the resources required by fully supervised models.

2 Methods

Following the original work of Schütze (1995), we begin by constructing a right context matrix, R , and a left context matrix, L . R_{ij} counts the number of times in the corpus a token of word type i is immediately followed by a token of word type j . Similarly, L_{ij} counts the number of times a token of type i is preceded by a token of type j . We truncate these matrices, including, in the right and left contexts, only the w_1 most frequent word types. The resulting L and R are of dimension $N_{\text{types}} \times w_1$, where N_{types} is the number of word types (spelling forms) in the corpus, and w_1 is set to 1000. (The full $N_{\text{types}} \times N_{\text{types}}$ context matrices satisfy $R = L^T$.)

* These authors contributed equally.

Next, both context matrices are factored using singular value decomposition:

$$\begin{aligned} L &= U_L S_L V_L^T \\ R &= U_R S_R V_R^T. \end{aligned}$$

The diagonal matrices S_L and S_R (each of rank 1000) are reduced down to rank $r_1 = 100$ by replacing the 900 smallest singular values in each matrix with zeros, yielding S_L^* and S_R^* . We then form a pair of latent-descriptor matrices defined by:

$$\begin{aligned} L^* &= U_L S_L^* \\ R^* &= U_R S_R^*. \end{aligned}$$

Row i in matrix L^* (resp. R^*) is the left (resp. right) latent descriptor for word type i . We next include a normalization step in which each row in each of L^* and R^* is scaled to unit length, yielding matrices L^{**} and R^{**} . Finally, we form a single descriptor matrix D by concatenating these matrices into $D = [L^{**} R^{**}]$. Row i in matrix D is the complete latent descriptor for word type i ; this latent descriptor sits on the Cartesian product of two 100-dimensional unit spheres, hereafter the *2-sphere*.

We next categorize these descriptors into $k_1 = 500$ groups, using a k -means clustering algorithm. Centroid initialization is done by placing the k initial centroids on the descriptors of the k most frequent words in the corpus. As the descriptors sit on the 2-sphere, we measure the proximity of a descriptor to a centroid by the dot product between them; this is equal to the sum of the cosines of the angles—computed on the left and right parts—between them. We update each cluster’s centroid as the weighted average of its constituents, the weight being the frequency of the word type; the centroids are then scaled, so they sit on the 2-sphere. Typically, only a few dozen iterations are required for full convergence of the clustering algorithm.

We then apply a second pass of this entire SVD-and-clustering procedure. In this second pass, we use the $k_1 = 500$ clusters from the first iteration to assemble a new pair of context matrices. Now, R_{ij} counts all the cluster- j ($j=1 \dots k_1$) words to the right of word i , and L_{ij} counts all the cluster- j words to the left of word i . The new matrices L and R have dimension $N_{\text{types}} \times k_1$.

As in the first pass, we perform reduced-rank SVD, this time down to rank $r_2 = 300$, and we again normalize the descriptors to unit length, yielding a new pair of latent descriptor matrices L^{**} and R^{**} . Finally, we concatenate L^{**} and R^{**} into a single matrix of descriptors, and cluster these descriptors into k_2 groups, where k_2 is the desired number of induced tags. We use the same

weighted k -means algorithm as in the first pass, again placing the k initial centroids on the descriptors of the k most frequent words in the corpus. The final tag of any token in the corpus is the cluster number of its type.

3 Data and Evaluation

We ran the SVD2 algorithm described above on the full Wall Street Journal part of the Penn Treebank (1,173,766 tokens). Capitalization was ignored, resulting in $N_{\text{types}} = 43,766$, with only a minor effect on accuracy. Evaluation was done against the POS-tag annotations of the 45-tag PTB tagset (hereafter PTB45), and against the Smith and Eisner (2005) coarse version of the PTB tagset (hereafter PTB17). We selected the three evaluation criteria of Gao and Johnson (2008): M-to-1, 1-to-1, and VI. M-to-1 and 1-to-1 are the tagging accuracies under the best many-to-one map and the greedy one-to-one map respectively; VI is a map-free information-theoretic criterion—see Gao and Johnson (2008) for details. Although we find M-to-1 to be the most reliable criterion of the three, we include the other two criteria for completeness.

In addition to the best M-to-1 map, we also employ here, for large values of k_2 , a *prototype-based M-to-1 map*. To construct this map, we first find, for each induced tag t , the word type with which it co-occurs most frequently; we call this word type the *prototype* of t . We then query the annotated data for the most common gold tag for each prototype, and we map induced tag t to this gold tag. This prototype-based M-to-1 map produces accuracy scores no greater—typically lower—than the best M-to-1 map. We discuss the value of this approach as a minimally-supervised post-processing step in Section 5.

4 Results

Low- k performance. Here we present the performance of the SVD2 model when k_2 , the number of induced tags, is the same or roughly the same as the number of tags in the gold standard—hence small. Table 1 compares the performance of SVD2 to other leading models. Following Gao and Johnson (2008), the number of induced tags is 17 for PTB17 evaluation and 50 for PTB45 evaluation. Thus, with the exception of Graça et al. (2009) who use 45 induced tags for PTB45, the number of induced tags is the same across each column of Table 1.

Model	M-to-1		1-to-1		VI	
	PTB17	PTB45	PTB17	PTB45	PTB17	PTB45
SVD2	0.730	0.660	0.513	0.467	3.02	3.84
HMM-EM	0.647	0.621	0.431	0.405	3.86	4.48
HMM-VB	0.637	0.605	0.514	0.461	3.44	4.28
HMM-GS	0.674	0.660	0.466	0.499	3.46	4.04
HMM-Sparse(32)	0.702(2.2)	0.654(1.0)	0.495	0.445		
VEM ($10^{-1}, 10^{-1}$)	0.682(0.8)	0.546(1.7)	0.528	0.460		

Table 1. Tagging accuracy under the best M-to-1 map, the greedy 1-to-1 map, and VI, for the full PTB45 tagset and the reduced PTB17 tagset. HMM-EM, HMM-VB and HMM-GS show the best results from Gao and Johnson (2008); HMM-Sparse(32) and VEM ($10^{-1}, 10^{-1}$) show the best results from Graça et al. (2009).

The performance of SVD2 compares favorably to the HMM models. Note that SVD2 is a deterministic algorithm. The table shows, in parentheses, the standard deviations reported in Graça et al. (2009). For the sake of comparison with Graça et al. (2009), we also note that, with $k_2 = 45$, SVD2 scores 0.659 on PTB45. The NVI scores (Reichart and Rappoport 2009) corresponding to the VI scores for SVD2 are 0.938 for PTB17 and 0.885 for PTB45. To examine the sensitivity of the algorithm to its four parameters, w_1 , r_1 , k_1 , and r_2 , we changed each of these parameters separately by a multiplicative factor of either 0.5 or 2; in neither case did M-to-1 accuracy drop by more than 0.014.

This performance was achieved despite the fact that the SVD2 tagger is mathematically much simpler than the other models. Our MATLAB implementation of SVD2 takes only a few minutes to run on a desktop computer, in contrast to HMM training times of several hours or days (Gao and Johnson 2008; Johnson 2007).

High- k performance. Not suffering from the same computational limitations as other models, SVD2 can easily accommodate high numbers of induced tags, resulting in fine-grained labelings. The value of this flexibility is discussed in the next section. Figure 1 shows, as a function of k_2 , the tagging accuracy of SVD2 under both the best and the prototype-based M-to-1 maps (see Section 3), for both the PTB45 and the PTB17 tagsets. The horizontal one-tag-per-word-type line in each panel is the theoretical upper limit for tagging accuracy in non-disambiguating models (such as SVD2). This limit is the fraction of all tokens in the corpus whose gold tag is the most frequent for their type.

5 Discussion

At the heart of the algorithm presented here is the reduced-rank SVD method of Schütze (1995), which transforms bigram counts into latent descriptors. In view of the present work,

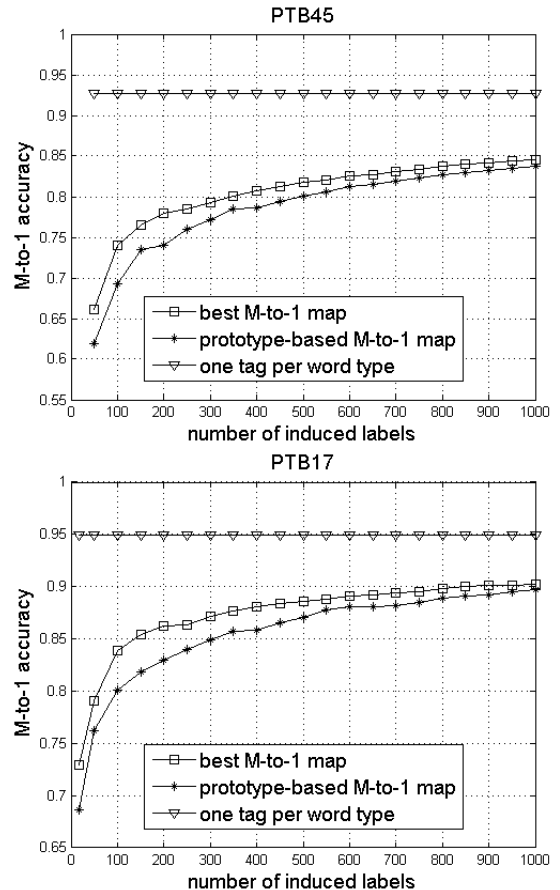


Figure 1. Performance of the SVD2 algorithm as a function of the number of induced tags. Top: PTB45; bottom: PTB17. Each plot shows the tagging accuracy under the best and the prototype-based M-to-1 maps, as well as the upper limit for non-disambiguating taggers.

which achieves state-of-the-art performance when evaluation is done with the criteria now in common use, Schütze's original work should rightly be praised as ahead of its time. The SVD2 model presented here differs from Schütze's work in many details of implementation—not all of which are explicitly specified in Schütze (1995). In what follows, we discuss the features of SVD2 that are most critical to its performance. Failure to incorporate any one of them signifi-

cantly reduces the performance of the algorithm (M-to-1 reduced by 0.04 to 0.08).

First, the reduced-rank left-singular vectors (for the right and left context matrices) are scaled, *i.e.*, multiplied, by the singular values. While the resulting descriptors, the rows of L^* and R^* , live in a much lower-dimensional space than the original context vectors, they are mapped by an angle-preserving map (defined by the matrices of right-singular vectors V_L and V_R) into vectors in the original space. These mapped vectors best approximate (in the least-squares sense) the original context vectors; they have the same geometric relationships as their equivalent high-dimensional images, making them good candidates for the role of word-type descriptors.

A second important feature of the SVD2 algorithm is the unit-length normalization of the latent descriptors, along with the computation of cluster centroids as the weighted averages of their constituent vectors. Thanks to this combined device, rare words are treated equally to frequent words regarding the length of their descriptor vectors, yet contribute less to the placement of centroids.

Finally, while the usual drawback of k -means-clustering algorithms is the dependency of the outcome on the initial—usually random—placement of centroids, our initialization of the k centroids as the descriptors of the k most frequent word types in the corpus makes the algorithm fully deterministic, and improves its performance substantially: M-to-1 PTB45 by 0.043, M-to-1 PTB17 by 0.063.

As noted in the Results section, SVD2 is fairly robust to changes in all four parameters w_1 , r_1 , k_1 , and r_2 . The values used here were obtained by a coarse, greedy strategy, where each parameter was optimized independently. It is worth noting that dispensing with the second pass altogether, *i.e.*, clustering directly the latent descriptor vectors obtained in the first pass into the desired number of induced tags, results in a drop of Many-to-1 score of only 0.021 for the PTB45 tagset and 0.009 for the PTB17 tagset.

Disambiguation. An obvious limitation of SVD2 is that it is a non-disambiguating tagger, assigning the same label to all tokens of a type. However, this limitation *per se* is unlikely to be the main obstacle to the improvement of low- k performance, since, as is well known, the theoretical upper limit for the tagging accuracy of non-disambiguating models (shown in Fig. 1) is much higher than the current state-of-the-art for

unsupervised taggers, whether disambiguating or not.

To further gain insight into how successful current models are at disambiguating when they have the power to do so, we examined a collection of HMM-VB runs (Gao and Johnson 2008) and asked how the accuracy scores would change if, after training was completed, the model were forced to assign the same label to all tokens of the same type. To answer this question, we determined, for each word type, the *modal* HMM state, *i.e.*, the state most frequently assigned by the HMM to tokens of that type. We then re-labeled all words with their modal label. The effect of thus eliminating the disambiguation capacity of the model was to slightly *increase* the tagging accuracy under the best M-to-1 map for every HMM-VB run (the average increase was 0.026 for PTB17, and 0.015 for PTB45). We view this as a further indication that, in the current state of the art and with regards to tagging accuracy, limiting oneself to non-disambiguating models may not adversely affect performance.

To the contrary, this limitation may actually benefit an approach such as SVD2. Indeed, on difficult learning tasks, simpler models often behave better than more powerful ones (Geman et al. 1992). HMMs are powerful since they can, in theory, induce both a system of tags and a system of contextual patterns that allow them to disambiguate word types in terms of these tags. However, carrying out both of these unsupervised learning tasks *at once* is problematic in view of the very large number of parameters to be estimated compared to the size of the training data set.

The POS-tagging subtask of disambiguation may then be construed as a challenge in its own right: demonstrate *effective* disambiguation in an unsupervised model. Specifically, show that tagging accuracy *decreases* when the model's disambiguation capacity is removed, by re-labeling all tokens with their modal label, defined above.

We believe that the SVD2 algorithm presented here could provide a launching pad for an approach that would successfully address the disambiguation challenge. It would do so by allowing a gradual and carefully controlled amount of ambiguity into an initially non-disambiguating model. This is left for future work.

Fine-grained labeling. An important feature of the SVD2 algorithm is its ability to produce a fine-grained labeling of the data, using a number of clusters much larger than the number of tags

in a syntax-motivated POS-tag system. Such fine-grained labelings can capture additional linguistic features. To achieve a fine-grained labeling, only the final clustering step in the SVD2 algorithm needs to be changed; the computational cost this entails is negligible. A high-quality fine-grained labeling, such as achieved by the SVD2 approach, may be of practical interest as an input to various types of unsupervised grammar-induction algorithms (Headden et al. 2008). This application is left for future work.

Prototype-based tagging. One potentially important practical application of a high-quality fine-grained labeling is its use for languages which lack any kind of annotated data. By first applying the SVD2 algorithm, word types are grouped together into a few hundred clusters. Then, a prototype word is automatically extracted from each cluster. This produces, in a completely unsupervised way, a list of only a few hundred words that need to be hand-tagged by an expert. The results shown in Fig. 1 indicate that these prototype tags can then be used to tag the entire corpus with only a minor decrease in accuracy compared to the best M-to-1 map—the construction of which requires a fully annotated corpus. Fig. 1 also indicates that, with only a few hundred prototypes, the gap left between the accuracy thus achieved and the upper bound for non-disambiguating models is fairly small.

References

- Alexander Clark. 2000. Inducing syntactic categories by context distribution clustering. In *The Fourth Conference on Natural Language Learning*.
- Alexander Clark. 2003. Combining distributional and morphological information for part of speech induction. In *10th Conference of the European Chapter of the Association for Computational Linguistics*, pages 59–66.
- Michael Collins. 2002. Discriminative training methods for hidden markov models: Theory and experiments with perceptron algorithms. In *Proceedings of the ACL-02 conference on Empirical methods in natural language processing – Volume 10*.
- Jianfeng Gao and Mark Johnson. 2008. A comparison of bayesian estimators for unsupervised Hidden Markov Model POS taggers. In *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*, pages 344–352.
- Stuart Geman, Elie Bienenstock and René Doursat. 1992. Neural Networks and the Bias/Variance Dilemma. *Neural Computation*, 4 (1), pages 1–58.
- Sharon Goldwater and Tom Griffiths. 2007. A fully Bayesian approach to unsupervised part-of-speech tagging. In *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, pages 744–751.
- João V. Graça, Kuzman Ganchev, Ben Taskar, and Fernando Pereira. 2009. Posterior vs. Parameter Sparsity in Latent Variable Models. In *Neural Information Processing Systems Conference (NIPS)*.
- Aria Haghighi and Dan Klein. 2006. Prototype-driven learning for sequence models. In *Proceedings of the Human Language Technology Conference of the NAACL, Main Conference*, pages 320–327, New York City, USA, June. Association for Computational Linguistics.
- William P. Headden, David McClosky, and Eugene Charniak. 2008. Evaluating unsupervised part-of-speech tagging for grammar induction. In *Proceedings of the International Conference on Computational Linguistics (COLING '08)*.
- Mark Johnson. 2007. Why doesn't EM find good HMM POS-taggers? In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 296–305.
- Marina Meilă. 2003. Comparing clusterings by the variation of information. In Bernhard Schölkopf and Manfred K. Warmuth, editors, *COLT 2003: The Sixteenth Annual Conference on Learning Theory*, volume 2777 of *Lecture Notes in Computer Science*, pages 173–187. Springer.
- Roi Reichart and Ari Rappoport. 2009. The NVI Clustering Evaluation Measure. In *Proceedings of the Thirteenth Conference on Computational Natural Language Learning (CoNLL)*, pages 165–173.
- Hinrich Schütze. 1995. Distributional part-of-speech tagging. In *Proceedings of the seventh conference on European chapter of the Association for Computational Linguistics*, pages 141–148.
- Noah A. Smith and Jason Eisner. 2005. Contrastive estimation: Training log-linear models on unlabeled data. In *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL '05)*, pages 354–362.
- Kristina Toutanova, Dan Klein, Christopher D. Manning and Yoram Singer. 2003. Feature-rich part-of-speech tagging with a cyclic dependency network. In *Proceedings of HLT-NAACL 2003*, pages 252–259.

Intelligent Selection of Language Model Training Data

Robert C. Moore William Lewis

Microsoft Research

Redmond, WA 98052, USA

{bobmoore,wilewis}@microsoft.com

Abstract

We address the problem of selecting non-domain-specific language model training data to build auxiliary language models for use in tasks such as machine translation. Our approach is based on comparing the cross-entropy, according to domain-specific and non-domain-specific language models, for each sentence of the text source used to produce the latter language model. We show that this produces better language models, trained on less data, than both random data selection and two other previously proposed methods.

1 Introduction

Statistical N-gram language models are widely used in applications that produce natural-language text as output, particularly speech recognition and machine translation. It seems to be a universal truth that output quality can always be improved by using more language model training data, but only if the training data is reasonably well-matched to the desired output. This presents a problem, because in virtually any particular application the amount of in-domain data is limited.

Thus it has become standard practice to combine in-domain data with other data, either by combining N-gram counts from in-domain and other data (usually weighting the counts in some way), or building separate language models from different data sources, interpolating the language model probabilities either linearly or log-linearly. Log-linear interpolation is particularly popular in statistical machine translation (e.g., Brants et al., 2007), because the interpolation weights can easily be discriminatively trained to optimize an end-to-end translation objective function (such as BLEU) by making the log probability according to each language model a separate feature function in the overall translation model.

The normal practice when using multiple languages models in machine translation seems to be to train models on as much data as feasible from each source, and to depend on feature weight optimization to down-weight the impact of data that is less well-matched to the translation application. In this paper, however, we show that for a data source that is not entirely in-domain, we can improve the match between the language model from that data source and the desired application output by intelligently selecting a subset of the available data as language model training data. This not only produces a language model better matched to the domain of interest (as measured in terms of perplexity on held-out in-domain data), but it reduces the computational resources needed to exploit a large amount of non-domain-specific data, since the resources needed to filter a large amount of data are much less (especially in terms of memory) than those required to build a language model from all the data.

2 Approaches to the Problem

Our approach to the problem assumes that we have enough in-domain data to train a reasonable in-domain language model, which we then use to help score text segments from other data sources, and we select segments based on a score cutoff optimized on held-out in-domain data.

We are aware of two comparable previous approaches. Lin et al. (1997) and Gao et al. (2002) both used a method similar to ours, in which the metric used to score text segments is their perplexity according to the in-domain language model. The candidate text segments with perplexity less than some threshold are selected.

The second previous approach does not explicitly make use of an in-domain language model, but is still applicable to our scenario. Klakow (2000) estimates a unigram language model from the entire non-domain-specific corpus to be selected

from, and scores each candidate text segment from that corpus by the change in the log likelihood of the in-domain data according to the unigram model, if that segment were removed from the corpus used to estimate the unigram model. Those segments whose removal would decrease the log likelihood of the in-domain data more than some threshold are selected.

Our method is a fairly simple variant of scoring by perplexity according to an in-domain language model. First, note that selecting segments based on a perplexity threshold is equivalent to selecting based on a cross-entropy threshold. Perplexity and cross-entropy are monotonically related, since the perplexity of a string s according to a model M is simply $b^{H_M(s)}$, where $H_M(s)$ is the cross-entropy of s according to M and b is the base with respect to which the cross-entropy is measured (e.g., bits or nats). However, instead of scoring text segments by perplexity or cross-entropy according to the in-domain language model, we score them by the difference of the cross-entropy of a text segment according to the in-domain language model and the cross-entropy of the text segment according to a language model trained on a random sample of the data source from which the text segment is drawn.

To state this formally, let I be an in-domain data set and N be a non-domain-specific (or otherwise not entirely in-domain) data set. Let $H_I(s)$ be the per-word cross-entropy, according to a language model trained on I , of a text segment s drawn from N . Let $H_N(s)$ be the per-word cross-entropy of s according to a language model trained on a random sample of N . We partition N into text segments (e.g., sentences), and score the segments according to $H_I(s) - H_N(s)$, selecting all text segments whose score is less than a threshold T .

This method can be justified by reasoning similar to that used to derive methods for training binary text classifiers without labeled negative examples (Denis et al., 2002; Elkin and Noto, 2008). Let us imagine that our non-domain-specific corpus N contains an in-domain subcorpus N_I , drawn from the same distribution as our in-domain corpus I . Since N_I is statistically just like our in-domain data I , it would seem to be a good candidate for the data that we want to extract from N . By a simple variant of Bayes rule, the probability $P(N_I|s, N)$ of a text segment s , drawn randomly from N , being in N_I is given by

$$P(N_I|s, N) = \frac{P(s|N_I, N)P(N_I|N)}{P(s|N)}$$

Since N_I is a subset of N , $P(s|N_I, N) = P(s|N_I)$, and by our assumption about the relationship of I and N_I , $P(s|N_I) = P(s|I)$. Hence,

$$P(N_I|s, N) = \frac{P(s|I)P(N_I|N)}{P(s|N)}$$

If we could estimate all the probabilities in the right-hand side of this equation, we could use it to select text segments that have a high probability of being in N_I .

We can estimate $P(s|I)$ and $P(s|N)$ by training language models on I and a sample of N , respectively. That leaves us only $P(N_I|N)$, to estimate, but we really don't care what $P(N_I|N)$ is, because knowing that would still leave us wondering what threshold to set on $P(N_I|s, N)$. We don't care about classification accuracy; we care only about the quality of the resulting language model, so we might as well just attempt to find a threshold on $P(s|I)/P(s|N)$ that optimizes the fit of the resulting language model to held-out in-domain data.

Equivalently, we can work in the log domain with the quantity $\log(P(s|I)) - \log(P(s|N))$. This gets us very close to working with the difference in cross-entropies, because $H_I(s) - H_N(s)$ is just a length-normalized version of $\log(P(s|I)) - \log(P(s|N))$, with the sign reversed. The reason that we need to normalize for length is that the value of $\log(P(s|I)) - \log(P(s|N))$ tends to correlate very strongly with text segment length. If the candidate text segments vary greatly in length—e.g., if we partition N into sentences—this correlation can be a serious problem.

We estimated this effect on a 1000-sentence sample of our experimental data described below, and found the correlation between sentence log probability difference and sentence length to be $r = -0.92$, while the cross-entropy difference was almost uncorrelated with sentence length ($r = 0.04$). Hence, using sentence probability ratios or log probability differences as our scoring function would result in selecting disproportionately very short sentences. We tested this in an experiment not described here in detail, and found it not to be significantly better as a selection criterion than random selection.

Corpus	Sentence count	Token count
Gigaword	133,310,562	3,445,946,266
Europarl train	1,651,392	48,230,859
Europarl test	2,000	55,566

Table 1: Corpus size statistics

3 Experiments

We have empirically evaluated our proposed method for selecting data from a non-domain-specific source to model text in a specific domain. For the in-domain corpus, we chose the English side of the English-French parallel text from release v5 of the Europarl corpus (Koehn, 2005). This consists of proceedings of the European Parliament from 1999 through 2009. We used the text from 1999 through 2008 as in-domain training data, and we used the first 2000 sentences from January 2009 as test data. For the non-domain-specific corpus, we used the LDC English Gigaword Third Edition (LDC Catalog No.: LDC2007T07).

We used a simple tokenization scheme on all data, splitting on white space and on boundaries between alphanumeric and nonalphanumeric (e.g., punctuation) characters. With this tokenization, the sizes of our data sets in terms of sentences and tokens are shown in Table 1. The token counts include added end-of-sentence tokens.

To implement our data selection method we required one language model trained on the Europarl training data and one trained on the Gigaword data. To make these language models comparable, and to show the feasibility of optimizing the fit to the in-domain data without training a model on the entire Gigaword corpus, we trained the Gigaword language model for data selection on a random sample of the Gigaword corpus of a similar size to that of the Europarl training data: 1,874,051 sentences, 48,459,945 tokens.

To further increase the comparability of these Europarl and Gigaword language models, we restricted the vocabulary of both models to the tokens appearing at least twice in the Europarl training data, treating all other tokens as instances of $\langle \text{UNK} \rangle$. With this vocabulary, 4-gram language models were trained on both the Europarl training data and the Gigaword random sample using back-off absolute discounting (Ney et al. 1994), with a discount of 0.7 used for all N-gram lengths. The

discounted probability mass at the unigram level was added to the probability of $\langle \text{UNK} \rangle$. A count cutoff of 2 occurrences was applied to the trigrams and 4-grams in estimating these models.

We computed the cross-entropy of each sentence in the Gigaword corpus according to both models, and scored each sentence by the difference in cross-entropy, $H_{Ep}(s) - H_{Gw}(s)$. We then selected subsets of the Gigaword data corresponding to 8 cutoff points in the cross-entropy difference scores, and trained 4-gram models (again using absolute discounting with a discount of 0.7) on each of these subsets and on the full Gigaword corpus. These language models were estimated without restricting the vocabulary or applying count cutoffs, but the only parameters computed were those needed to determine the perplexity of the held-out Europarl test set, which saves a substantial amount of computation in determining the optimal selection threshold.

We compared our selection method to three other methods. As a baseline, we trained language models on random subsets of the Gigaword corpus of approximately equal size to the data sets produced by the cutoffs we selected for the cross-entropy difference scores. Next, we scored all the Gigaword sentences by the cross-entropy according to the Europarl-trained model alone. As we noted above, this is equivalent to the in-domain perplexity scoring method used by Lin et al. (1997) and Gao et al. (2002). Finally, we implemented Klakow’s (2000) method, scoring each Gigaword sentence by removing it from the Gigaword corpus and computing the difference in the log likelihood of the Europarl corpus according to unigram models trained on the Gigaword corpus with and without that sentence. With the latter two methods, we chose cutoff points in the resulting scores to produce data sets approximately equal in size to those obtained using our selection method.

4 Results

For all four selection methods, plots of test set perplexity vs. the number of training data tokens selected are displayed in Figure 1. (Note that the training data token counts are displayed on a logarithmic scale.) The test set perplexity for the language model trained on the full Gigaword corpus is 135. As we might expect, reducing training data by random sampling always increases perplexity. Selecting Gigaword sentences by their

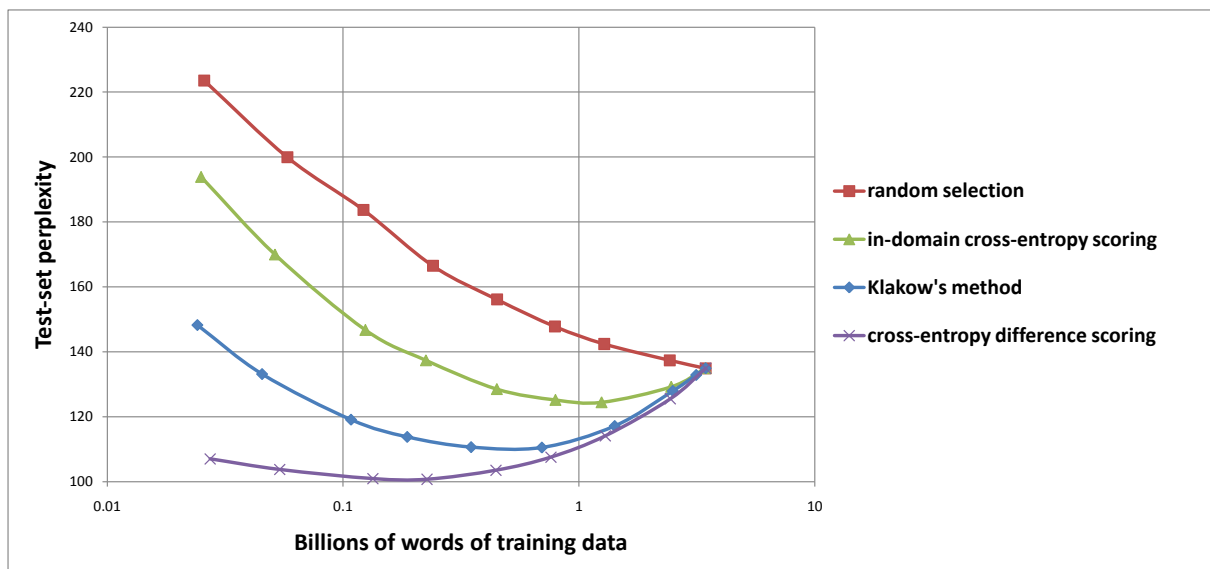


Figure 1: Test set perplexity vs. training set size

Selection Method	Original LM PPL	Modified LM PPL
in-domain cross-entropy scoring	124.4	124.8
Klakow's method	110.5	110.8
cross-entropy difference scoring	100.7	101.9

Table 2: Results adjusted for vocabulary coverage

cross-entropy according to the Europarl-trained model is effective in reducing both test set perplexity and training corpus size, with an optimum perplexity of 124, obtained with a model built from 36% of the Gigaword corpus. Klakow's method is even more effective, with an optimum perplexity of 111, obtained with a model built from 21% of the Gigaword corpus. The cross-entropy difference selection method, however, is yet more effective, with an optimum perplexity of 101, obtained with a model built from less than 7% of the Gigaword corpus.

The comparisons implied by Figure 1, however, are only approximate, because each perplexity (even along the same curve) is computed with respect to a different vocabulary, resulting in a different out-of-vocabulary (OOV) rate. OOV tokens in the test data are excluded from the perplexity computation, so the perplexity measurements are not strictly comparable.

Out of the 55566 test set tokens, the number of OOV tokens ranges from 418 (0.75%), for the smallest training set based on in-domain cross-entropy scoring, to 20 (0.03%), for training on the full Gigaword corpus. If we consider only

the training sets that appear to produce the lowest perplexity for each selection method, however, the spread of OOV counts is much narrower, ranging 53 (0.10%) for best training set based on cross-entropy difference scoring, to 20 (0.03%), for random selection.

To control for the difference in vocabulary, we estimated a modified 4-gram language model for each selection method (other than random selection) using the training set that appeared to produce the lowest perplexity for that selection method in our initial experiments. In the modified language models, the unigram model based on the selected training set is smoothed by absolute discounting, and backed-off to an unsmoothed unigram model based on the full Gigaword corpus. This produces language models that are normalized over the same vocabulary as a model trained on the full Gigaword corpus; thus the test set has the same OOVs for each model.

Test set perplexity for each of these modified language models is compared to that of the original version of the model in Table 2. It can be seen that adjusting the vocabulary in this way, so that all models are based on the same vocabulary,

yields only very small changes in the measured test-set perplexity, and these differences are much smaller than the differences between the different selection methods, whichever way the vocabulary of the language models is determined.

5 Conclusions

The cross-entropy difference selection method introduced here seems to produce language models that are both a better match to texts in a restricted domain, and require less data for training, than any of the other data selection methods tested. This study is preliminary, however, in that we have not yet shown improved end-to-end task performance applying this approach, such as improved BLEU scores in a machine translation task. However, we believe there is reason to be optimistic about this. When a language model trained on non-domain-specific data is used in a statistical translation model as a separate feature function (as is often the case), lower perplexity on in-domain target language test data derived from reference translations corresponds directly to assigning higher language model feature scores to those reference translations, which should in turn lead to translation system output that matches reference translations better.

References

- Thorsten Brants, Ashok C. Popat, Peng Xu, Franz J. Och, and Jeffrey Dean. 2007. Large language models in machine translation. In *Proceedings of the Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, June 28–30, Prague, Czech Republic, 858–867.
- François Denis, Remi Gilleron, and Marc Tommasi. 2002. Text classification from positive and unlabeled examples. In *The 9th International Conference on Information Processing and Management of Uncertainty in Knowledge-Based Systems (IPMU 2002)*, 1927–1934.
- Charles Elkin and Keith Noto. 2008. Learning classifiers from only positive and unlabeled data. In *KDD 2008*, August 24–27, Las Vegas, Nevada, USA, 213–220.
- Jianfeng Gao, Joshua Goodman, Mingjing Li, and Kai-Fu Lee. 2002. Toward a unified approach to statistical language modeling for Chinese. *ACM Transactions on Asian Language Information Processing*, 1(1):3–33.
- Dietrich Klakow. 2000. Selecting articles from the language model training corpus. In *ICASSP 2000*, June 5–9, Istanbul, Turkey, vol. 3, 1695–1698.
- Philipp Koehn. 2005. Europarl: a parallel corpus for statistical machine translation. In *MT Summit X*, September 12–16, Phuket, Thailand, 79–86.
- Sung-Chien Lin, Chi-Lung Tsai, Lee-Feng Chien, Ker-Jiann Chen, and Lin-Shan Lee. 1997. Chinese language model adaptation based on document classification and multiple domain-specific language models. In *EUROSPEECH-1997*, 1463–1466.
- Hermann Ney, Ute Essen, and Reinhard Kneser. 1994. On structuring dependencies in stochastic language modelling. *Computer Speech and Language*, 8:1–38.

Blocked Inference in Bayesian Tree Substitution Grammars

Trevor Cohn

Department of Computer Science
University of Sheffield
T.Cohn@dcs.shef.ac.uk

Phil Blunsom

Computing Laboratory
University of Oxford
Phil.Blunsom@comlab.ox.ac.uk

Abstract

Learning a tree substitution grammar is very challenging due to derivational ambiguity. Our recent approach used a Bayesian non-parametric model to induce good derivations from treebanked input (Cohn et al., 2009), biasing towards small grammars composed of small generalisable productions. In this paper we present a novel training method for the model using a blocked Metropolis-Hastings sampler in place of the previous method's local Gibbs sampler. The blocked sampler makes considerably larger moves than the local sampler and consequently converges in less time. A core component of the algorithm is a grammar transformation which represents an infinite tree substitution grammar in a finite context free grammar. This enables efficient blocked inference for training and also improves the parsing algorithm. Both algorithms are shown to improve parsing accuracy.

1 Introduction

Tree Substitution Grammar (TSG) is a compelling grammar formalism which allows nonterminal rewrites in the form of trees, thereby enabling the modelling of complex linguistic phenomena such as argument frames, lexical agreement and idiomatic phrases. A fundamental problem with TSGs is that they are difficult to estimate, even in the supervised scenario where treebanked data is available. This is because treebanks are typically not annotated with their TSG derivations (how to decompose a tree into elementary tree fragments); instead the derivation needs to be inferred.

In recent work we proposed a TSG model which infers an optimal decomposition under a non-parametric Bayesian prior (Cohn et al., 2009).

This used a Gibbs sampler for training, which repeatedly samples for every node in every training tree a binary value indicating whether the node is or is not a substitution point in the tree's derivation. Aggregated over the whole corpus, these values and the underlying trees specify the weighted grammar. Local Gibbs samplers, although conceptually simple, suffer from slow convergence (a.k.a. poor mixing). The sampler can get easily stuck because many locally improbable decisions are required to escape from a locally optimal solution. This problem manifests itself both locally to a sentence and globally over the training sample. The net result is a sampler that is non-convergent, overly dependent on its initialisation and cannot be said to be sampling from the posterior.

In this paper we present a blocked Metropolis-Hastings sampler for learning a TSG, similar to Johnson et al. (2007). The sampler jointly updates all the substitution variables in a tree, making much larger moves than the local single-variable sampler. A critical issue when developing a Metropolis-Hastings sampler is choosing a suitable proposal distribution, which must have the same support as the true distribution. For our model the natural proposal distribution is a MAP point estimate, however this cannot be represented directly as it is infinitely large. To solve this problem we develop a grammar transformation which can succinctly represent an infinite TSG in an equivalent finite Context Free Grammar (CFG). The transformed grammar can be used as a proposal distribution, from which samples can be drawn in polynomial time. Empirically, the blocked sampler converges in fewer iterations and in less time than the local Gibbs sampler. In addition, we also show how the transformed grammar can be used for parsing, which yields theoretical and empirical improvements over our previous method which truncated the grammar.

2 Background

A *Tree Substitution Grammar* (TSG; Bod et al. (2003)) is a 4-tuple, $G = (T, N, S, R)$, where T is a set of *terminal symbols*, N is a set of *non-terminal symbols*, $S \in N$ is the distinguished *root nonterminal* and R is a set of productions (rules). The productions take the form of tree fragments, called *elementary trees* (ETs), in which each internal node is labelled with a nonterminal and each leaf is labelled with either a terminal or a nonterminal. The *frontier nonterminal* nodes in each ET form the sites into which other ETs can be substituted. A *derivation* creates a tree by recursive substitution starting with the root symbol and finishing when there are no remaining frontier nonterminals. Figure 1 (left) shows an example derivation where the arrows denote substitution. A *Probabilistic Tree Substitution Grammar* (PTSG) assigns a probability to each rule in the grammar, where each production is assumed to be conditionally independent given its root nonterminal. A derivation’s probability is the product of the probabilities of the rules therein.

In this work we employ the same non-parametric TSG model as Cohn et al. (2009), which we now summarise. The inference problem within this model is to identify the posterior distribution of the elementary trees \mathbf{e} given whole trees \mathbf{t} . The model is characterised by the use of a Dirichlet Process (DP) prior over the grammar. We define the distribution over elementary trees e with root nonterminal symbol c as

$$G_c | \alpha_c, P_0 \sim \text{DP}(\alpha_c, P_0(\cdot | c))$$

$$e | c \sim G_c$$

where $P_0(\cdot | c)$ (the *base distribution*) is a distribution over the infinite space of trees rooted with c , and α_c (the *concentration parameter*) controls the model’s tendency towards either reusing elementary trees or creating novel ones as each training instance is encountered.

Rather than representing the distribution G_c explicitly, we integrate over all possible values of G_c . The key result required for inference is that the conditional distribution of e_i , given $\mathbf{e}_{-i} = e_1 \dots e_n \setminus e_i$ and the root category c is:

$$p(e_i | \mathbf{e}_{-i}, c, \alpha_c, P_0) = \frac{n_{e_i, c}^{-i}}{n_{\cdot, c}^{-i} + \alpha_c} + \frac{\alpha_c P_0(e_i | c)}{n_{\cdot, c}^{-i} + \alpha_c} \quad (1)$$

where $n_{e_i, c}^{-i}$ is the number number of times e_i has been used to rewrite c in \mathbf{e}_{-i} , and $n_{\cdot, c}^{-i} = \sum_e n_{e, c}^{-i}$

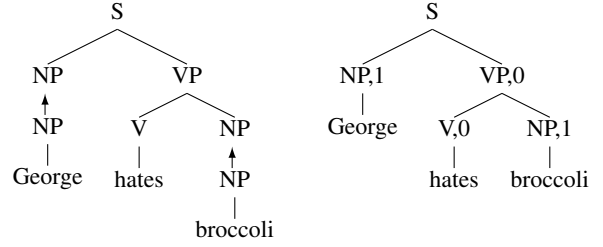


Figure 1: TSG derivation and its corresponding Gibbs state for the local sampler, where each node is marked with a binary variable denoting whether it is a substitution site.

is the total count of rewriting c . Henceforth we omit the $-i$ sub-/super-script for brevity.

A primary consideration is the definition of P_0 . Each e_i can be generated in one of two ways: by drawing from the base distribution, where the probability of any particular tree is proportional to $\alpha_c P_0(e_i | c)$, or by drawing from a cache of previous expansions of c , where the probability of any particular expansion is proportional to the number of times that expansion has been used before. In Cohn et al. (2009) we presented base distributions that favour small elementary trees which we expect will generalise well to unseen data. In this work we show that if P_0 is chosen such that it decomposes with the CFG rules contained within each elementary tree,¹ then we can use a novel dynamic programming algorithm to sample derivations without ever enumerating all the elementary trees in the grammar.

The model was trained using a local Gibbs sampler (Geman and Geman, 1984), a Markov chain Monte Carlo (MCMC) method in which random variables are repeatedly sampled conditioned on the values of all other random variables in the model. To formulate the local sampler, we associate a binary variable with each non-root internal node of each tree in the training set, indicating whether that node is a substitution point or not (illustrated in Figure 1). The sampler then visits each node in a random schedule and resamples that node’s substitution variable, where the probability of the two different configurations are given by (1). Parsing was performed using a Metropolis-Hastings sampler to draw derivation samples for a string, from which the best tree was recovered. However the sampler used for parsing was biased

¹Both choices of base distribution in Cohn et al. (2009) decompose into CFG rules. In this paper we focus on the better performing one, P_0^C , which combines a PCFG applied recursively with a stopping probability, s , at each node.

because it used as its proposal distribution a truncated grammar which excluded all but a handful of the unseen elementary trees. Consequently the proposal had smaller support than the true model, voiding the MCMC convergence proofs.

3 Grammar Transformation

We now present a blocked sampler using the Metropolis-Hastings (MH) algorithm to perform sentence-level inference, based on the work of Johnson et al. (2007) who presented a MH sampler for a Bayesian PCFG. This approach repeats the following steps for each sentence in the training set: 1) run the inside algorithm (Lari and Young, 1990) to calculate marginal expansion probabilities under a MAP approximation, 2) sample an analysis top-down and 3) accept or reject using a Metropolis-Hastings (MH) test to correct for differences between the MAP proposal and the true model. Though our model is similar to Johnson et al. (2007)'s, we have an added complication: the MAP grammar cannot be estimated directly. This is a consequence of the base distribution having infinite support (assigning non-zero probability to infinitely many unseen tree fragments), which means the MAP has an infinite rule set. For example, if our base distribution licences the CFG production $\text{NP} \rightarrow \text{NP PP}$ then our TSG grammar will contain the infinite set of elementary trees $\text{NP} \rightarrow \text{NP PP}$, $\text{NP} \rightarrow (\text{NP NP PP}) \text{PP}$, $\text{NP} \rightarrow (\text{NP (NP NP PP) PP}) \text{PP}$, ... with decreasing but non-zero probability.

However, we can represent the infinite MAP using a grammar transformation inspired by Goodman (2003), which represents the MAP TSG in an equivalent finite PCFG.² Under the transformed PCFG inference is efficient, allowing its use as the proposal distribution in a blocked MH sampler. We represent the MAP using the grammar transformation in Table 1 which separates the $n_{e,c}$ and P_0 terms in (1) into two separate CFGs, A and B. Grammar A has productions for every ET with $n_{e,c} \geq 1$ which are assigned unsmoothed probabilities: omitting the P_0 term from (1).³ Grammar B has productions for every CFG production licensed under P_0 ; its productions are denoted using

²Backoff DOP uses a similar packed representation to encode the set of smaller subtrees for a given elementary tree (Sima'an and Buratto, 2003), which are used to smooth its probability estimate.

³The transform assumes inside inference. For Viterbi replace the probability for $c \rightarrow \text{sign}(e)$ with $\frac{n_{e,c} + \alpha_c P_0(e|c)}{n_{\cdot,c} + \alpha_c}$.

For every ET, e , rewriting c with non-zero count:	
$c \rightarrow \text{sign}(e)$	$\frac{n_{e,c}}{n_{\cdot,c} + \alpha_c}$
For every internal node e_i in e with children $e_{i,1}, \dots, e_{i,n}$	
$\text{sign}(e_i) \rightarrow \text{sign}(e_{i,1}) \dots \text{sign}(e_{i,n})$	1
For every nonterminal, c :	
$c \rightarrow c'$	$\frac{\alpha_c}{n_{\cdot,c} + \alpha_c}$
For every pre-terminal CFG production, $c \rightarrow t$:	
$c' \rightarrow t$	$P_{CFG}(c \rightarrow t)$
For every unary CFG production, $c \rightarrow a$:	
$c' \rightarrow a$	$P_{CFG}(c \rightarrow a)s_a$
$c' \rightarrow a'$	$P_{CFG}(c \rightarrow a)(1 - s_a)$
For every binary CFG production, $c \rightarrow ab$:	
$c' \rightarrow ab$	$P_{CFG}(c \rightarrow ab)s_a s_b$
$c' \rightarrow ab'$	$P_{CFG}(c \rightarrow ab)s_a(1 - s_b)$
$c' \rightarrow a'b$	$P_{CFG}(c \rightarrow ab)(1 - s_a)s_b$
$c' \rightarrow a'b'$	$P_{CFG}(c \rightarrow ab)(1 - s_a)(1 - s_b)$

Table 1: Grammar transformation rules to map a MAP TSG into a CFG. Production probabilities are shown to the right of each rule. The $\text{sign}(e)$ function creates a unique string signature for an ET e (where the signature of a frontier node is itself) and s_c is the Bernoulli probability of c being a substitution variable (and stopping the P_0 recursion).

primed ($'$) nonterminals. The rule $c \rightarrow c'$ bridges from A to B, weighted by the smoothing term excluding P_0 , which is computed recursively via child productions. The remaining rules in grammar B correspond to every CFG production in the underlying PCFG base distribution, coupled with the binary decision whether or not nonterminal children should be substitution sites (frontier nonterminals). This choice affects the rule probability by including a s or $1 - s$ factor, and child substitution sites also function as a bridge back from grammar B to A. In this way there are often two equivalent paths to reach the same chart cell using the same elementary tree – via grammar A using observed TSG productions and via grammar B using P_0 backoff; summing these yields the desired net probability.

Figure 2 shows an example of the transformation of an elementary tree with non-zero count, $n_{e,c} \geq 1$, into the two types of CFG rules. Both parts are capable of parsing the string NP, saw, NP into a S, as illustrated in Figure 3; summing the probability of both analyses gives the model probability from (1). Note that although the probabilities exactly match the true model for a single elementary tree, the probability of derivations composed of many elementary trees may not match because the model's caching behaviour has been suppressed, i.e., the counts, n , are not incremented during the course of a derivation.

For training we define the MH sampler as follows. First we estimate the MAP grammar over

$S \rightarrow NP VP\{V\{saw\},NP\}$	$\frac{n_{e,S}^-}{n_{e,S}^- + \alpha_S}$
$VP\{V\{saw\},NP\} \rightarrow V\{saw\} NP$	1
$V\{saw\} \rightarrow saw$	1
<hr/>	
$S \rightarrow S'$	$\frac{\alpha_S}{n_{e,S}^- + \alpha_S}$
$S' \rightarrow NP VP'$	$P_{CFG}(S \rightarrow NP VP)_{SNP}(1 - s_{VP})$
$VP' \rightarrow V' NP$	$P_{CFG}(VP \rightarrow V NP)(1 - s_V)_{SNP}$
$V' \rightarrow saw$	$P_{CFG}(V \rightarrow saw)$

Figure 2: Example of the transformed grammar for the ET (S NP (VP (V saw) NP)). Taking the product of the rule scores above the line yields the left term in (1), and the product of the scores below the line yields the right term.

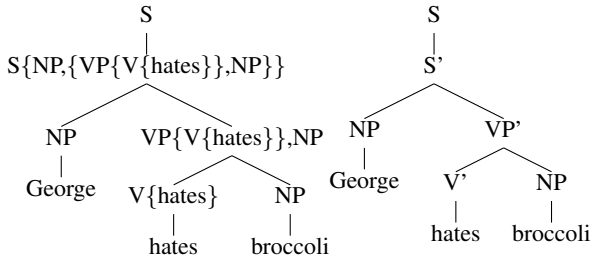


Figure 3: Example trees under the grammar transform, which both encode the same TSG derivation from Figure 1. The left tree encodes that the $S \rightarrow NP (VP (V hates) NP)$ elementary tree was drawn from the cache, while for the right tree this same elementary tree was drawn from the base distribution (the left and right terms in (1), respectively).

the derivations of training corpus excluding the current tree, which we represent using the PCFG transformation. The next step is to sample derivations for a given tree, for which we use a constrained variant of the inside algorithm (Lari and Young, 1990). We must ensure that the TSG derivation produces the given tree, and therefore during inside inference we only consider spans that are constituents in the tree and are labelled with the correct nonterminal. Nonterminals are said to match their primed and signed counterparts, e.g., NP' and $NP\{DT, NN\{car\}\}$ both match NP . Under the tree constraints the time complexity of inside inference is linear in the length of the sentence. A derivation is then sampled from the inside chart using a top-down traversal (Johnson et al., 2007), and converted back into its equivalent TSG derivation. The derivation is scored with the true model and accepted or rejected using the MH test; accepted samples then replace the current derivation for the tree, and rejected samples leave the previous derivation unchanged. These steps are then repeated for another tree in the training set, and the process is then repeated over the full training set many times.

Parsing The grammar transform is not only useful for training, but also for parsing. To parse a sentence we sample a number of TSG derivations from the MAP which are then accepted or rejected into the full model using a MH step. The samples are obtained from the same transformed grammar but adapting the algorithm for an unsupervised setting where parse trees are not available. For this we use the standard inside algorithm applied to the sentence, omitting the tree constraints, which has time complexity cubic in the length of the sentence. We then sample a derivation from the inside chart and perform the MH acceptance test. This setup is theoretically more appealing than our previous approach in which we truncated the approximation grammar to exclude most of the zero count rules (Cohn et al., 2009). We found that both the maximum probability derivation and tree were considerably worse than a tree constructed to maximise the expected number of correct CFG rules (MER), based on Goodman’s (2003) algorithm for maximising labelled recall. For this reason we use the MER parsing algorithm using sampled Monte Carlo estimates for the marginals over CFG rules at each sentence span.

4 Experiments

We tested our model on the Penn treebank using the same data setup as Cohn et al. (2009). Specifically, we used only section 2 for training and section 22 (devel) for reporting results. Our models were all sampled for 5k iterations with hyperparameter inference for α_c and $s_c \forall c \in N$, but in contrast to our previous approach we did not use annealing which we did not find to help generalisation accuracy. The MH acceptance rates were in excess of 99% across both training and parsing. All results are averages over three runs.

For training the blocked MH sampler exhibits faster convergence than the local Gibbs sampler, as shown in Figure 4. Irrespective of the initialisation the blocked sampler finds higher likelihood states in many fewer iterations (the same trend continues until iteration 5k). To be fair, the blocked sampler is slower per iteration (roughly 50% worse) due to the higher overheads of the grammar transform and performing dynamic programming (despite nominal optimisation).⁴ Even after accounting for the time differ-

⁴The speed difference diminishes with corpus size: on sections 2–22 the blocked sampler is only 19% slower per

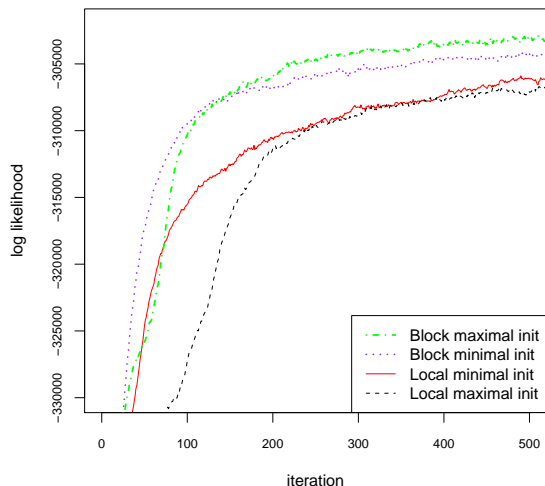


Figure 4: Training likelihood vs. iteration. Each sampling method was initialised with both minimal and maximal elementary trees.

Training	truncated	transform
Local minimal init	77.63	77.98
Local maximal init	77.19	77.71
Blocked minimal init	77.98	78.40
Blocked maximal init	77.67	78.24

Table 2: Development F1 scores using the truncated parsing algorithm and the novel grammar transform algorithm for four different training configurations.

ence the blocked sampler is more effective than the local Gibbs sampler. Training likelihood is highly correlated with generalisation F1 (Pearson’s correlation efficient of 0.95), and therefore improving the sampler convergence will have immediate effects on performance.

Parsing results are shown in Table 2.⁵ The blocked sampler results in better generalisation F1 scores than the local Gibbs sampler, irrespective of the initialisation condition or parsing method used. The use of the grammar transform in parsing also yields better scores irrespective of the underlying model. Together these results strongly advocate the use of the grammar transform for inference in infinite TSGs.

We also trained the model on the standard Penn treebank training set (sections 2–21). We initialised the model with the final sample from a run on the small training set, and used the blocked sampler for 6500 iterations. Averaged over three runs, the test F1 (section 23) was 85.3 an improvement over the local sampler.

⁵Our baseline ‘Local maximal init’ slightly exceeds previously reported score of 76.89% (Cohn et al., 2009).

ment over our earlier 84.0 (Cohn et al., 2009) although still well below state-of-the-art parsers. We conjecture that the performance gap is due to the model using an overly simplistic treatment of unknown words, and also a further mixing problems with the sampler. For the full data set the counts are much larger in magnitude which leads to stronger modes. The sampler has difficulty escaping such modes and therefore is slower to mix. One way to solve the mixing problem is for the sampler to make more global moves, e.g., with table label resampling (Johnson and Goldwater, 2009) or split-merge (Jain and Neal, 2000). Another way is to use a variational approximation instead of MCMC sampling (Wainwright and Jordan, 2008).

5 Discussion

We have demonstrated how our grammar transformation can implicitly represent an exponential space of tree fragments efficiently, allowing us to build a sampler with considerably better mixing properties than a local Gibbs sampler. The same technique was also shown to improve the parsing algorithm. These improvements are in no way limited to our particular choice of a TSG parsing model, many hierarchical Bayesian models have been proposed which would also permit similar optimised samplers. In particular models which induce segmentations of complex structures stand to benefit from this work; Examples include the word segmentation model of Goldwater et al. (2006) for which it would be trivial to adapt our technique to develop a blocked sampler. Hierarchical Bayesian segmentation models have also become popular in statistical machine translation where there is a need to learn phrasal translation structures that can be decomposed at the word level (DeNero et al., 2008; Blunsom et al., 2009; Cohn and Blunsom, 2009). We envisage similar representations being applied to these models to improve their mixing properties.

A particularly interesting avenue for further research is to employ our blocked sampler for unsupervised grammar induction. While it is difficult to extend the local Gibbs sampler to the case where the tree is not observed, the dynamic program for our blocked sampler can be easily used for unsupervised inference by omitting the tree matching constraints.

References

- Phil Blunsom, Trevor Cohn, Chris Dyer, and Miles Osborne. 2009. A Gibbs sampler for phrasal synchronous grammar induction. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP (ACL-IJCNLP)*, pages 782–790, Suntec, Singapore, August.
- Rens Bod, Remko Scha, and Khalil Sima'an, editors. 2003. *Data-oriented parsing*. Center for the Study of Language and Information - Studies in Computational Linguistics. University of Chicago Press.
- Trevor Cohn and Phil Blunsom. 2009. A Bayesian model of syntax-directed tree to string grammar induction. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 352–361, Singapore, August.
- Trevor Cohn, Sharon Goldwater, and Phil Blunsom. 2009. Inducing compact but accurate tree-substitution grammars. In *Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics (HLT-NAACL)*, pages 548–556, Boulder, Colorado, June.
- John DeNero, Alexandre Bouchard-Côté, and Dan Klein. 2008. Sampling alignment structure under a Bayesian translation model. In *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*, pages 314–323, Honolulu, Hawaii, October.
- Stuart Geman and Donald Geman. 1984. Stochastic relaxation, Gibbs distributions and the Bayesian restoration of images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 6:721–741.
- Sharon Goldwater, Thomas L. Griffiths, and Mark Johnson. 2006. Contextual dependencies in unsupervised word segmentation. In *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics*, pages 673–680, Sydney, Australia, July.
- Joshua Goodman. 2003. Efficient parsing of DOP with PCFG-reductions. In Bod et al. (Bod et al., 2003), chapter 8.
- Sonia Jain and Radford M. Neal. 2000. A split-merge Markov chain Monte Carlo procedure for the Dirichlet process mixture model. *Journal of Computational and Graphical Statistics*, 13:158–182.
- Mark Johnson and Sharon Goldwater. 2009. Improving nonparameteric bayesian inference: experiments on unsupervised word segmentation with adaptor grammars. In *Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 317–325, Boulder, Colorado, June.
- Mark Johnson, Thomas Griffiths, and Sharon Goldwater. 2007. Bayesian inference for PCFGs via Markov chain Monte Carlo. In *Proceedings of Human Language Technologies 2007: The Conference of the North American Chapter of the Association for Computational Linguistics*, pages 139–146, Rochester, NY, April.
- Karim Lari and Steve J. Young. 1990. The estimation of stochastic context-free grammars using the inside-outside algorithm. *Computer Speech and Language*, 4:35–56.
- Khalil Sima'an and Luciano Buratto. 2003. Backoff parameter estimation for the dop model. In Nada Lavrac, Dragan Gamberger, Ljupco Todorovski, and Hendrik Blockeel, editors, *ECML*, volume 2837 of *Lecture Notes in Computer Science*, pages 373–384. Springer.
- Martin J Wainwright and Michael I Jordan. 2008. *Graphical Models, Exponential Families, and Variational Inference*. Now Publishers Inc., Hanover, MA, USA.

Online Generation of Locality Sensitive Hash Signatures

Benjamin Van Durme

HLTCOE

Johns Hopkins University
Baltimore, MD 21211 USA

Ashwin Lall

College of Computing
Georgia Institute of Technology
Atlanta, GA 30332 USA

Abstract

Motivated by the recent interest in streaming algorithms for processing large text collections, we revisit the work of Ravichandran et al. (2005) on using the Locality Sensitive Hash (LSH) method of Charikar (2002) to enable fast, approximate comparisons of vector cosine similarity. For the common case of feature updates being additive over a data stream, we show that LSH signatures can be maintained *online*, without additional approximation error, and with lower memory requirements than when using the standard *offline* technique.

1 Introduction

There has been a surge of interest in adapting results from the streaming algorithms community to problems in processing large text collections. The term *streaming* refers to a model where data is made available sequentially, and it is assumed that resource limitations preclude storing the entirety of the data for offline (batch) processing. Statistics of interest are approximated via online, randomized algorithms. Examples of text applications include: collecting approximate counts (Talbot, 2009; Van Durme and Lall, 2009a), finding top- n elements (Goyal et al., 2009), estimating term co-occurrence (Li et al., 2008), adaptive language modeling (Levenberg and Osborne, 2009), and building top- k ranklists based on pointwise mutual information (Van Durme and Lall, 2009b).

Here we revisit the work of Ravichandran et al. (2005) on building word similarity measures from large text collections by using the Locality Sensitive Hash (LSH) method of Charikar (2002). For the common case of feature updates being additive over a data stream (such as when tracking lexical co-occurrence), we show that LSH signatures can be maintained online, without additional

approximation error, and with lower memory requirements than when using the standard offline technique.

We envision this method being used in conjunction with dynamic clustering algorithms, for a variety of applications. For example, Petrovic et al. (2010) made use of LSH signatures generated over individual *tweets*, for the purpose of *first story detection*. Streaming LSH should allow for the clustering of Twitter *authors*, based on the tweets they generate, with signatures continually updated over the Twitter stream.

2 Locality Sensitive Hashing

We are concerned with computing the *cosine similarity* of feature vectors, defined for a pair of vectors \vec{u} and \vec{v} as the dot product normalized by their lengths:

$$\text{cosine-similarity}(\vec{u}, \vec{v}) = \frac{\vec{u} \cdot \vec{v}}{|\vec{u}||\vec{v}|}.$$

This similarity is the cosine of the angle between these high-dimensional vectors and attains a value of one (i.e., $\cos(0)$) when the vectors are parallel and zero (i.e., $\cos(\pi/2)$) when orthogonal.

Building on the seminal work of Indyk and Motwani (1998) on *locality sensitive hashing* (LSH), Charikar (2002) presented an LSH that maps high-dimensional vectors to a much smaller dimensional space while still preserving (cosine) similarity between vectors in the original space. The LSH algorithm computes a succinct signature of the feature set of the words in a corpus by computing d independent dot products of each feature vector \vec{v} with a random unit vector \vec{r} , i.e., $\sum_i v_i r_i$, and retaining the sign of the d resulting products. Each entry of \vec{r} is drawn from the distribution $N(0, 1)$, the normal distribution with zero mean and unit variance. Charikar's algorithm makes use of the fact (proved by Goemans and Williamson

(1995) for an unrelated application) that the angle between any two vectors summarized in this fashion is proportional to the expected Hamming distance of their signature vectors. Hence, we can retain length d bit-signatures in the place of high dimensional feature vectors, while preserving the ability to (quickly) approximate cosine similarity in the original space.

Ravichandran et al. (2005) made use of this algorithm to reduce the computation in searching for similar nouns by first computing signatures for each noun and then computing similarity over the signatures rather than the original feature space.

3 Streaming Algorithm

In this work, we focus on features that can be maintained additively, such as raw frequencies.¹ Our streaming algorithm for this problem makes use of the simple fact that the dot product of the feature vector with random vectors is a linear operation. This permits us to replace the $v_i \cdot r_i$ operation by v_i individual additions of r_i , once for each time the feature is encountered in the stream (where v_i is the frequency of a feature and r_i is the randomly chosen Gaussian-distributed value associated with this feature). The result of the final computation is identical to the dot products computed by the algorithm of Charikar (2002), but the processing can now be done online. A similar technique, for stable random projections, was independently discussed by Li et al. (2008).

Since each feature may appear multiple times in the stream, we need a consistent way to retrieve the random values drawn from $N(0, 1)$ associated with it. To avoid the expense of computing and storing these values explicitly, as is the norm, we propose the use of a precomputed pool of random values drawn from this distribution that we can then hash into. Hashing into a fixed pool ensures that the same feature will consistently be associated with the same value drawn from $N(0, 1)$. This introduces some weak dependence in the random vectors, but we will give some analysis showing that this should have very limited impact on the cosine similarity computation, which we further support with experimental evidence (see Table 3).

Our algorithm traverses a stream of words and

¹Note that Ravichandran et al. (2005) used pointwise mutual information features, which are not additive since they require a global statistic to compute.

Algorithm 1 STREAMING LSH ALGORITHM

Parameters:

m : size of pool
 d : number of bits (size of resultant signature)
 s : a random seed
 h_1, \dots, h_d : hash functions mapping $\langle s, f_i \rangle$ to $\{0, \dots, m-1\}$

INITIALIZATION:

- 1: Initialize floating point array $P[0, \dots, m-1]$
- 2: Initialize H , a hashtable mapping words to floating point arrays of size d
- 3: **for** $i := 0 \dots m-1$ **do**
- 4: $P[i] :=$ random sample from $N(0, 1)$, using s as seed

ONLINE:

- 1: **for** each word w in the stream **do**
- 2: **for** each feature f_i associated with w **do**
- 3: **for** $j := 1 \dots d$ **do**
- 4: $H[w][j] := H[w][j] + P[h_j(s, f_i)]$

SIGNATURECOMPUTATION:

- 1: **for** each $w \in H$ **do**
- 2: **for** $i := 1 \dots d$ **do**
- 3: **if** $H[w][i] > 0$ **then**
- 4: $S[w][i] := 1$
- 5: **else**
- 6: $S[w][i] := 0$

maintains some state for each possible word that it encounters (cf. Algorithm 1). In particular, the state maintained for each word is a vector of floating point numbers of length d . Each element of the vector holds the (partial) dot product of the feature vector of the word with a random unit vector. Updating the state for a feature seen in the stream for a given word simply involves incrementing each position in the word's vector by the random value associated with the feature, accessed by hash functions h_1 through h_d . At any point in the stream, the vector for each word can be processed (in time $O(d)$) to create a signature computed by checking the sign of each component of its vector.

3.1 Analysis

The update cost of the streaming algorithm, per word in the stream, is $O(df)$, where d is the target signature size and f is the number of features associated with each word in the stream.² This results in an overall cost of $O(ndf)$ for the streaming algorithm, where n is the length of the stream. The memory footprint of our algorithm is $O(n_0d+m)$, where n_0 is the number of distinct words in the stream and m is the size of the pool of normally distributed values. In comparison, the original LSH algorithm computes signatures at a cost of $O(nf + n_0dF)$ updates and $O(n_0F + dF + n_0d)$ memory, where F is the (large) number of unique

²For the bigram features used in § 4, $f = 2$.

features. Our algorithm is superior in terms of memory (because of the pooling trick), and has the benefit of supporting similarity queries online.

3.2 Pooling Normally-distributed Values

We now discuss why it is possible to use a fixed pool of random values instead of generating unique ones for each feature. Let g be the c.d.f. of the distribution $N(0, 1)$. It is easy to see that picking $x \in (0, 1)$ uniformly results in $g^{-1}(x)$ being chosen with distribution $N(0, 1)$. Now, if we select for our pool the values

$$g^{-1}(1/m), g^{-1}(2/m), \dots, g^{-1}(1 - 1/m),$$

for some sufficiently large m , then this is identical to sampling from $N(0, 1)$ with the caveat that the accuracy of the sample is limited. More precisely, the deviation from sampling from this pool is off from the actual value by at most

$$\max_{i=1, \dots, m-2} \{g^{-1}((i+1)/m) - g^{-1}(i/m)\}.$$

By choosing m to be sufficiently large, we can bound the error of the approximate sample from a true sample (i.e., the loss in precision expressed above) to be a small fraction (e.g., 1%) of the actual value. This would result in the same relative error in the computation of the dot product (i.e., 1%), which would almost never affect the sign of the final value. Hence, pooling as above should give results almost identical to the case where all the random values were chosen independently. Finally, we make the observation that, for large m , randomly choosing m values from $N(0, 1)$ results in a set of values that are distributed very similarly to the pool described above. An interesting avenue for future work is making this analysis more mathematically precise.

3.3 Extensions

Decay The algorithm can be extended to support *temporal decay* in the stream, where recent observations are given higher relative weight, by multiplying the current sums by a decay value (e.g., 0.9) on a regular interval (e.g., once an hour, once a day, once a week, etc.).

Distributed The algorithm can be easily distributed across multiple machines in order to process different parts of a stream, or multiple different streams, in parallel, such as in the context of the MapReduce framework (Dean and Ghemawat,

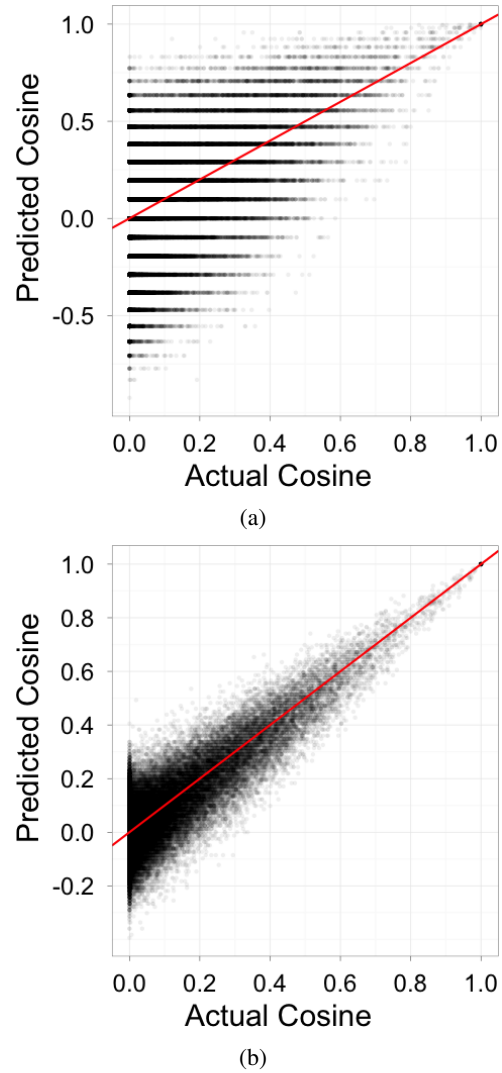


Figure 1: Predicted versus actual cosine values for 50,000 pairs, using LSH signatures generated online, with $d = 32$ in Fig. 1(a) and $d = 256$ in Fig. 1(b).

2004). The underlying operation is a linear operator that is easily composed (i.e., via addition), and the randomness between machines can be tied based on a shared seed s . At any point in processing the stream(s), current results can be aggregated by summing the d -dimensional vectors for each word, from each machine.

4 Experiments

Similar to the experiments of Ravichandran et al. (2005), we evaluated the fidelity of signature generation in the context of calculating distributional similarity between words across a large text collection: in our case, articles taken from the NYTimes portion of the Gigaword corpus (Graff, 2003). The collection was processed as a stream, sentence by sentence, using bigram fea-

d	16	32	64	128	256
SLSH	0.2885	0.2112	0.1486	0.1081	0.0769
LSH	0.2892	0.2095	0.1506	0.1083	0.0755

Table 1: Mean absolute error when using signatures generated online (StreamingLSH), compared to offline (LSH).

tures. This gave a stream of 773,185,086 tokens, with 1,138,467 unique types. Given the number of types, this led to a (sparse) feature space with dimension on the order of 2.5 million.

After compiling signatures, fifty-thousand $\langle x, y \rangle$ pairs of types were randomly sampled by selecting x and y each independently, with replacement, from those types with at least 10 tokens in the stream (where 310,327 types satisfied this constraint). The true cosine values between each such x and y was computed based on offline calculation, and compared to the cosine similarity predicted by the Hamming distance between the signatures for x and y . Unless otherwise specified, the random pool size was fixed at $m = 10,000$.

Figure 1 visually reaffirms the trade-off in LSH between the number of bits and the accuracy of cosine prediction across the range of cosine values. As the underlying vectors are strictly positive, the true cosine is restricted to $[0, 1]$. Figure 2 shows the absolute error between truth and prediction for a similar sample, measured using signatures of a variety of bit lengths. Here we see horizontal bands arising from truly orthogonal vectors leading to step-wise absolute error values tracked to Hamming distance.

Table 1 compares the online and batch LSH algorithms, giving the mean absolute error between predicted and actual cosine values, computed for the fifty-thousand element sample, using signatures of various lengths. These results confirm that we achieve the same level of accuracy with online updates as compared to the standard method.

Figure 3 shows how a pool size as low as $m = 100$ gives reasonable variation in random values, and that $m = 10,000$ is sufficient. When using a standard 32 bit floating point representation, this is just 40 KBytes of memory, as compared to, e.g., the 2.5 GBytes required to store 256 random vectors each containing 2.5 million elements.

Table 2 is based on taking an example for each of three part-of-speech categories, and reporting the resultant top-5 words as according to approximated cosine similarity. Depending on the intended application, these results indicate a range

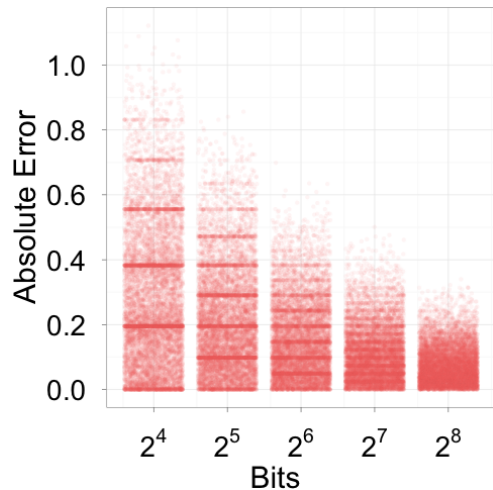


Figure 2: Absolute error between predicted and true cosine for a sample of pairs, when using signatures of length $\log_2(d) \in \{4, 5, 6, 7, 8\}$, drawn with added *jitter* to avoid overplotting.

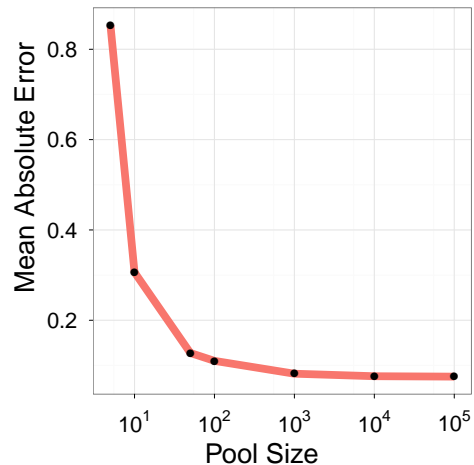


Figure 3: Error versus pool size, when using $d = 256$.

of potentially sufficient signature lengths.

5 Conclusions

We have shown that when updates to a feature vector are additive, it is possible to convert the offline LSH signature generation method into a streaming algorithm. In addition to allowing for online querying of signatures, our approach leads to space efficiencies, as it does not require the explicit representation of either the feature vectors, nor the random matrix. Possibilities for future work include the pairing of this method with algorithms for dynamic clustering, as well as exploring algorithms for different distances (e.g., L_2) and estimators (e.g., asymmetric estimators (Dong et al., 2009)).

London

Milan.₉₇, **Madrid**.₉₆, **Stockholm**.₉₆, **Manila**.₉₅, **Moscow**.₉₅
ASHER₀, Champaign₀, MANS₀, NOBLE₀, come₀
Prague₁, Vienna₁, suburban₁, synchronism₁, Copenhagen₂
Frankfurt₄, Prague₄, Tazsar₅, Brussels₆, Copenhagen₆
Prague₁₂, Stockholm₁₂, Frankfurt₁₄, Madrid₁₄, Manila₁₄
Stockholm₂₀, Milan₂₂, Madrid₂₄, Taipei₂₄, Frankfurt₂₅

in

during.₉₉, **on**.₉₈, **beneath**.₉₈, **from**.₉₈, **onto**.₉₇
Across₀, Addressing₀, Addy₀, Against₀, Allmon₀
aboard₀, mishandled₀, overlooking₀, Addressing₁, Rejecting₁
Rejecting₂, beneath₂, during₂, from₃, hamstringing₃
during₄, beneath₅, of₆, on₇, overlooking₇
during₁₀, on₁₃, beneath₁₅, of₁₇, overlooking₁₇

sold

deployed.₈₄, **presented**.₈₃, **sacrificed**.₈₂, **held**.₈₂, **installed**.₈₂
Bustino₀, Diors₀, Draining₀, Kosses₀, UNA₀
delivered₂, held₂, marks₂, seared₂, Ranked₃
delivered₅, rendered₅, presented₆, displayed₇, exhibited₇
held₁₈, rendered₁₈, presented₁₉, deployed₂₀, displayed₂₀
presented₄₁, rendered₄₂, held₄₇, leased₄₇, reopened₄₇

Table 2: Top-5 items based on true cosine (bold), then using minimal Hamming distance, given in top-down order when using signatures of length $\log_2(d) \in \{4, 5, 6, 7, 8\}$. Ties broken lexicographically. Values given as subscripts.

Acknowledgments

Thanks to Deepak Ravichandran, Miles Osborne, Sasa Petrovic, Ken Church, Glen Coppersmith, and the anonymous reviewers for their feedback. This work began while the first author was at the University of Rochester, funded by NSF grant IIS-1016735. The second author was supported in part by NSF grant CNS-0905169, funded under the American Recovery and Reinvestment Act of 2009.

References

- Moses Charikar. 2002. Similarity estimation techniques from rounding algorithms. In *Proceedings of STOC*.
- Jeffrey Dean and Sanjay Ghemawat. 2004. MapReduce: Simplified Data Processing on Large Clusters. In *Proceedings of OSDI*.
- Wei Dong, Moses Charikar, and Kai Li. 2009. Asymmetric distance estimation with sketches for similarity search in high-dimensional spaces. In *Proceedings of SIGIR*.
- Michel X. Goemans and David P. Williamson. 1995. Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming. *JACM*, 42:1115–1145.
- Amit Goyal, Hal Daumé III, and Suresh Venkatasubramanian. 2009. Streaming for large scale NLP: Language Modeling. In *Proceedings of NAACL*.
- David Graff. 2003. English Gigaword. Linguistic Data Consortium, Philadelphia.
- Piotr Indyk and Rajeev Motwani. 1998. Approximate nearest neighbors: towards removing the curse of dimensionality. In *Proceedings of STOC*.
- Abby Levenberg and Miles Osborne. 2009. Stream-based Randomised Language Models for SMT. In *Proceedings of EMNLP*.
- Ping Li, Kenneth W. Church, and Trevor J. Hastie. 2008. One Sketch For All: Theory and Application of Conditional Random Sampling. In *Advances in Neural Information Processing Systems 21*.
- Sasa Petrovic, Miles Osborne, and Victor Lavrenko. 2010. Streaming First Story Detection with application to Twitter. In *Proceedings of NAACL*.
- Deepak Ravichandran, Patrick Pantel, and Eduard Hovy. 2005. Randomized Algorithms and NLP: Using Locality Sensitive Hash Functions for High Speed Noun Clustering. In *Proceedings of ACL*.
- David Talbot. 2009. Succinct approximate counting of skewed data. In *Proceedings of IJCAI*.
- Benjamin Van Durme and Ashwin Lall. 2009a. Probabilistic Counting with Randomized Storage. In *Proceedings of IJCAI*.
- Benjamin Van Durme and Ashwin Lall. 2009b. Streaming Pointwise Mutual Information. In *Advances in Neural Information Processing Systems 22*.

Optimizing Question Answering Accuracy by Maximizing Log-Likelihood

Matthias H. Heie, Edward W. D. Whittaker and Sadaoki Furui

Department of Computer Science

Tokyo Institute of Technology

Tokyo 152-8552, Japan

{heie,edw,furui}@furui.cs.titech.ac.jp

Abstract

In this paper we demonstrate that there is a strong correlation between the Question Answering (QA) accuracy and the log-likelihood of the answer typing component of our statistical QA model. We exploit this observation in a clustering algorithm which optimizes QA accuracy by maximizing the log-likelihood of a set of question-and-answer pairs. Experimental results show that we achieve better QA accuracy using the resulting clusters than by using manually derived clusters.

1 Introduction

Question Answering (QA) distinguishes itself from other information retrieval tasks in that the system tries to return accurate answers to queries posed in natural language. Factoid QA limits itself to questions that can usually be answered with a few words. Typically factoid QA systems employ some form of question type analysis, so that a question such as *What is the capital of Japan?* will be answered with a geographical term. While many QA systems use hand-crafted rules for this task, such an approach is time-consuming and doesn't generalize well to other languages. Machine learning methods have been proposed, such as question classification using support vector machines (Zhang and Lee, 2003) and language modeling (Merkel and Klakow, 2007). In these approaches, question categories are predefined and a classifier is trained on manually labeled data. This is an example of supervised learning. In this paper we present an unsupervised method, where we attempt to cluster question-and-answer (q-a) pairs without any predefined question categories, hence no manually class-labeled questions are used.

We use a statistical QA framework, described in Section 2, where the system is trained with clusters

of q-a pairs. This framework was used in several TREC evaluations where it placed in the top 10 of participating systems (Whittaker et al., 2006). In Section 3 we show that answer accuracy is strongly correlated with the log-likelihood of the q-a pairs computed by this statistical model. In Section 4 we propose an algorithm to cluster q-a pairs by maximizing the log-likelihood of a disjoint set of q-a pairs. In Section 5 we evaluate the QA accuracy by training the QA system with the resulting clusters.

2 QA system

In our QA framework we choose to model only the probability of an answer A given a question Q , and assume that the answer A depends on two sets of features: $W = W(Q)$ and $X = X(Q)$:

$$P(A|Q) = P(A|W, X), \quad (1)$$

where W represents a set of $|W|$ features describing the question-type part of Q such as *who*, *when*, *where*, *which*, etc., and X is a set of features which describes the “information-bearing” part of Q , i.e. what the question is actually about and what it refers to. For example, in the questions *Where is Mount Fuji?* and *How high is Mount Fuji?*, the question type features W differ, while the information-bearing features X are identical. Finding the best answer \hat{A} involves a search over all A for the one which maximizes the probability of the above model, i.e.:

$$\hat{A} = \arg \max_A P(A|W, X). \quad (2)$$

Given the correct probability distribution, this will give us the optimal answer in a maximum likelihood sense. Using Bayes' rule, assuming uniform $P(A)$ and that W and X are independent of each other given A , in addition to ignoring $P(W, X)$ since it is independent of A , enables us to rewrite Eq. (2) as

$$\hat{A} = \arg \max_A \underbrace{P(A | X)}_{\text{retrieval model}} \cdot \underbrace{P(W | A)}_{\text{filter model}}. \quad (3)$$

2.1 Retrieval Model

The retrieval model $P(A|X)$ is essentially a language model which models the probability of an answer sequence A given a set of information-bearing features $X = \{x_1, \dots, x_{|X|}\}$. This set is constructed by extracting single-word features from Q that are not present in a stop-list of high-frequency words. The implementation of the retrieval model used for the experiments described in this paper, models the proximity of A to features in X . It is not examined further here; see (Whittaker et al., 2005) for more details.

2.2 Filter Model

The question-type feature set $W = \{w_1, \dots, w_{|W|}\}$ is constructed by extracting n -tuples ($n = 1, 2, \dots$) such as *where*, *in what* and *when were* from the input question Q . We limit ourselves to extracting single-word features. The 2522 most frequent words in a collection of example questions are considered in-vocabulary words; all other words are out-of-vocabulary words, and substituted with $\langle \text{UNK} \rangle$.

Modeling the complex relationship between W and A directly is non-trivial. We therefore introduce an intermediate variable $C_E = \{c_1, \dots, c_{|C_E|}\}$, representing a set of classes of example q-a pairs. In order to construct these classes, given a set $E = \{t_1, \dots, t_{|E|}\}$ of example q-a pairs, we define a mapping function $f : E \mapsto C_E$ which maps each example q-a pair t_j for $j = 1 \dots |E|$ into a particular class $f(t_j) = c_e$. Thus each class c_e may be defined as the union of all component q-a features from each t_j satisfying $f(t_j) = c_e$. Hence each class c_e constitutes a cluster of q-a pairs. Finally, to facilitate modeling we say that W is conditionally independent of A given c_e so that,

$$P(W | A) = \sum_{e=1}^{|C_E|} P(W | c_W^e) \cdot P(c_A^e | A), \quad (4)$$

where c_W^e and c_A^e refer to the subsets of question-type features and example answers for the class c_e , respectively.

$P(W | c_W^e)$ is implemented as trigram language models with backoff smoothing using absolute discounting (Huang et al., 2001).

Due to data sparsity, our set of example q-a pairs cannot be expected to cover all the possible answers to questions that may ever be asked. We therefore employ answer class modeling rather than answer word modeling by expanding Eq. (4) as follows:

$$P(W | A) = \sum_{e=1}^{|C_E|} P(W | c_W^e) \cdot \sum_{a=1}^{|K_A|} P(c_A^e | k_a) P(k_a | A), \quad (5)$$

where k_a is a concrete class in the set of $|K_A|$ answer classes K_A . These classes are generated using the Kneser-Ney clustering algorithm, commonly used for generating class definitions for class language models (Kneser and Ney, 1993).

In this paper we restrict ourselves to single-word answers; see (Whittaker et al., 2005) for the modeling of multi-word answers. We estimate $P(c_A^e | k_A)$ as

$$P(c_A^e | k_A) = \frac{f(k_A, c_A^e)}{\sum_{g=1}^{|C_E|} f(k_A, c_A^g)}, \quad (6)$$

where

$$f(k_A, c_A^e) = \frac{\sum_{\forall i: i \in c_A^e} \delta(i \in k_A)}{|c_A^e|}, \quad (7)$$

and $\delta(\cdot)$ is a discrete indicator function which equals 1 if its argument evaluates true and 0 if false.

$P(k_a | A)$ is estimated as

$$P(k_a | A) = \frac{1}{\sum_{\forall j: j \in K_a} \delta(A \in j)}. \quad (8)$$

3 The Relationship between Mean Reciprocal Rank and Log-Likelihood

We use Mean Reciprocal Rank (MRR) as our metric when evaluating the QA accuracy on a set of questions $G = \{g_1 \dots g_{|G|}\}$:

$$MRR = \frac{\sum_{i=1}^{|G|} 1/R_i}{|G|}, \quad (9)$$

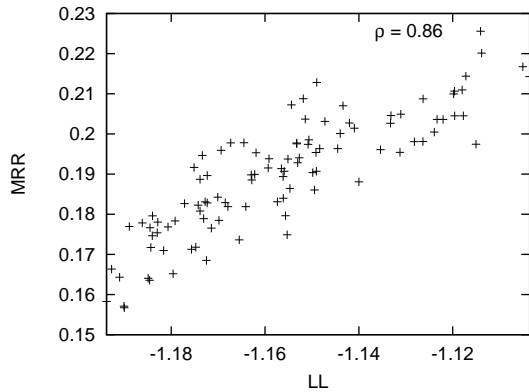


Figure 1: MRR vs. LL (average per q-a pair) for 100 random cluster configurations.

where R_i is the rank of the highest ranking correct candidate answer for g_i .

Given a set $D = (d_1 \dots d_{|D|})$ of q-a pairs disjoint from the q-a pairs in C_E , we can, using Eq. (5), calculate the log-likelihood as

$$\begin{aligned}
 LL &= \sum_{d=1}^{|D|} \log P(W_d | A_d) \\
 &= \sum_{d=1}^{|D|} \log \sum_{e=1}^{|C_E|} P(W_d | c_W^e) \cdot \sum_{a=1}^{|K_A|} P(c_A^e | k_a) P(k_a | A_d). \quad (10)
 \end{aligned}$$

To examine the relationship between MRR and LL , we randomly generate configurations C_E , with a fixed cluster size of 4, and plot the resulting MRR and LL , computed on the same data set D , as data points in a scatter plot, as seen in Figure 1. We find that LL and MRR are strongly correlated, with a correlation coefficient $\rho = 0.86$.

This observation indicates that we should be able to improve the answer accuracy of the QA system by optimizing the LL of the filter model in isolation, similar to how, in automatic speech recognition, the LL of the language model can be optimized in isolation to improve the speech recognition accuracy (Huang et al., 2001).

4 Clustering algorithm

Using the observation that LL is correlated with MRR on the same data set, we expect that optimizing LL on a development set (LL_{dev}) will also improve MRR on an evaluation set (MRR_{eval}). Hence we propose the following greedy algorithm to maximize LL_{dev} :

```

init:  $c_1 \in C_E$  contains all training pairs  $|E|$ 
while  $improvement > threshold$  do
   $best\_LL_{dev} \leftarrow -\infty$ 
  for all  $j = 1 \dots |E|$  do
     $original\_cluster = f(t_j)$ 
    Take  $t_j$  out of  $f(t_j)$ 
    for  $e = -1, 1 \dots |C_E|, |C_E| + 1$  do
      Put  $t_j$  in  $c_e$ 
      Calculate  $LL_{dev}$ 
      if  $LL_{dev} > best\_LL_{dev}$  then
         $best\_LL_{dev} \leftarrow LL_{dev}$ 
         $best\_cluster \leftarrow e$ 
         $best\_pair \leftarrow j$ 
      end if
      Take  $t_j$  out of  $c_e$ 
    end for
    Put  $t_j$  back in  $original\_cluster$ 
  end for
  Take  $t_{best\_pair}$  out of  $f(t_{best\_pair})$ 
  Put  $t_{best\_pair}$  into  $C_{best\_cluster}$ 
end while

```

In this algorithm, c_{-1} indicates the set of training pairs outside the cluster configuration, thus every training pair will not necessarily be included in the final configuration. $c_{|C_E|+1}$ refers to a new, empty cluster, hence this algorithm automatically finds the optimal number of clusters as well as the optimal configuration of them.

5 Experiments

5.1 Experimental Setup

For our data sets, we restrict ourselves to questions that start with *who*, *when* or *where*. Furthermore, we only use q-a pairs which can be answered with a single word. As training data we use questions and answers from the Knowledge-Master collection¹. Development/evaluation questions are the questions from TREC QA evaluations from TREC 2002 to TREC 2006, the answers to which are to be retrieved from the AQUAINT corpus. In total we have 2016 q-a pairs for training and 568 questions for development/evaluation. We are able to retrieve the correct answer for 317 of the development/evaluation questions, thus the theoretical upper bound for our experiments is an answer accuracy of $MRR = 0.558$.

Accuracy is evaluated using 5-fold (rotating) cross-validation, where in each fold the TREC QA data is partitioned into a development set of

¹<http://www.greatauk.com/>

Configuration	LL_{eval}	MRR_{eval}	#clusters
manual	-1.18	0.262	3
all-in-one	-1.32	0.183	1
one-in-each	-0.87	0.263	2016
automatic	-0.24	0.281	4

Table 1: LL_{eval} (average per q-a pair) and MRR_{eval} (over all held-out TREC years), and number of clusters (median of the cross-evaluation folds) for the various configurations.

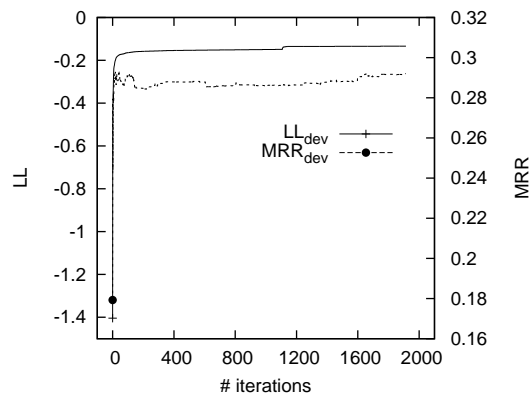
4 years’ data and an evaluation set of one year’s data. For each TREC question the top 50 documents from the AQUAINT corpus are retrieved using Lucene². We use the QA system described in Section 2 for QA evaluation. Our evaluation metric is MRR_{eval} , and LL_{dev} is our optimization criterion, as motivated in Section 3.

Our baseline system uses manual clusters. These clusters are obtained by putting all *who* q-a pairs in one cluster, all *when* pairs in a second and all *where* pairs in a third. We compare this baseline with using clusters resulting from the algorithm described in Section 4. We run this algorithm until there are no further improvements in LL_{dev} . Two other cluster configurations are also investigated: all q-a pairs in one cluster (all-in-one), and each q-a pair in its own cluster (one-in-each). The all-in-one configuration is equivalent to not using the filter model, i.e. answer candidates are ranked solely by the retrieval model. The one-in-each configuration was shown to perform well in the TREC 2006 QA evaluation (Whittaker et al., 2006), where it ranked 9th among 27 participants on the factoid QA task.

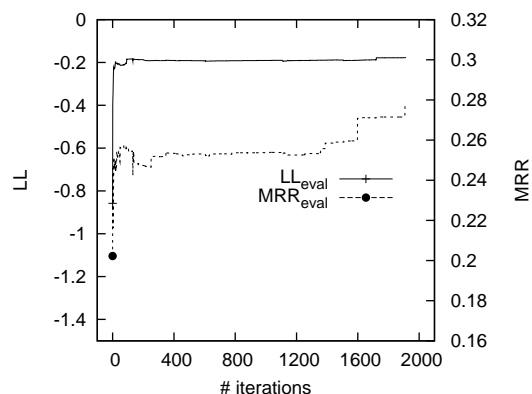
5.2 Results

In Table 1, we see that the manual clusters (baseline) achieves an MRR_{eval} of 0.262, while the clusters resulting from the clustering algorithm give an MRR_{eval} of 0.281, which is a relative improvement of 7%. This improvement is statistically significant at the 0.01 level using the Wilcoxon signed-rank test. The one-in-each cluster configuration achieves an MRR_{eval} of 0.263, which is not a statistically significant improvement over the baseline. The all-in-one cluster configuration (i.e. no filter model) has the lowest accuracy, with an MRR_{eval} of 0.183.

²<http://lucene.apache.org/>



(a) Development set, 4 year’s TREC.



(b) Evaluation set, 1 year’s TREC.

Figure 2: MRR and LL (average per q-a pair) vs. number of algorithm iterations for one cross-validation fold.

6 Discussion

Manual inspection of the automatically derived clusters showed that the algorithm had constructed configurations where typically *who*, *when* and *where* q-a pairs were put in separate clusters, as in the manual configuration. However, in some cases both *who* and *where* q-a pairs occurred in the same cluster, so as to better answer questions like *Who won the World Cup?*, where the answer could be a country name.

As can be seen from Table 1, there are only 4 clusters in the automatic configuration, compared to 2016 in the one-in-each configuration. Since the computational complexity of the filter model described in Section 2.2 is linear in the number of clusters, a beneficial side effect of our clustering procedure is a significant reduction in the computational requirement of the filter model.

In Figure 2 we plot LL and MRR for one of the cross-validation folds over multiple iterations (the *while* loop) of the clustering algorithm in Sec-

tion 4. It can clearly be seen that the optimization of LL_{dev} leads to improvement in MRR_{eval} , and that LL_{eval} is also well correlated with MRR_{eval} .

7 Conclusions and Future Work

In this paper we have shown that the log-likelihood of our statistical model is strongly correlated with answer accuracy. Using this information, we have clustered training q-a pairs by maximizing log-likelihood on a disjoint development set of q-a pairs. The experiments show that with these clusters we achieve better QA accuracy than using manually clustered training q-a pairs.

In future work we will extend the types of questions that we consider, and also allow for multi-word answers.

Acknowledgements

The authors wish to thank Dietrich Klakow for his discussion at the concept stage of this work. The anonymous reviewers are also thanked for their constructive feedback.

References

- [Huang et al.2001] Xuedong Huang, Alex Acero and Hsiao-Wuen Hon. 2001. *Spoken Language Processing*. Prentice-Hall, Upper Saddle River, NJ, USA.
- [Kneser and Ney1993] Reinhard Kneser and Hermann Ney. 1993. Improved Clustering Techniques for Class-based Statistical Language Modelling. *Proceedings of the European Conference on Speech Communication and Technology (EUROSPEECH)*.
- [Merkel and Klakow2007] Andreas Merkel and Dietrich Klakow. 2007. Language Model Based Query Classification. *Proceedings of the European Conference on Information Retrieval (ECIR)*.
- [Whittaker et al.2005] Edward Whittaker, Sadaoki Furui and Dietrich Klakow. 2005. A Statistical Classification Approach to Question Answering using Web Data. *Proceedings of the International Conference on Cyberworlds*.
- [Whittaker et al.2006] Edward Whittaker, Josef Novak, Pierre Chatain and Sadaoki Furui. 2006. TREC 2006 Question Answering Experiments at Tokyo Institute of Technology. *Proceedings of The Fifteenth Text REtrieval Conference (TREC)*.
- [Zhang and Lee2003] Dell Zhang and Wee Sun Lee. 2003. Question Classification using Support Vector Machines. *Proceedings of the Special Interest Group on Information Retrieval (SIGIR)*.

Generating Entailment Rules from FrameNet

Roni Ben Aharon
Department of Computer Science
Bar-Ilan University
Ramat Gan, Israel
r.ben.aharon@gmail.com

Idan Szpektor
Yahoo! Research
Haifa, Israel
idan@yahoo-inc.com

Ido Dagan
Department of Computer Science
Bar-Ilan University
Ramat Gan, Israel
dagan@cs.biu.ac.il

Abstract

Many NLP tasks need accurate knowledge for semantic inference. To this end, mostly WordNet is utilized. Yet WordNet is limited, especially for inference between predicates. To help filling this gap, we present an algorithm that generates inference rules between predicates from FrameNet. Our experiment shows that the novel resource is effective and complements WordNet in terms of rule coverage.

1 Introduction

Many text understanding applications, such as Question Answering (QA) and Information Extraction (IE), need to infer a target textual meaning from other texts. This need was proposed as a generic semantic inference task under the Textual Entailment (TE) paradigm (Dagan et al., 2006).

A fundamental component in semantic inference is the utilization of knowledge resources. However, a major obstacle to improving semantic inference performance is the lack of such knowledge (Bar-Haim et al., 2006; Giampiccolo et al., 2007). We address one prominent type of inference knowledge known as *entailment rules*, focusing specifically on rules between predicates, such as ‘*cure X ⇒ X recover*’.

We aim at highly accurate rule acquisition, for which utilizing manually constructed sources seem appropriate. The most widely used manual resource is WordNet (Fellbaum, 1998). Yet it is incomplete for generating entailment rules between predicates (Section 2.1). Hence, other manual resources should also be targeted.

In this work¹, we explore how FrameNet (Baker et al., 1998) could be effectively used for generating entailment rules between predicates.

¹The detailed description of our work can be found in (Ben Aharon, 2010).

FrameNet is a manually constructed database based on Frame Semantics. It models the semantic argument structure of predicates in terms of prototypical situations called *frames*.

Prior work utilized FrameNet’s argument mapping capabilities but took entailment relations from other resources, namely WordNet. We propose a novel method for generating entailment rules from FrameNet by detecting the entailment relations implied in FrameNet. We utilize FrameNet’s annotated sentences and relations between frames to extract both the entailment relations and their argument mappings.

Our analysis shows that the rules generated by our algorithm have a reasonable “per-rule” accuracy of about 70%². We tested the generated rule-set on an entailment testbed derived from an IE benchmark and compared it both to WordNet and to state-of-the-art rule generation from FrameNet. Our experiment shows that our method outperforms prior work. In addition, our rule-set’s performance is comparable to WordNet and it is complementary to WordNet when uniting the two resources. Finally, additional analysis shows that our rule-set accuracy is 90% in practical use.

2 Background

2.1 Entailment Rules and their Acquisition

To generate entailment rules, two issues should be addressed: a) identifying the lexical entailment relations between predicates, e.g. ‘*cure ⇒ recover*’; b) mapping argument positions, e.g. ‘*cure X ⇒ X recover*’. The main approach for generating highly accurate rule-sets is to use manually constructed resources. To this end, most systems mainly utilize WordNet (Fellbaum, 1998), being the most prominent lexical resource with broad coverage of predicates. Furthermore, some of its

²The rule-set is available at: <http://www.cs.biu.ac.il/~nlp/downloads>

relations capture types of entailment relations, including *synonymy*, *hyponymy*, *morphologically-derived*, *entailment* and *cause*.

Yet, WordNet is limited for entailment rule generation. First, many entailment relations, notably for the WordNet *entailment* and *cause* relation types, are missing, e.g. ‘*elect* \Rightarrow *vote*’. Furthermore, WordNet does not include argument mapping between related predicates. Thus, only *substitutable* WordNet relations (*synonymy* and *hyponymy*), for which argument positions are preserved, could be used to generate entailment rules. The other *non-substitutable* relations, e.g. *cause* (‘*kill* \Rightarrow *die*’) and *morphologically-derived* (‘*meet.v* \Leftrightarrow *meeting.n*’), cannot be used.

2.2 FrameNet

FrameNet (Baker et al., 1998) is a knowledge-base of *frames*, describing prototypical situations. Frames can be related to each other by inter-frame relations, e.g. *Inheritance*, *Precedence*, *Usage* and *Perspective*.

For each frame, several semantic roles are specified, called *frame elements* (FEs), denoting the participants in the situation described. Each FE may be labeled as *core* if it is central to the frame. For example, some core FEs of the *Commerce_pay* frame are *Buyer* and *Goods*, while a non-core FE is *Place*. Each FE may also be labeled with a semantic type, e.g. *Sentient*, *Event*, and *Time*.

A frame includes a list of predicates that can evoke the described situation, called *lexical units* (LUs). LUs are mainly verbs but may also be nouns or adjectives. For example, the frame *Commerce_pay* lists the LUs *pay.v* and *payment.n*.

Finally, FrameNet contains annotated sentences that represent typical LU occurrences in texts. Each annotation refers to one LU in a specific frame and the FEs of the frame that occur in the sentence. An example sentence is “*I_{Buyer} have to pay the bills_{Money}*”. Each sentence is accompanied by a *valence pattern*, which provides, among other info, grammatical functions of the core FEs with respect to the LU. The valence pattern of the above sentence is [(*Buyer Subj*), (*Money Obj*)].

2.3 Using FrameNet for Semantic Inference

To the best of our knowledge, the only work that utilized FrameNet for entailment rule generation is LexPar (Coyne and Rambow, 2009). LexPar first identifies lexical entailment relations by going over all LU pairs which are either in the

same frame or whose frames are related by one of FrameNet’s inter-frame relations. Each candidate pair is considered entailing if the two LUs are either synonyms or in a direct hyponymy relation in WordNet (providing the vast majority of LexPar’s relations), or if their related frames are connected via the *Perspective* relation in FrameNet.

Then, argument mappings between each entailing LU pair are extracted based on the core FEs that are shared between the two LUs. The syntactic positions of the shared FEs are taken from the valence patterns of the LUs. A LexPar rule example is presented in Figure 3 (top part).

Since most of LexPar’s entailment relations are based on WordNet’s relations, LexPar’s rules could be viewed as an intersection of WordNet and FrameNet lexical relations, accompanied with argument mappings taken from FrameNet.

3 Rule Extraction from FrameNet

The above prior work identified lexical entailment relations mainly from WordNet, which limits the use of FrameNet in two ways. First, some relations that appear in FrameNet are missed because they do not appear in WordNet. Second, unlike FrameNet, WordNet does not include argument mappings for its relations. Thus, prior work for rule generation considered only substitutable relations from WordNet (*synonyms* and *hyponyms*), not utilizing FrameNet’s capability to map arguments of non-substitutable relations.

Our goal in this paper is to generate entailment rules solely from the information within FrameNet. We present a novel algorithm for generating entailment rules from FrameNet, called *FRED* (**F**rameNet **E**ntailment-rule **D**erivation), which operates in three steps: a) extracting templates for each LU; b) detecting lexical entailment relations between pairs of LUs; c) generating entailment rules by mapping the arguments between two LUs in each entailing pair.

3.1 Template Extraction

Many LUs in FrameNet are accompanied by annotated sentences (Section 2.2). From each sentence of a given LU, we extract one template for each annotated FE in the sentence. Each template includes the LU, one argument corresponding to the target FE and their syntactic relation in the sentence parse-tree. We focus on extracting unary templates, as they can describe any ar-

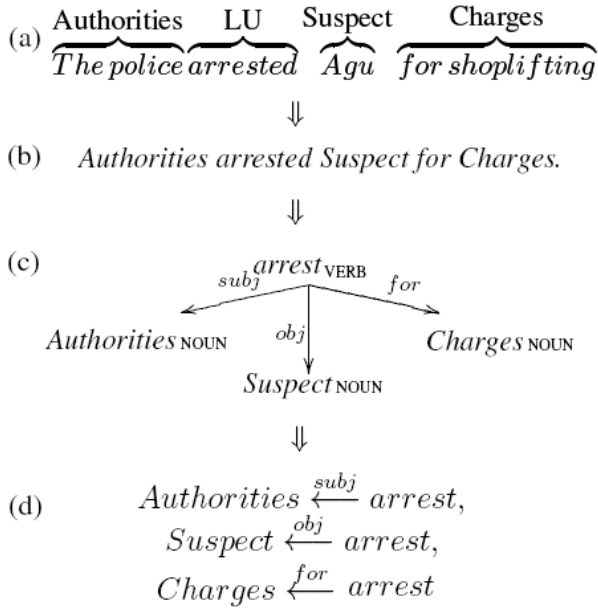


Figure 1: Template extraction for a sentence containing the LU ‘arrest’.

argument mapping by decomposing templates with several arguments into unary ones (Szpektor and Dagan, 2008). Figure 1 exemplifies this process.

As a pre-parsing step, all FE phrases in a given sentence are replaced by their related FE names, excluding syntactic information such as prepositions or possessives (step (b) in Figure 1). Then, the sentence is parsed using the Minipar dependency parser (Lin, 1998) (step (c)). Finally, a path in the parse-tree is extracted between each FE node and the node of the LU (step (d)). Each extracted path is converted into a template by replacing the FE node with an argument variable.

We simplify each extracted path by removing nodes along the path that are not part of the syntactic relation between the LU and the FE, such as conjunctions and other FE nodes. For example, ‘*Authorities* $\xleftarrow{\text{subj}}$ enter $\xrightarrow{\text{conj}}$ arrest’ is simplified into ‘*Authorities* $\xleftarrow{\text{subj}}$ arrest’.

Some templates originated from different annotated sentences share the same LU and syntactic structure, but differ in their FEs. Usually, one of these templates is incorrect, due to erroneous parse (e.g. ‘*Suspect* $\xleftarrow{\text{obj}}$ arrest’ is a correct template, in contrast to ‘*Charges* $\xleftarrow{\text{obj}}$ arrest’). We thus keep only the most frequently annotated template out of the identical templates, assuming it is the correct one.

3.2 Identifying Lexical Entailment Relations

FrameNet groups LUs in frames and describes relations between frames. However, relations between LUs are not explicitly defined. We next describe how we automatically extract several types of lexical entailment relations between LUs using two approaches.

In the first approach, LUs in the same frame that are morphological derivations of each other, e.g. ‘*negotiation.n*’ and ‘*negotiate.v*’, are marked as paraphrases. We take morphological derivation information from the CATVAR database (Habash and Dorr, 2003).

The second approach is based on our observation that some LUs express the prototypical situation that their frame describes, which we denote *dominant LUs*. For example, the LU ‘*recover*’ is dominant for the *Recovery* frame. We mark LUs as dominant if they are morphologically derived from the frame’s name.

Our assumption is that since dominant LUs express the frame’s generic meaning, their meaning is likely to be entailed by the other LUs in this frame. Consequently, we generate such lexical rules between any dominant LU and any other LU in a given frame, e.g. ‘*heal* \Rightarrow *recover*’ and ‘*convalescence* \Rightarrow *recover*’ for the *Recovery* frame.

In addition, we assume that if two frames are related by some type of entailment relation, their dominant LUs are also related by the same relation. Accordingly, we extract entailment relations between dominant LUs of frames that are connected via the *Inheritance*, *Cause* and *Perspective* relations, where *Inheritance* and *Cause* generate directional entailment relations (e.g. ‘*choose* \Rightarrow *decide*’ and ‘*cure* \Rightarrow *recover*’, respectively) while *Perspective* generates bidirectional paraphrase relations (e.g. ‘*transfer* \Leftrightarrow *receive*’).

Finally, we generate the transitive closure of the set of lexical relations identified by the above methods. For example, the combination of ‘*sell* \Leftrightarrow *buy*’ and ‘*buy* \Rightarrow *get*’ generates ‘*sell* \Rightarrow *get*’.

3.3 Generating Entailment Rules

The final step in the FRED algorithm generates lexical syntactic entailment rules from the extracted templates and lexical entailment relations.

For each identified lexical relation ‘*left* \Rightarrow *right*’ between two LUs, the set of FEs that are shared by both LUs is collected. Then, for each shared FE, we take the list of templates that connect this FE

Lexical Relation:cure \Rightarrow recovery**Templates:**

Patient \xleftarrow{obj} *cure* (cure *Patient*)
Affliction \xleftarrow{of} *cure* (cure of *Affliction*)
Patient \xleftarrow{gen} *recovery* (*Patient*'s recovery)
Patient \xleftarrow{of} *recovery* (recovery of *Patient*)
Affliction \xleftarrow{from} *recovery* (recovery from *Affliction*)

Intra-LU Entailment Rules:
 $Patient \xleftarrow{gen} recovery \iff Patient \xleftarrow{of} recovery$
Inter-LU Entailment Rules:

$Patient \xleftarrow{obj} cure \implies Patient \xleftarrow{gen} recovery$
 $Patient \xleftarrow{obj} cure \implies Patient \xleftarrow{of} recovery$
 $Affliction \xleftarrow{of} cure \implies Affliction \xleftarrow{from} recovery$

Figure 2: Some entailment rules generated for the lexical relation ‘*cure.v* \Rightarrow *recovery.n*’.

Configuration	R (%)	P (%)	F1
No-Rules	13.8	57.7	20.9
LexPar	14.1	42.9	17.4
WordNet	18.3	32.2	17.8
FRED	17.6	55.1	24.6
FRED \cup WordNet	21.8	33.3	20.9

Table 1: Macro average Recall (R), Precision (P) and F1 results for the tested configurations.

to each of the LUs, denoted by T_{left}^{fe} and T_{right}^{fe} . Finally, for each template pair, $l \in T_{left}^{fe}$ and $r \in T_{right}^{fe}$, the rule ‘ $l \Rightarrow r$ ’ is generated. In addition, we generate paraphrase rules between the various templates including the same FE and the same LU. Figure 2 illustrates this process.

To improve rule quality, we filter out rules that map FEs of adjunct-like semantic types, such as *Time* and *Location*, since different templates of such FEs may have different semantic meanings (e.g. ‘*Time* \xleftarrow{before} *arrive*’ ‘*Time* \xleftarrow{after} *arrive*’). Thus, it is hard to identify those template pairs that correctly map these FEs for entailment.

We manually evaluated a random sample of 250 rules from the resulting rule-set, out of which we judged 69% as correct.

4 Application-based Evaluation

4.1 Experimental Setup

We would like to evaluate the overall utility of our resource for NLP applications, assessing the correctness of the actual rule applications performed

in practice, as well as to compare its performance to related resources. To this end, we follow the experimental setup presented in (Szpektor and Dagan, 2009), which utilized the ACE 2005 event dataset³ as a testbed for entailment rule-sets. We briefly describe this setup here.

The task is to extract argument mentions for 26 events, such as *Sue* and *Attack*, from the ACE annotated corpus, using a given tested entailment rule-set. Each event is represented by a set of unary *seed templates*, one for each event argument. Some seed templates for *Attack* are ‘*Attacker* \xleftarrow{subj} *attack*’ and ‘*attack* \xrightarrow{obj} *Target*’.

Argument mentions are found in the ACE corpus by matching either the seed templates or templates entailing them found in the tested rule-set. We manually added for each event its relevant WordNet synset-ids and FrameNet frame-ids, so only rules fitting the event target meaning will be extracted from the tested rule-sets.

4.2 Tested Configurations

We evaluated several rule-set configurations:

No-Rules The system matches only the seed templates directly, without any additional rules.

WordNet Rules are generated from WordNet 3.0, using only the *synonymy* and *hypernymy* relations (see Section 2.1). Transitive chaining of relations is allowed (Moldovan and Novischi, 2002).

LexPar Rules are generated from the publicly available LexPar database. We generated unary rules from each LexPar rule based on a manually constructed mapping from FrameNet grammatical functions to Minipar dependency relations. Figure 3 presents an example of this procedure.

FRED Rules are generated by our algorithm.

FRED \cup WordNet The union of the rule-sets of FRED and WordNet.

4.3 Results

Each configuration was tested on each ACE event. We measured *recall*, *precision* and *F1*. Table 1 reports macro averages of the three measures over the 26 ACE events.

As expected, using *No-Rules* achieves the highest precision and the lowest recall compared to all other configurations. When adding LexPar rules,

³<http://projects.ldc.upenn.edu/ace/>

LexPar rule:Lexemes: arrest \rightarrow apprehendValencies: [(Authorities Subj), (Suspect Obj), (Offense (for))] \Rightarrow [(Authorities Subj), (Suspect Obj), (Offense (in))]**Generated unary rules:** $X \xleftarrow{subj} \text{arrest} \Rightarrow X \xleftarrow{subj} \text{apprehend}$, $\text{arrest} \xrightarrow{obj} Y \Rightarrow \text{apprehend} \xrightarrow{obj} Y$, $\text{arrest} \xrightarrow{for} Z \Rightarrow \text{apprehend} \xrightarrow{in} Z$

Figure 3: An example for generation of unary entailment rules from a LexPar rule.

only a slight increase in recall is gained. This shows that the subset of WordNet rules captured by LexPar (Section 2.3) might be too small for the ACE application setting.

When using all WordNet’s substitutable relations, a substantial relative increase in recall is achieved (32%). Yet, precision decreases dramatically (relative decrease of 44%), causing an overall decrease in F1. Most errors are due to correct WordNet rules whose LHS is ambiguous. Since we do not apply a WSD module, these rules are also incorrectly applied to other senses of the LHS. While this phenomenon is common to all rule-sets, WordNet suffers from it the most since it contains many infrequent word senses.

Our main result is that using FRED’s rule-set, recall increases significantly, a relative increase of 27% compared to No-Rules, while precision hardly decreases. Hence, overall F1 is the highest compared to all other configurations (a relative increase of 17% compared to No-Rules). The improvement in F1 is statistically significant compared to all other configurations, according to the two-sided Wilcoxon signed rank test at the level of 0.01 (Wilcoxon, 1945).

FRED performs significantly better than LexPar in both recall, precision and F1 (a relative increase of 25%, 28% and 41% respectively). For example, LexPar hardly utilizes FrameNet’s argument mapping capabilities since most of its rules are based on a sub-set of WordNet’s substitutable relations.

FRED’s precision is substantially higher than WordNet. This mostly results from the fact that FrameNet mainly contains common senses of predicates while WordNet includes many rare word senses; which, as said above, harms precision when WSD is not applied. Error analysis showed that only 7.5% of incorrect extractions are due to erroneous rules in FRED, while the majority of errors are due to sense mismatch or syntactic matching errors of the seed templates of entailing templates in texts.

FRED’s Recall is somewhat lower than Word-

Net, since FrameNet is a much smaller resource. Yet, its rules are mostly complementary to those from WordNet. This added value is demonstrated by the 19% recall increase for the union of FRED and WordNet rule-sets compared to WordNet alone. FRED provides mainly argument mappings for non-substitutable WordNet relations, e.g. ‘*attack.n on X \Rightarrow attack.v X*’, but also lexical relations that are missing from WordNet, e.g. ‘*ambush.v \Rightarrow attack.v*’.

Overall, our experiment shows that the rule-base generated by FRED seems an appropriate complementary resource to the widely used WordNet-based rules in semantic inference and expansion over predicates. This suggestion is especially appealing since our rule-set performs well even when a WSD module is not applied.

5 Conclusions

We presented FRED, a novel algorithm for generating entailment rules solely from the information contained in FrameNet. Our experiment showed that FRED’s rules perform substantially better than LexPar, the only prior rule-set derived from FrameNet. In addition, FRED’s rule-set largely complements the rules generated from WordNet because it contains argument mappings between non-substitutable predicates, which are missing from WordNet, as well as lexical relations that are not included in WordNet.

In future work we plan to investigate combining FrameNet and WordNet rule-sets in a transitive manner, instead of their simple union.

Acknowledgments

This work was partially supported by the Rector’s research grant of Bar-Ilan University, the PASCAL-2 Network of Excellence of the European Community FP7-ICT-2007-1-216886 and the Israel Science Foundation grant 1112/08.

References

- Collin Baker, Charles Fillmore, and John Lowe. 1998. The berkeley framenet project. In *Proceedings of COLING-ACL*, Montreal, Canada.
- Roy Bar-Haim, Ido Dagan, Bill Dolan, Lisa Ferro, Danilo Giampiccolo, Bernardo Magnini, and Idan Szpektor. 2006. The second pascal recognising textual entailment challenge. In *Second PASCAL Challenge Workshop for Recognizing Textual Entailment*.
- Roni Ben Aharon. 2010. Generating entailment rules from framenet. Master's thesis, Bar-Ilan University.
- Robert Coyne and Owen Rambow. 2009. Lexpar: A freely available english paraphrase lexicon automatically extracted from framenet. In *Proceedings of the Third IEEE International Conference on Semantic Computing*.
- Ido Dagan, Oren Glickman, and Bernardo Magnini. 2006. The pascal recognising textual entailment challenge. In *Lecture Notes in Computer Science*, volume 3944, pages 177–190.
- Christiane Fellbaum, editor. 1998. *WordNet: An Electronic Lexical Database*. MIT Press, Cambridge, Massachusetts.
- Danilo Giampiccolo, Bernardo Magnini, Ido Dagan, and Bill Dolan. 2007. The third pascal recognizing textual entailment challenge. In *Proceedings of the ACL-PASCAL Workshop on Textual Entailment and Paraphrasing*.
- Nizar Habash and Bonnie Dorr. 2003. A categorial variation database for english. In *Proceedings of the North American Association for Computational Linguistics (NAACL '03)*, pages 96–102, Edmonton, Canada. Association for Computational Linguistics.
- Dekang Lin. 1998. Dependency-based evaluation of minipar. In *Proceedings of the Workshop on Evaluation of Parsing Systems at LREC*.
- Dan Moldovan and Adrian Novischi. 2002. Lexical chains for question answering. In *Proceedings of COLING*.
- Idan Szpektor and Ido Dagan. 2008. Learning entailment rules for unary templates. In *Proceedings of the 22nd International Conference on Computational Linguistics (Coling 2008)*, pages 849–856, Manchester, UK, August.
- Idan Szpektor and Ido Dagan. 2009. Augmenting wordnet-based inference with argument mapping. In *Proceedings of the 2009 Workshop on Applied Textual Inference*, pages 27–35, Suntec, Singapore, August.
- Frank Wilcoxon. 1945. Individual comparisons by ranking methods. *Biometrics Bulletin*, 1(6):80–83.

Don't 'have a clue'?

Unsupervised co-learning of downward-entailing operators

Cristian Danescu-Niculescu-Mizil and Lillian Lee

Department of Computer Science, Cornell University
cristian@cs.cornell.edu, lle@cs.cornell.edu

Abstract

Researchers in textual entailment have begun to consider inferences involving *downward-entailing operators*, an interesting and important class of lexical items that change the way inferences are made. Recent work proposed a method for learning English downward-entailing operators that requires access to a high-quality collection of *negative polarity items* (NPIs). However, English is one of the very few languages for which such a list exists. We propose the first approach that can be applied to the many languages for which there is no pre-existing high-precision database of NPIs. As a case study, we apply our method to Romanian and show that our method yields good results. Also, we perform a cross-linguistic analysis that suggests interesting connections to some findings in linguistic typology.

1 Introduction

Cristi: "Nicio" ... is that adjective you've mentioned.

Anca: A negative pronominal adjective.

Cristi: You mean there are people who analyze that kind of thing?

Anca: The Romanian Academy.

Cristi: They're crazy.

—From the movie *Police*, adjective

Downward-entailing operators are an interesting and varied class of lexical items that change the default way of dealing with certain types of inferences. They thus play an important role in understanding natural language [6, 18–20, etc.].

We explain what downward entailing means by first demonstrating the "default" behavior, which is *upward entailing*. The word 'observed' is an example upward-entailing operator: the statement

(i) 'Witnesses observed opium use.'

implies

(ii) 'Witnesses observed narcotic use.'

but not vice versa (we write $i \Rightarrow (\not\Leftarrow) ii$). That is, the truth value is preserved if we replace the

argument of an upward-entailing operator by a superset (a more general version); in our case, the set 'opium use' was replaced by the superset 'narcotic use'.

Downward-entailing (DE) (also known as *downward monotonic* or *monotone decreasing*) operators violate this default inference rule: with DE operators, reasoning instead goes from "sets to subsets". An example is the word 'bans':

'The law bans opium use'

$\not\Leftarrow (\Leftarrow)$

'The law bans narcotic use'.

Although DE behavior represents an exception to the default, DE operators are as a class rather common. They are also quite diverse in sense and even part of speech. Some are simple negations, such as 'not', but some other English DE operators are 'without', 'reluctant to', 'to doubt', and 'to allow'.¹ This variety makes them hard to extract automatically.

Because DE operators violate the default "sets to supersets" inference, identifying them can potentially improve performance in many NLP tasks. Perhaps the most obvious such tasks are those involving textual entailment, such as question answering, information extraction, summarization, and the evaluation of machine translation [4]. Researchers are in fact beginning to build textual-entailment systems that can handle inferences involving downward-entailing operators other than simple negations, although these systems almost all rely on small handcrafted lists of DE operators [1–3, 15, 16].² Other application areas are natural-language generation and human-computer interaction, since downward-entailing inferences induce

¹Some examples showing different constructions for analyzing these operators: 'The defendant does not own a blue car' $\not\Leftarrow (\Leftarrow)$ 'The defendant does not own a car'; 'They are reluctant to tango' $\not\Leftarrow (\Leftarrow)$ 'They are reluctant to dance'; 'Police doubt Smith threatened Jones' $\not\Leftarrow (\Leftarrow)$ 'Police doubt Smith threatened Jones or Brown'; 'You are allowed to use Mastercard' $\not\Leftarrow (\Leftarrow)$ 'You are allowed to use any credit card'.

²The exception [2] employs the list automatically derived by Danescu-Niculescu-Mizil, Lee, and Ducott [5], described later.

greater cognitive load than inferences in the opposite direction [8].

Most NLP systems for the applications mentioned above have only been deployed for a small subset of languages. A key factor is the lack of relevant resources for other languages. While one approach would be to separately develop a method to acquire such resources for each language individually, we instead aim to ameliorate the resource-scarcity problem in the case of DE operators wholesale: we propose a single unsupervised method that can extract DE operators in any language for which raw text corpora exist.

Overview of our work Our approach takes the English-centric work of Danescu-Niculescu-Mizil et al. [5] — DLD09 for short — as a starting point, as they present the first and, until now, only algorithm for automatically extracting DE operators from data. However, our work departs significantly from DLD09 in the following key respect.

DLD09 critically depends on access to a high-quality, carefully curated collection of *negative polarity items (NPIs)* — lexical items such as ‘any’, ‘ever’, or the idiom ‘have a clue’ that tend to occur only in negative environments (see §2 for more details). DLD09 use NPIs as signals of the occurrence of downward-entailing operators. However, almost every language other than English lacks a high-quality accessible NPI list.

To circumvent this problem, we introduce a knowledge-lean *co-learning* approach. Our algorithm is initialized with a very small seed set of NPIs (which we describe how to generate), and then iterates between (a) discovering a set of DE operators using a collection of *pseudo-NPIs* — a concept we introduce — and (b) using the newly-acquired DE operators to detect new pseudo-NPIs.

Why this isn’t obvious Although the algorithmic idea sketched above seems quite simple, it is important to note that prior experiments in that direction have not proved fruitful. Preliminary work on learning (German) NPIs using a small list of simple known DE operators did not yield strong results [14]. Hoeksema [10] discusses why NPIs might be hard to learn from data.³ We circumvent this problem because we are not interested in learning NPIs per se; rather, for our pur-

³In fact, humans can have trouble agreeing on NPI-hood; for instance, Lichte and Soehn [14] mention doubts about over half of Kürschner [12]’s 344 manually collected German NPIs.

poses, pseudo-NPIs suffice. Also, our preliminary work determined that one of the most famous co-learning algorithms, *hubs and authorities* or *HITS* [11], is poorly suited to our problem.⁴

Contributions To begin with, we apply our algorithm to produce the first large list of DE operators for a language other than English. In our case study on Romanian (§4), we achieve quite high precisions at k (for example, iteration achieves a precision at 30 of 87%).

Auxiliary experiments explore the effects of using a large but noisy NPI list, should one be available for the language in question. Intriguingly, we find that co-learning new pseudo-NPIs provides better results.

Finally (§5), we engage in some cross-linguistic analysis based on the results of applying our algorithm to English. We find that there are some suggestive connections with findings in linguistic typology.

Appendix available A more complete account of our work and its implications can be found in a version of this paper containing appendices, available at www.cs.cornell.edu/~cristian/acl2010/.

2 DLD09: successes and challenges

In this section, we briefly summarize those aspects of the DLD09 method that are important to understanding how our new co-learning method works.

DE operators and NPIs Acquiring DE operators is challenging because of the complete lack of annotated data. DLD09’s insight was to make use of *negative polarity items (NPIs)*, which are words or phrases that tend to occur only in negative contexts. The reason they did so is that Ladusaw’s hypothesis [7, 13] asserts that *NPIs only occur within the scope of DE operators*. Figure 1 depicts examples involving the English NPIs ‘any’⁵ and ‘have a clue’ (in the idiomatic sense) that illustrate this relationship. Some other English NPIs are ‘ever’, ‘yet’ and ‘give a damn’.

Thus, NPIs can be treated as clues that a DE operator might be present (although DE operators may also occur without NPIs).

⁴We explored three different edge-weighting schemes based on co-occurrence frequencies and seed-set membership, but the results were extremely poor; HITS invariably retrieved very frequent words.

⁵The *free-choice* sense of ‘any’, as in ‘I can skim any paper in five minutes’, is a known exception.

DE operators	NPIs	
	<i>any</i> ³	<i>have a clue</i> , idiomatic sense
not or n't	✓ We do n't have <i>any</i> apples	✓ We do n't <i>have a clue</i>
doubt	✓ I doubt they have <i>any</i> apples	✓ I doubt they <i>have a clue</i>
no DE operator	× They have <i>any</i> apples	× They <i>have a clue</i>

Figure 1: Examples consistent with Ladusaw’s hypothesis that NPIs can only occur within the scope of DE operators. A ✓ denotes an acceptable sentence; a × denotes an unacceptable sentence.

DLD09 algorithm Potential DE operators are collected by extracting those words that appear in an NPI’s context at least once.⁶ Then, the potential DE operators x are ranked by

$$f(x) := \frac{\text{fraction of NPI contexts that contain } x}{\text{relative frequency of } x \text{ in the corpus}},$$

which compares x ’s probability of occurrence conditioned on the appearance of an NPI with its probability of occurrence overall.⁷

The method just outlined requires access to a list of NPIs. DLD09’s system used a subset of John Lawler’s carefully curated and “moderately complete” list of English NPIs.⁸ The resultant rankings of candidate English DE operators were judged to be of high quality.

The challenge in porting to other languages: cluelessness Can the unsupervised approach of DLD09 be successfully applied to languages other than English? Unfortunately, for most other languages, it does not seem that large, high-quality NPI lists are available.

One might wonder whether one can circumvent the NPI-acquisition problem by simply translating a known English NPI list into the target language. However, NPI-hood need not be preserved under translation [17]. Thus, for most languages, we lack the critical clues that DLD09 depends on.

3 Getting a clue

In this section, we develop an iterative co-learning algorithm that can extract DE operators in the many languages where a high-quality NPI

⁶DLD09 policies: (a) “NPI context” was defined as the part of the sentence to the left of the NPI up to the first comma, semi-colon or beginning of sentence; (b) to encourage the discovery of new DE operators, those sentences containing one of a list of 10 well-known DE operators were discarded. For Romanian, we treated only negations (‘nu’ and ‘n-’) and questions as well-known environments.

⁷DLD09 used an additional *distilled* score, but we found that the distilled score performed worse on Romanian.

⁸<http://www-personal.umich.edu/~jlawler/ae/npi.html>

database is not available, using Romanian as a case study.

3.1 Data and evaluation paradigm

We used Rada Mihalcea’s corpus of ≈ 1.45 million sentences of raw Romanian newswire articles.

Note that we cannot evaluate impact on textual inference because, to our knowledge, no publicly available textual-entailment system or evaluation data for Romanian exists. We therefore examine the system outputs directly to determine whether the top-ranked items are actually DE operators or not. Our evaluation metric is precision at k of a given system’s ranked list of candidate DE operators; it is not possible to evaluate recall since no list of Romanian DE operators exists (a problem that is precisely the motivation for this paper).

To evaluate the results, two native Romanian speakers labeled the system outputs as being “DE”, “not DE” or “Hard (to decide)”. The labeling protocol, which was somewhat complex to prevent bias, is described in the [externally-available appendices](#) (§7.1). The complete system output and annotations are publicly available at: <http://www.cs.cornell.edu/~cristian/acl2010/>.

3.2 Generating a seed set

Even though, as discussed above, the translation of an NPI need not be an NPI, a preliminary review of the literature indicates that in many languages, there is some NPI that can be translated as ‘any’ or related forms like ‘anybody’. Thus, with a small amount of effort, one can form a minimal NPI seed set for the DLD09 method by using an appropriate target-language translation of ‘any’. For Romanian, we used ‘vreo’ and ‘vreun’, which are the feminine and masculine translations of English ‘any’.

3.3 DLD09 using the Romanian seed set

We first check whether DLD09 with the two-item seed set described in §3.2 performs well on Romanian. In fact, the results are fairly poor:

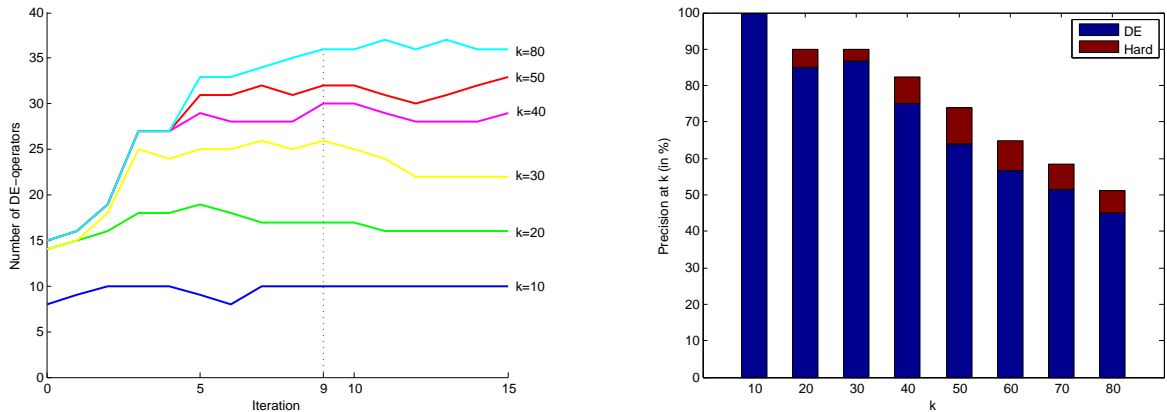


Figure 2: **Left:** Number of DE operators in the top k results returned by the co-learning method at each iteration. Items labeled “Hard” are not included. Iteration 0 corresponds to DLD09 applied to $\{\text{‘vreo’}, \text{‘vreun’}\}$. Curves for $k = 60$ and 70 omitted for clarity. **Right:** Precisions at k for the results of the 9th iteration. The bar divisions are: DE (blue/darkest/largest) and Hard (red/lighter, sometimes non-existent).

for example, the precision at 30 is below 50%. (See blue/dark bars in figure 3 in the [externally-available appendices](#) for detailed results.)

This relatively unsatisfactory performance may be a consequence of the very small size of the NPI list employed, and may therefore indicate that it would be fruitful to investigate automatically extending our list of clues.

3.4 Main idea: a co-learning approach

Our main insight is that not only can NPIs be used as clues for finding DE operators, as shown by DLD09, but conversely, DE operators (if known) can potentially be used to discover new NPI-like clues, which we refer to as *pseudo-NPIs* (or *pNPIs* for short). By “NPI-like” we mean, “serve as possible indicators of the presence of DE operators, regardless of whether they are actually restricted to negative contexts, as true NPIs are”. For example, in English newswire, the words ‘allegation’ or ‘rumor’ tend to occur mainly in DE contexts, like ‘denied’ or ‘dismissed’, even though they are clearly not true NPIs (the sentence ‘I heard a rumor’ is fine). Given this insight, we approach the problem using an *iterative co-learning* paradigm that integrates the search for new DE operators with a search for new pNPIs.

First, we describe an algorithm that is the “reverse” of DLD09 (henceforth *rDLD*), in that it retrieves and ranks pNPIs assuming a given list of DE operators. Potential pNPIs are collected by extracting those words that appear in a DE context (defined here, to avoid the problems of parsing or scope determination, as the part of the sentence to

the right of a DE operator, up to the first comma, semi-colon or end of sentence); these candidates x are then ranked by

$$f_r(x) := \frac{\text{fraction of DE contexts that contain } x}{\text{relative frequency of } x \text{ in the corpus}}.$$

Then, our *co-learning* algorithm consists of the iteration of the following two steps:

- (*DE learning*) Apply DLD09 using a set \mathcal{N} of pseudo-NPIs to retrieve a list of candidate DE operators ranked by f (defined in Section 2). Let \mathcal{D} be the top n candidates in this list.
- (*pNPI learning*) Apply rDLD using the set \mathcal{D} to retrieve a list of pNPIs ranked by f_r ; extend \mathcal{N} with the top n_r pNPIs in this list. Increment n .

Here, \mathcal{N} is initialized with the NPI seed set. At each iteration, we consider the output of the algorithm to be the ranked list of DE operators retrieved in the DE-learning step. In our experiments, we initialized n to 10 and set n_r to 1.

4 Romanian results

Our results show that there is indeed favorable synergy between DE-operator and pNPI retrieval. Figure 2 plots the number of correctly retrieved DE operators in the top k outputs at each iteration. The point at iteration 0 corresponds to a datapoint already discussed above, namely, DLD09 applied to the two ‘any’-translation NPIs. Clearly, we see general substantial improvement over DLD09, although the increases level off in later iterations.

(Determining how to choose the optimal number of iterations is a subject for future research.)

Additional experiments, described in the [externally-available appendices](#) (§7.2), suggest that pNPIs can even be more effective clues than a noisy list of NPIs. (Thus, a larger seed set does not necessarily mean better performance.) pNPIs also have the advantage of being derivable automatically, and might be worth investigating from a linguistic perspective in their own right.

5 Cross-linguistic analysis

Applying our algorithm to English: connections to linguistic typology So far, we have made no assumptions about the language on which our algorithm is applied. A valid question is, does the quality of the results vary with choice of application language? In particular, what happens if we run our algorithm on English?

Note that in some sense, this is a perverse question: the motivation behind our algorithm is the non-existence of a high-quality list of NPIs for the language in question, and English is essentially the only case that does not fit this description. On the other hand, the fact that DLD09 applied their method for extraction of DE operators to English necessitates some form of comparison, for the sake of experimental completeness.

We thus ran our algorithm on the English BLLIP newswire corpus with seed set {‘any’}. We observe that, surprisingly, the iterative addition of pNPIs has very little effect: the precisions at k are good at the beginning and stay about the same across iterations (for details see figure 5 in the [externally-available appendices](#)). Thus, on English, co-learning does not hurt performance, which is good news; but unlike in Romanian, it does not lead to improvements.

Why is English ‘any’ seemingly so “powerful”, in contrast to Romanian, where iterating beyond the initial ‘any’ translations leads to better results? Interestingly, findings from linguistic typology may shed some light on this issue. Haspelmath [9] compares the functions of indefinite pronouns in 40 languages. He shows that English is one of the minority of languages (11 out of 40)⁹ in which there exists an indefinite pronoun series that occurs in all (Haspelmath’s) classes of DE contexts, and thus can constitute a sufficient seed on

⁹English, Ancash Quechua, Basque, Catalan, French, Hindi/Urdu, Irish, Portuguese, Swahili, Swedish, Turkish.

its own. In the other languages (including Romanian),¹⁰ no indirect pronoun can serve as a sufficient seed. So, we expect our method to be viable for all languages; while the iterative discovery of pNPIs is not necessary (although neither is it harmful) for the subset of languages for which a sufficient seed exists, such as English, it is essential for the languages for which, like Romanian, ‘any’-equivalents do not suffice.

Using translation Another interesting question is whether directly translating DE operators from English is an alternative to our method. First, we emphasize that there exists no complete list of English DE operators (the largest available collection is the one extracted by DLD09). Second, we do not know whether DE operators in one language translate into DE operators in another language. Even if that were the case, and we somehow had access to ideal translations of DLD09’s list, there would still be considerable value in using our method: 14 (39%) of our top 36 highest-ranked Romanian DE operators for iteration 9 do not, according to the Romanian-speaking author, have English equivalents appearing on DLD09’s 90-item list. Some examples are: ‘abținut’ (abstained), ‘criticat’ (criticized) and ‘reacționat’ (reacted). Therefore, a significant fraction of the DE operators derived by our co-learning algorithm would have been missed by the translation alternative even under ideal conditions.

6 Conclusions

We have introduced the first method for discovering downward-entailing operators that is universally applicable. Previous work on automatically detecting DE operators assumed the existence of a high-quality collection of NPIs, which renders it inapplicable in most languages, where such a resource does not exist. We overcome this limitation by employing a novel *co-learning* approach, and demonstrate its effectiveness on Romanian.

Also, we introduce the concept of *pseudo-NPIs*. Auxiliary experiments described in the [externally-available appendices](#) show that pNPIs are actually more effective seeds than a noisy “true” NPI list.

Finally, we noted some cross-linguistic differences in performance, and found an interesting connection between these differences and Haspelmath’s [9] characterization of cross-linguistic variation in the occurrence of indefinite pronouns.

¹⁰Examples: Chinese, German, Italian, Polish, Serbian.

Acknowledgments We thank Tudor Marian for serving as an annotator, Rada Mihalcea for access to the Romanian newswire corpus, and Claire Cardie, Yejin Choi, Effi Georgala, Mark Liberman, Myle Ott, João Paula Muchado, Stephen Purpura, Mark Yatskar, Ainur Yessenalina, and the anonymous reviewers for their helpful comments. Supported by NSF grant IIS-0910664.

References

- [1] Roy Bar-Haim, Jonathan Berant, Ido Dagan, Iddo Greental, Shachar Mirkin, Eyal Shnarch, and Idan Szpektor. Efficient semantic deduction and approximate matching over compact parse forests. In *Proceedings of the Text Analysis Conference (TAC)*, 2008.
- [2] Eric Breck. A simple system for detecting non-entailment. In *Proceedings of the Text Analysis Conference (TAC)*, 2009.
- [3] Christos Christodoulopoulos. Creating a natural logic inference system with combinatory categorial grammar. Master’s thesis, University of Edinburgh, 2008.
- [4] Ido Dagan, Oren Glickman, and Bernardo Magnini. The PASCAL Recognising Textual Entailment challenge. In *Machine Learning Challenges, Evaluating Predictive Uncertainty, Visual Object Classification and Recognizing Textual Entailment, First PASCAL Machine Learning Challenges Workshop*, pages 177–190. Springer, 2006.
- [5] Cristian Danescu-Niculescu-Mizil, Lillian Lee, and Richard Duccott. Without a ‘doubt’? Unsupervised discovery of downward-entailing operators. In *Proceedings of NAACL HLT*, 2009.
- [6] David Dowty. The role of negative polarity and concord marking in natural language reasoning. In Mandy Harvey and Lynn Santelmann, editors, *Proceedings of SALT IV*, pages 114–144, 1994.
- [7] Gilles Fauconnier. Polarity and the scale principle. In *Proceedings of the Chicago Linguistic Society (CLS)*, pages 188–199, 1975. Reprinted in Javier Gutierrez-Rexach (ed.), *Semantics: Critical Concepts in Linguistics*, 2003.
- [8] Bart Geurts and Frans van der Slik. Monotonicity and processing load. *Journal of Semantics*, 22(1):97–117, 2005.
- [9] Martin Haspelmath. *Indefinite Pronouns*. Oxford University Press, 2001.
- [10] Jack Hoeksema. Corpus study of negative polarity items. *IV-V Jornades de corpus linguistics 1996-1997*, 1997. <http://odur.let.rug.nl/~hoeksema/docs/barcelona.html>.
- [11] Jon Kleinberg. Authoritative sources in a hyperlinked environment. In *Proceedings of the 9th ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 668–677, 1998. Extended version in *Journal of the ACM*, 46:604–632, 1999.
- [12] Wilfried Kürschner. *Studien zur Negation im Deutschen*. Narr, 1983.
- [13] William A. Ladusaw. *Polarity Sensitivity as Inherent Scope Relations*. Garland Press, New York, 1980. Ph.D. thesis date 1979.
- [14] Timm Lichte and Jan-Philipp Soehn. The retrieval and classification of Negative Polarity Items using statistical profiles. In Sam Featherston and Wolfgang Sternefeld, editors, *Roots: Linguistics in Search of its Evolutionary Base*, pages 249–266. Mouton de Gruyter, 2007.
- [15] Bill MacCartney and Christopher D. Manning. Modeling semantic containment and exclusion in natural language inference. In *Proceedings of COLING*, pages 521–528, 2008.
- [16] Rowan Nairn, Cleo Condoravdi, and Lauri Karttunen. Computing relative polarity for textual inference. In *Proceedings of Inference in Computational Semantics (ICoS)*, 2006.
- [17] Frank Richter, Janina Radó, and Manfred Sailer. Negative polarity items: Corpus linguistics, semantics, and psycholinguistics: Day 2: Corpus linguistics. Tutorial slides: <http://www.sfs.uni-tuebingen.de/~fr/esslli/08/byday/day2/day2-part1.pdf>, 2008.
- [18] Víctor Sánchez Valencia. *Studies on natural logic and categorial grammar*. PhD thesis, University of Amsterdam, 1991.
- [19] Johan van Benthem. *Essays in Logical Semantics*. Reidel, Dordrecht, 1986.
- [20] Ton van der Wouden. *Negative contexts: Collocation, polarity and multiple negation*. Routledge, 1997.

Vocabulary Choice as an Indicator of Perspective

Beata Beigman Klebanov, Eyal Beigman, Daniel Diermeier

Northwestern University and Washington University in St. Louis

beata, d-diermeier@northwestern.edu, beigman@wustl.edu

Abstract

We establish the following characteristics of the task of perspective classification: (a) using term frequencies in a document does not improve classification achieved with absence/presence features; (b) for datasets allowing the relevant comparisons, a small number of top features is found to be as effective as the full feature set and indispensable for the best achieved performance, testifying to the existence of perspective-specific keywords. We relate our findings to research on word frequency distributions and to discourse analytic studies of perspective.

1 Introduction

We address the task of perspective classification. Apart from the spatial sense not considered here, perspective can refer to an agent's role (doctor vs patient in a dialogue), or understood as "a particular way of thinking about something, especially one that is influenced by one's beliefs or experiences," stressing the manifestation of one's broader perspective in some specific issue, or "the state of one's ideas, the facts known to one, etc., in having a meaningful interrelationship," stressing the meaningful connectedness of one's stances and pronouncements on possibly different issues.¹

Accordingly, one can talk about, say, opinion on a particular proposed legislation on abortion within pro-choice or pro-life perspectives; in this case, perspective essentially boils down to opinion in a particular debate. Holding the issue constant but relaxing the requirement of a debate on a specific document, we can consider writings from pro- and con- perspective, in, for example, the death penalty controversy over a course of a period of time. Relaxing the issue specificity somewhat,

one can talk about perspectives of people on two sides of a conflict; this is not opposition or support for any particular proposal, but ideas about a highly related cluster of issues, such as Israeli and Palestinian perspectives on the conflict in all its manifestations. Zooming out even further, one can talk about perspectives due to certain life contingencies, such as being born and raised in a particular culture, region, religion, or political tradition, such perspectives manifesting themselves in certain patterns of discourse on a wide variety of issues, for example, views on political issues in the Middle East from Arab vs Western observers.

In this article, we consider perspective at all the four levels of abstraction. We apply the same types of models to all, in order to discover any common properties of perspective classification. We contrast it with text categorization and with opinion classification by employing models routinely used for such tasks. Specifically, we consider models that use term frequencies as features (usually found to be superior for text categorization) and models that use term absence/presence (usually found to be superior for opinion classification). We motivate our hypothesis that presence/absence features would be as good as or better than frequencies, and test it experimentally. Secondly, we investigate the question of feature redundancy often observed in text categorization.

2 Vocabulary Selection

A line of inquiry going back at least to Zipf strives to characterize word frequency distributions in texts and corpora; see Baayen (2001) for a survey. One of the findings in this literature is that a multinomial (called "urn model" by Baayen) is not a good model for word frequency distributions. Among the many proposed remedies (Baayen, 2001; Jansche, 2003; Baroni and Evert, 2007; Bhat and Sproat, 2009), we would like to draw attention to the following insight articulated

¹Google English Dictionary, Dictionary.com

most clearly in Jansche (2003). Estimation is improved if texts are construed as being generated by two processes, one choosing which words would appear at all in the text, and then, for words that have been chosen to appear, how many times they would in fact appear. Jansche (2003) describes a two-stage generation process: (1) Toss a z -biased coin; if it comes up heads, generate 0; if it comes up tails, (2) generate according to $F(\theta)$, where $F(\theta)$ is a negative binomial distribution and z is a parameter controlling the extent of zero-inflation.

The postulation of two separate processes is effective for predicting word frequencies, but is there any meaning to the two processes? The first process of deciding on the vocabulary, or word types, for the text – what is its function? Jansche (2003) suggests that the zero-inflation component takes care of the multitude of vocabulary words that are not “on topic” for the given text, including taboo words, technical jargon, proper names. This implies that words that are chosen to appear are all “on topic”. Indeed, text segmentation studies show that tracing recurrence of words in a text permits topical segmentation (Hearst, 1997; Hoey, 1991). Yet, if a person compares abortion to *infanticide* – are we content with describing this word as being merely “on topic,” that is, having a certain probability of occurrence once the topic of abortion comes up? In fact, it is only likely to occur if the speaker holds a pro-life perspective, while a pro-choicer would avoid this term.

We therefore hypothesize that the choice of vocabulary is not only a matter of topic but also of perspective, while word *recurrence* has mainly to do with the topical composition of the text. Therefore, tracing word frequencies is not going to be effective for perspective classification beyond noting the mere presence/absence of words, differently from the findings in text categorization, where frequency-based features usually do better than boolean features for sufficiently large vocabulary sizes (McCallum and Nigam, 1998).

3 Data

Partial Birth Abortion (PBA) debates: We use transcripts of the debates on Partial Birth Abortion Ban Act on the floors of the US House and Senate in 104-108 Congresses (1995-2003). Similar legislation was proposed multiple times, passed the legislatures, and, after having initially been vetoed by President Clinton, was signed into law

by President Bush in 2003. We use data from 278 legislators, with 669 speeches in all. We take only one speech per speaker per year; since many serve multiple years, each speaker is represented with 1 to 5 speeches. We perform 10-fold cross-validation splitting by speakers, so that all speeches by the same speaker are assigned to the same fold and testing is always inter-speaker.

When deriving the label for perspective, it is important to differentiate between a particular legislation and a pro-choice / pro-life perspective. A pro-choice person might still support the bill: “I am pro-choice, but believe late-term abortions are wrong. Abortion is a very personal decision and a woman’s right to choose whether to terminate a pregnancy subject to the restrictions of *Roe v. Wade* must be protected. In my judgment, however, the use of this particular procedure cannot be justified.” (Rep. Shays, R-CT, 2003). To avoid inconsistency between vote and perspective, we use data from pro-choice and pro-life non-governmental organizations, NARAL and NRLC, that track legislators’ votes on abortion-related bills, showing the percentage of times a legislator supported the side the organization deems consistent with its perspective. We removed 22 legislators with a mixed record, that is, those who gave 20-60% support to one of the positions.²

Death Penalty (DP) blogs: We use University of Maryland Death Penalty Corpus (Greene and Resnik, 2009) of 1085 texts from a number of pro- and anti-death penalty websites. We report 4-fold cross-validation (DP-4) using the folds in Greene and Resnik (2009), where training and testing data come from different websites for each of the sides, as well as 10-fold cross-validation performance on the entire corpus, irrespective of the site.³

Bitter Lemons (BL): We use the GUEST part of the BitterLemons corpus (Lin et al., 2006), containing 296 articles published in 2001-2005 on <http://www.bitterlemons.org> by more than 200 different Israeli and Palestinian writers on issues related to the conflict.

Bitter Lemons International (BL-I): We collected 150 documents each by a different per-

²Ratings are from: <http://www.OnTheIssues.org/>. We further excluded data from Rep. James Moran, D-VA, as he changed his vote over the years. For legislators rated by neither NRLC nor NARAL, we assumed the vote aligns with the perspective.

³The 10-fold setting yields almost perfect performance likely due to site-specific features beyond perspective per se, hence we do not use this setting in subsequent experiments.

son from either Arab or Western perspectives on Middle Eastern affairs in 2003-2009 from <http://www.bitterlemons-international.org/>. The writers and interviewees on this site are usually former diplomats or government officials, academics, journalists, media and political analysts.⁴ The specific issues cover a broad spectrum, including public life, politics, wars and conflicts, education, trade relations in and between countries like Lebanon, Jordan, Iraq, Egypt, Yemen, Morocco, Saudi Arabia, as well as their relations with the US and members of the European Union.

3.1 Pre-processing

We are interested in perspective manifestations using common English vocabulary. To avoid the possibility that artifacts such as names of senators or states drive the classification, we use as features words that contain only lowercase letters, possibly hyphenated. No stemming is performed, and no stopwords are excluded.⁵

Table 1: Summary of corpora

Data	#Docs	#Features	# CV folds
PBA	669	9.8 K	10
BL	296	10 K	10
BL-I	150	9 K	10
DP	1085	25 K	4

4 Models

For generative models, we use two versions of Naive Bayes models termed *multi-variate Bernoulli* (here, NB-BOOL) and *multinomial* (here, NB-COUNT), respectively, in McCallum and Nigam (1998) study of event models for text categorization. The first records presence/absence of a word in a text, while the second records the number of occurrences. McCallum and Nigam (1998) found NB-COUNT to do better than NB-BOOL for sufficiently large vocabulary sizes for text categorization by topic. For discriminative models, we use linear SVM, with presence-absence, normalized frequency, and tfidf feature weighting. Both types of models are commonly used for text classification tasks. For example, Lin et al. (2006) use

⁴We excluded Israeli, Turkish, Iranian, Pakistani writers as not clearly representing either perspective.

⁵We additionally removed words containing *support*, *oppos*, *sustain*, *override* from the PBA data, in order not to inflate the performance on perspective classification due to the explicit reference to the upcoming vote.

NB-COUNT and SVM-NORMF for perspective classification; Pang et al. (2002) consider most and Yu et al. (2008) all of the above for related tasks of movie review and political party classification. We use SVM^{light} (Joachims, 1999) for SVM and WEKA toolkit (Witten and Frank, 2005; Hall et al., 2009) for both version of Naive Bayes. Parameter optimization for all SVM models is performed using grid search on the training data separately for each partition into train and test data.⁶

5 Results

Table 2 summarizes the cross-validation results for the four datasets discussed above. Notably, the SVM-BOOL model is either the best or not significantly different from the best performing model, although the competitors use more detailed textual information, namely, the count of each word’s appearance in the text, either raw (NB-COUNT), normalized (SVM-NORMF), or combined with document frequency (SVM-TFIDF).

Table 2: Classification accuracy. Scores significantly different from the best performance ($p_{2t} < 0.05$ on paired t-test) are given an asterisk.

Data	NB		SVM		
	BOOL	COUNT	BOOL	NORMF	TFIDF
PBA	*0.93	0.96	0.96	0.96	0.97
DP-4	0.82	0.82	0.83	0.82	0.72 ⁷
DP-10	*0.88	*0.93	0.98	*0.97	*0.97
BL	0.89	0.88	0.89	0.86	0.84
BL-I	0.68	0.66	0.73	0.65	0.65

We conclude that there is no evidence for the relevance of the frequency composition of the text for perspective classification, for all levels of venue- and topic-control, from the tightest (PBA debates) to the loosest (Western vs Arab authors on Middle Eastern affairs). This result is a clear indication that perspective classification is quite different from text categorization by topic, where count-based features usually perform better than boolean features. On the other hand, we have not

⁶Parameter c controlling the trade-off between errors on training data and margin is optimized for all datasets, with the grid $c = \{10^{-6}, 10^{-5}, \dots, 10^5\}$. On the DP data parameter j controlling penalties for misclassification of positive and negative cases is optimized as well ($j = \{10^{-2}, 10^{-1}, \dots, 10^2\}$), since datasets are unbalanced (for example, there is a fold with 27%-73% split).

⁷Here SVM-TFIDF is doing somewhat better than SVM-BOOL on one of the folds and much worse on two other folds; paired t-test with just 4 pairs of observations does not detect a significant difference.

observed that boolean features are reliably better than count-based features, as reported for the sentiment classification task in the movie review domain (Pang et al., 2002).

We note the low performance on BL-I, which could testify to a low degree of lexical consolidation in the Arab vs Western perspectives (more on this below). It is also possible that the small size of BL-I leads to overfitting and low accuracies. However, PBA subset with only 151 items (only 2002 and 2003 speeches) is still 96% classifiable, so size alone does not explain low BL-I performance.

6 Consolidation of perspective

We explore feature redundancy in perspective classification. We first investigate retention of only N best features, then elimination thereof. As a proxy of feature quality, we use the weight assigned to the feature by the SVM-BOOL model based on the training data. Thus, to get the performance with N best features, we take the $\frac{N}{2}$ highest and lowest weight features, for the positive and negative classes, respectively, and retrain SVM-BOOL with these features only.⁸

Table 3: Consolidation of perspective. Nbest shows the smallest N and its proportion out of all features for which the performance of SVM-BOOL with only the best N features is not significantly inferior ($p_{1t} > 0.1$) to that of the full feature set. No-Nbest shows the largest number N for which a model *without* N best features is not significantly inferior to the full model. $N = \{50, 100, 150, \dots, 1000\}$; for DP and BL-I, additionally $N = \{1050, 1100, \dots, 1500\}$; for PBA, additionally $N = \{10, 20, 30, 40\}$.

Data	Nbest		No-Nbest	
	N	%	N	%
PBA	250	2.6%	10	<1%
BL	500	4.9%	100	<1%
DP	100	<1%	1250	5.2%
BL-I	200	2.2%	950	11%

We observe that it is generally sufficient to use a small percentage of the available words to obtain the same classification accuracy as with the full feature set, even in high-accuracy cases such as PBA and BL. The effectiveness of a small subset of features is consistent with the observation in the discourse analysis studies that rivals

⁸We experimented with the mutual information based feature selection as well, with generally worse results.

in long-lasting controversies tend to consolidate their vocabulary and signal their perspective with certain *stigma words* and *banner words*, that is, specific keywords used by a discourse community to implicate adversaries and to create sympathy with own perspective, respectively (Teubert, 2001). Thus, in abortion debates, using *infanticide* as a synonym for abortion is a pro-life stigma. Note that this does not mean the rest of the features are not informative for classification, only that they are redundant with respect to a small percentage of top weight features.

When N best features are eliminated, performance goes down significantly with even smaller N for PBA and BL datasets. Thus, top features are not only effective, they are also crucial for accurate classification, as their discrimination capacity is not replicated by any of the other vocabulary words. This finding is consistent with Lin and Hauptmann (2006) study of perspective vs topic classification: While topical differences between two corpora are manifested in difference in distributions of great many words, they observed little perspective-based variation in distributions of most words, apart from certain words that are preferentially used by adherents of one or the other perspective on the given topic.

For DP and BL-I datasets, the results seem to suggest perspectives with more diffused keyword distribution (No-NBest figures are higher). We note, however, that feature redundancy experiments are confounded in these cases by either a low power of the paired t-test with only 4 pairs (DP) or by a high variance in performance among the 10 folds (BL-I), both of which lead to numerically large discrepancy in performance that is not deemed significant, making it easy to “match” the full set performance with small- N best features as well as without large- N best features. Better comparisons are needed in order to verify the hypothesis of low consolidation.

In future work, we plan to experiment with additional features. For example, Greene and Resnik (2009) reported higher classification accuracies for the DP-4 data using syntactic frames in which a selected group of words appeared, rather than mere presence/absence of the words. Another direction is exploring words as members of semantic fields – while word use might be insufficiently consistent within a perspective, selection of a semantic domain might show better consistency.

References

- Herald Baayen. 2001. *Word frequency distributions*. Dordrecht: Kluwer.
- Marco Baroni and Stefan Evert. 2007. Words and Echoes: Assessing and Mitigating the Non-Randomness Problem in Word Frequency Distribution Modeling. In *Proceedings of the ACL*, pages 904–911, Prague, Czech Republic.
- Suma Bhat and Richard Sproat. 2009. Knowing the Unseen: Estimating Vocabulary Size over Unseen Samples. In *Proceedings of the ACL*, pages 109–117, Suntec, Singapore, August.
- Stephan Greene and Philip Resnik. 2009. More than Words: Syntactic Packaging and Implicit Sentiment. In *Proceedings of HLT-NAACL*, pages 503–511, Boulder, CO, June.
- Mark Hall, Eibe Frank, Geoffrey Holmes, Bernhard Pfahringer, Peter Reutemann, and Ian H. Witten. 2009. The WEKA data mining software: An update. *SIGKDD Explorations*, 11(1).
- Marti Hearst. 1997. TextTiling: Segmenting Text into Multi-Paragraph Subtopic Passages. *Computational Linguistics*, 23(1):33–64.
- Michael Hoey. 1991. *Patterns of Lexis in Text*. Oxford University Press.
- Martin Jansche. 2003. Parametric Models of Linguistic Count Data. In *Proceedings of the ACL*, pages 288–295, Sapporo, Japan, July.
- Thorsten Joachims. 1999. Making large-scale SVM learning practical. In B. Schölkopf, C. Burges, and A. Smola, editors, *Advances in Kernel Methods - Support Vector Learning*. MIT Press.
- Wei-Hao Lin and Alexander Hauptmann. 2006. Are these documents written from different perspectives? A test of different perspectives based on statistical distribution divergence. In *Proceedings of the ACL*, pages 1057–1064, Morristown, NJ, USA.
- Wei-Hao Lin, Theresa Wilson, Janyce Wiebe, and Alexander Hauptmann. 2006. Which side are you on? Identifying perspectives at the document and sentence levels. In *Proceedings of CoNLL*, pages 109–116, Morristown, NJ, USA.
- Andrew McCallum and Kamal Nigam. 1998. A comparison of event models for Naive Bayes text classification. In *Proceedings of AAAI-98 Workshop on Learning for Text Categorization*, pages 41–48, Madison, WI, July.
- Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. 2002. Thumbs up? Sentiment Classification using Machine Learning Techniques. In *Proceedings of EMNLP*, Philadelphia, PA, July.
- Wolfgang Teubert. 2001. A Province of a Federal Superstate, Ruled by an Unelected Bureaucracy – Keywords of the Euro-Sceptic Discourse in Britain. In Andreas Musolff, Colin Good, Petra Points, and Ruth Wittlinger, editors, *Attitudes towards Europe: Language in the unification process*, pages 45–86. Ashgate Publishing Ltd, Hants, England.
- Ian H. Witten and Eibe Frank. 2005. *Data Mining: Practical Machine Learning Tools and Techniques*. Morgan Kaufmann, 2 edition.
- Bei Yu, Stefan Kaufmann, and Daniel Diermeier. 2008. Classifying party affiliation from political speech. *Journal of Information Technology and Politics*, 5(1):33–48.

Cross Lingual Adaptation: An Experiment on Sentiment Classifications

Bin Wei

University of Rochester
Rochester, NY, USA.
bwei@cs.rochester.edu

Christopher Pal

École Polytechnique de Montréal
Montréal, QC, Canada.
christopher.pal@polymtl.ca

Abstract

In this paper, we study the problem of using an annotated corpus in English for the same natural language processing task in another language. While various machine translation systems are available, automated translation is still far from perfect. To minimize the noise introduced by translations, we propose to use only key ‘reliable’ parts from the translations and apply structural correspondence learning (SCL) to find a low dimensional representation shared by the two languages. We perform experiments on an English-Chinese sentiment classification task and compare our results with a previous co-training approach. To alleviate the problem of data sparseness, we create extra pseudo-examples for SCL by making queries to a search engine. Experiments on real-world on-line review data demonstrate the two techniques can effectively improve the performance compared to previous work.

1 Introduction

In this paper we are interested in the problem of transferring knowledge gained from data gathered in one language to another language. A simple and straightforward solution for this problem might be to use automatic machine translations. However, while machine translation has been the subject of a great deal of development in recent years, many of the recent gains in performance manifest as syntactically as opposed to semantically correct sentences. For example, “PIANYI” is a word mainly used in positive comments in Chinese but its translation from the online Google translator is always “cheap”, a word typically used in a negative context in English. To reduce this kind of

error introduced by the translator, Wan in (Wan, 2009) applied a co-training scheme. In this setting classifiers are trained in both languages and the two classifiers teach each other for the unlabeled examples. The co-training approach manages to boost the performance as it allows the text similarity in the target language to compete with the “fake” similarity from the translated texts. However, the translated texts are still used as training data and thus can potentially mislead the classifier. As we are not really interested in predicting something on the language created by the translator, but rather on the real one, it may be better to further diminish the role of the translated texts in the learning process. Motivated by this observation, we suggest here to view this problem as a special case of domain adaptation, in the source domain, we mainly observe English features, while in the other domain mostly features from Chinese. The problem we address is how to associate the features under a unified setting.

There has been a lot of work in domain adaption for NLP (Dai et al., 2007)(Jiang and Zhai, 2007) and one suitable choice for our problem is the approach based on structural correspondence learning (SCL) as in (Blitzer et al., 2006) and (Blitzer et al., 2007b). The key idea of SCL is to identify a low-dimensional representations that capture correspondence between features from both domains (x_s and x_t in our case) by modeling their correlations with some special pivot features. The SCL approach is a good fit for our problem as it performs knowledge transfer through identifying important features. In the cross-lingual setting, we can restrict the translated texts by using them only through the pivot features. We believe this form is more robust to errors in the language produced by the translator.

Adapting language resources and knowledge to a new language was first studied for general text categorization and information retrieval as in (Bel

et al., 2003), where the authors translate a keyword lexicon to perform cross-lingual text categorization. In (Mihalcea et al., 2007), different shortcomings of lexicon-based translation scheme was discussed for the more semantic-oriented task subjective analysis, instead the authors proposed to use a parallel-corpus, apply the classifier in the source language and use the corresponding sentences in the target language to train a new classifier. With the rapid development of automatic machine translations, translating the whole corpus becomes a plausible option. One can either choose to translate a corpus in the target language and apply the classifier in the source language to obtain labeled data, or directly translated the existing data set to the new language. Various experiments of the first strategy are performed in (Banea et al., 2008) for the subjective analysis task and an average 65 F1 score was reported. In (Wan, 2008), the authors propose to combine both strategies with ensemble learning and train a bi-lingual classifier.

In this paper, we are also interested in exploring whether a search engine can be used to improve the performance of NLP systems through reducing the effect of data sparseness. As the SCL algorithm we use here is based on co-occurrence statistics, we adopt a simple approach of creating pseudo-examples from the query counts returned by Google.

2 Our Approach

To begin, we give a formal definition of the problem we are considering. Assume we have two languages l_s and l_t and denote features in these two languages as x_s and x_t respectively. We also have text-level translations and we use $x_{t'}$ for features in the translations from l_s to l_t and $x_{s'}$ for the other direction. Let y be the output variable we want to predict, we have labeled examples (y, x_s) and some unlabeled examples (x_t) . Our task is to train a classifier for (y, x_t) . In this paper, we consider the binary sentiment classification (positive or negative) problem where l_s and l_t correspond to English and Chinese (for general sentiment analysis, we refer the readers to the various previous studies as in (Turney, 2002), (Pang et al., 2002), and (McDonald et al., 2007)). With these definitions in place, we now describe our approach in further detail.

2.1 Structural Correspondence Learning(SCL)

Due to space limitations, we give a very brief overview of the SCL framework here. For a detailed illustration, please refer to (Ando and Zhang, 2005). When SCL is used in a domain adaptation problem, one first needs to find a set of pivot features x_p . These pivot features should behave in a similar manner in both domains, and can be used as “references” to estimate how much other features may contribute when used in a classifier to predict a target variable. These features can either be identified with heuristics (Blitzer et al., 2006) or by automatic selection (Blitzer et al., 2007b). Take sentiment classification as an example, “very good” and “awful” are good pivot features, if a certain feature in the target domain co-occurs often with “very good” but infrequently with “awful”, we could expect this feature will play a similar role as “very good” in the final classifier but a different role from “awful”. We can make this observation purely based on the co-occurrence between these features. No hand-labeling is required and this specific feature doesn’t need to be present in our labeled training data of the source domain.

The SCL approach of (Ando and Zhang, 2005) formulates the above idea by constructing a set of linear predictors for each of the pivot features. Each of these linear predictor is binary like whether “very good” occurs in the text and we have a set of training instances $(1|0, \{x_i\})$. The weight matrix of these linear predictors will encode the co-occurrence statistics between an ordinary feature and the pivot features. As the co-occurrence data are generally very sparse for a typical NLP task, we usually compress the weight matrix using the singular vector decomposition and only selects the top k eigenvectors v_k . This matrix w of the k vectors $\{v_k\}$ gives a mapping from the original feature space to a lower dimensional representation and is shown in (Ando and Zhang, 2005) to be the optimal choice of dimension k under common loss functions. In the next step we can then train a classifier on the extended feature $(x, w * x)$ in the source domain. As w groups the features from different domains with similar behavior relative to the pivot features together, if such a classifier has good performance on the source domain, it will likely do well on the target domain as well.

2.2 SCL for the Cross-lingual Adaptation

Viewing our task as a domain adaptation problem. The source domain correspond to English reviews and the target domain for Chinese ones. The full feature vector is (x_s, x_t) . The difficulty we are facing is, due to noise in the translations, the conditional probabilities $p(y|x_s)$ and the one in the translated texts $p(y|x_{s'})$ may be quite different. Consider the following two straightforward strategies of using automatic machine translations: one can translate the original English labeled data (y, x_s) into $(y, x_{t'})$ in Chinese and train a classifier, or one can train a classifier on (y, x_s) and translate x_t in Chinese into $x_{s'}$ in English so as to use the classifier. But as the conditional distribution can be quite different for the original language and the pseudo language produced by the machine translators, these two strategies give poor performance as reported in (Wan, 2009).

Our solution to this problem is simple: instead of using all the features as $(x_s, x_{t'})$ and $(x_{s'}, x_t)$, we only preserves the pivot features in the translated texts $x_{s'}$ and $x_{t'}$ respectively and discard the other features produced by the translator. So, now we will have (x_s, x_{tp}) and (x_{sp}, x_t) where $x_{(s|t)p}$ are pivot features in the source and the target languages. In other words, when we use the SCL on our problem, the translations are only used to decide if a certain pivot feature occurs or not in the training of the linear predictors. All the other non-pivot features in the translators are blocked to reduce the noise.

In the original SCL as we mentioned earlier, the final classifier is trained on the extended features $(x, w * x)$. However, as mentioned above we will only use the pivot features. To represent this constraint, we can modify the vector to be $(w_p * x, w * x)$ where w_p is a constant matrix that only selects the pivot features. This modification will not affect the deduction procedure and results in (Ando and Zhang, 2005). Experiments show that using only pivot features actually outperforms the full feature setting.

For the selection of the pivot features, we follow the automatic selection method proposed in (Blitzer et al., 2007a). We first select some candidates that occur at least some constant number of times in reviews of the two languages. Then, we rank these features according to their conditional entropy to the labels on the training set. In table 1, we give some of the pivot features with English

English Pivot Features
“poor quality”, “not buy”, “easy use”, “very easy” “excellent”, “perfect”, “still very”, “garbage”, “poor”, “not work”, “not to”, “very comfortable”
Chinese Pivot Features
wanmei(perfect), xiaoguo hen(effect is very...) tisheng(improve), feichang hao(very good), cha(poor), shushi(comfortable), chuse(excellent)

Table 1: Some pivot features.

translations associated with the Chinese pivot features. As we can see from the table, although we only have text-level translations we still get some features with similar meaning from different languages, just like performing an alignment of words.

2.3 Utilizing the Search Engine

Data sparseness is a common problem in NLP tasks. On the other hand, search engines nowadays usually index a huge amount of web pages. We now show how they can also be used as a valuable data source in a less obvious way. Previous studies like (Bollegala, 2007) have shown that search engine results can be comparable to language statistics from a large scale corpus for some NLP tasks like word sense disambiguation. For our problem, we use the query counts returned by a search engine to compute the correlations between a normal feature and the pivot features.

Consider the word “PIANYI” which is mostly used in positive comments, the query “CHAN-PIN(product) PING(comment) CHA(bad) PIANYI” has 2,900,000 results, while “CHAN-PIN(product) PING(comment) HAO(good) PIANYI” returns 57,400,000 pages. The results imply the word “PIANYI” is closer to the pivot feature “good” and it behaves less similar with the pivot feature “bad”.

To add the query counts into the SCL scheme, we create pseudo examples when training linear predictors for pivot features. To construct a pseudo-positive example between a certain feature x_i and a certain pivot feature x_p , we simply query the term $x_i x_p$ and get a count c_1 . We also query x_p alone and get another count c_2 . Then we can create an example $(1, \{0, \dots, 0, x_i = \frac{c_1}{c_2}, 0, \dots, 0\})$. The pseudo-negative examples are created similarly. These pseudo examples are equivalent to texts with a single word and the count is used to

approximate the empirical expectation. As an initial experiment, we select 10,000 Chinese features that occur more than once in the Chinese unlabeled data set but not frequent enough to be captured by the original SCL. And we also select the top 20 most informative Chinese pivot features to perform the queries.

3 Experiment

3.1 Data Set

For comparison, we use the same data set in (Wan, 2009):

Test Set(Labeled Chinese Reviews): The data set contains a total of 886 labeled product reviews in Chinese (451 positive reviews and 435 negative ones). These reviews are extracted from a popular Chinese IT product website IT168¹. The reviews are mainly about electronic devices like mp3 players, mobile phones, digital cameras and computers.

Training Set(Labeled English Reviews): This is the data set used in the domain adaption experiment of (Blitzer et al., 2007b). It contains four major categories: books, DVDs, electronics and kitchen appliances. The data set consists of 8000 reviews with 4000 positive and 4000 negative, It is a public data set available on the web².

Unlabeled Set (Unlabeled Chinese Reviews): 1000 Chinese reviews downloaded from the same website as the Chinese training set. They are of the same domain as the test set.

We translate each English review into Chinese and vice versus through the public Google Translation service. Also following the setting in (Wan, 2009), we only use the Chinese unlabeled data and English training sets for our SCL training procedures. The test set is blind to the training stage.

The features we used are bigrams and unigrams in the two languages as in (Wan, 2009). In Chinese, we first apply the stanford Chinese word segmenter³ to segment the reviews. Bigrams refers to a single Chinese word and a bigram refers to two adjacent Chinese words. The features are also pre-processed and normalized as in (Blitzer et al., 2007b).

¹<http://www.it168.com>

²<http://www.cis.upenn.edu/~mdredze/datasets/sentiment/>

³<http://nlp.stanford.edu/software/segmenter.shtml>

Models	Precision	Recall	F-Score
CoTrain	0.768	0.905	0.831
SCL-B	0.772	0.914	0.837
SCL-C	0.764	0.896	0.825
SCL-O	0.760	0.909	0.828
SCL-E	0.801	0.909	0.851

Table 2: Results on the Positive Reviews

Models	Precision	Recall	F-Score
CoTrain	0.879	0.717	0.790
SCL-B	0.931	0.752	0.833
SCL-C	0.908	0.743	0.817
SCL-O	0.928	0.739	0.823
SCL-E	0.928	0.796	0.857

Table 3: Results on the Negative Reviews

3.2 Comparisons

We compare our procedure with the co-training scheme reported in (Wan, 2009):

CoTrain: The method with the best performance in (Wan, 2009). Two standard SVMs are trained using the co-training scheme for the Chinese views and the English views. And the results of the two SVMs are combined to give the final output.

SCL-B: The basic SCL procedure as explained.

SCL-O: The basic SCL except that we use all features from the translated texts instead of only the pivot features.

SCL-C: The training procedure is still the same as **SCL-B** except in the test time we only use the Chinese pivot features and neglect the English pivot features from translations.

SCL-E: The same as **SCL-B** except that in the training of linear pivot predictors, we also use the pseudo examples constructed from queries of the search engine.

Table 2 and 3 give results measured on the positive labeled reviews and negative reviews separately. Table 4 gives the overall accuracy on the whole 886 reviews. Our basic SCL approach **SCL-B** outperforms the original **Co-Training** approach by 2.2% in the overall accuracy. We can

CoTrain	SCL-B	SCL-O	SCL-C	SCL-E
0.813	0.835	0.826	0.822	0.854

Table 4: Overall Accuracy of Different Methods

also notice that using all the features including the ones from translations actually deteriorate the performance from 0.835 to 0.826.

The model incorporating the co-occurrence count information from the search engine has the best overall performance of 0.857. It is interesting to note that the simple scheme we have adopted increased the recall performance on the negative reviews significantly. After examining the reviews, we find the negative part contains some idioms and words mainly used on the internet and the query count seems to be able to capture their usage.

Finally, as our final goal is to train a Chinese sentiment classifier, it will be best if our model can only rely on the Chinese features. The **SCL-C** model improves the performance from the **Co-Training** method a little but not as much as the *SCL - B* and the *SCL - O* approaches. This observation suggests that the translations are still helpful for the cross-lingual adaptation problem as the translators perform some implicit semantic mapping.

4 Conclusion

In this paper, we are interested in adapting existing knowledge to a new language. We show that instead of fully relying on automatic translation, which may be misleading for a highly semantic task like the sentiment analysis, using techniques like SCL to connect the two languages through feature-level mapping seems a more suitable choice. We also perform an initial experiment using the co-occurrence statistics from a search engine to handle the data sparseness problem in the adaptation process, and the result is encouraging.

As future research we believe a promising avenue of exploration is to construct a probabilistic version of the SCL approach which could offer a more explicit model of the relations between the two domains and the relations between the search engine results and the model parameters. Also, in the current work, we select the pivot features by simple ranking with mutual information, which only considers the distribution information. Incorporating the confidence from the translator may further improve the performance.

References

Rie Kubota Ando and Tong Zhang. 2005. A framework for learning predictive structures from multi-

ple tasks and unlabeled data. *Journal of Machine Learning Research*.

Carmen Banea, Rada Mihalcea, Janyce Wiebe, and Samer Hassan. 2008. Multilingual subjectivity analysis using machine translation. In *Proceedings of EMNLP*.

Nuria Bel, Cornelis H. A. Koster, and Marta Villegas. 2003. Cross-lingual text categorization. In *Research and Advanced Technology for Digital Libraries*.

John Blitzer, Ryan McDonald, and Fernando Pereira. 2006. Domain adaptation with structural correspondence learning. In *Proceedings of EMNLP*.

John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman. 2007a. Learning bounds for domain adaptation. In *Proceedings of NIPS*.

John Blitzer, Mark Dredze, and Fernando Pereira. 2007b. Biographies, bollywood, boom-boxes and blenders: Domain adaptation for sentiment classification. In *Proceedings of ACL*.

Danushka Bollegala. 2007. Measuring semantic similarity between words using web search engines. In *Proceedings of WWW 07*.

Wenyuan Dai, Gui-Rong Xue, Qiang Yang, and Yong Yu. 2007. Co-clustering based classification for out-of-domain documents. In *Proceedings of KDD*.

Jing Jiang and ChengXiang Zhai. 2007. A two-stage approach to domain adaptation for statistical classifiers. In *Proceedings of CIKM*.

Ryan T. McDonald, Kerry Hannan, Tyler Neylon, Mike Wells, and Jeffrey C. Reynar. 2007. Structured models for fine-to-coarse sentiment analysis. In *Proceedings of ACL*.

Rada Mihalcea, Carmen Banea, and Janyce Wiebe. 2007. Learning multilingual subjective language via cross-lingual projections. In *Proceedings of ACL*.

Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. 2002. Thumbs up? sentiment classification using machine learning techniques. In *Proceedings of EMNLP*.

Peter D. Turney. 2002. Thumbs up or thumbs down? semantic orientation applied to unsupervised classification of reviews. In *Proceedings of ACL*.

Xiaojun Wan. 2008. Using bilingual knowledge and ensemble techniques for unsupervised chinese sentiment analysis. In *Proceedings of EMNLP*.

Xiaojun Wan. 2009. Co-training for cross-lingual sentiment classification. In *Proceedings of ACL*.

Using Anaphora Resolution to Improve Opinion Target Identification in Movie Reviews

Niklas Jakob

Technische Universität Darmstadt
Hochschulstraße 10, 64289 Darmstadt

<http://www.ukp.tu-darmstadt.de/people>

Iryna Gurevych

Technische Universität Darmstadt
Hochschulstraße 10, 64289 Darmstadt

Abstract

Current work on automatic opinion mining has ignored opinion targets expressed by anaphorical pronouns, thereby missing a significant number of opinion targets. In this paper we empirically evaluate whether using an off-the-shelf anaphora resolution algorithm can improve the performance of a baseline opinion mining system. We present an analysis based on two different anaphora resolution systems. Our experiments on a movie review corpus demonstrate, that an unsupervised anaphora resolution algorithm significantly improves the opinion target extraction. We furthermore suggest domain and task specific extensions to an off-the-shelf algorithm which in turn yield significant improvements.

1 Introduction

Over the last years the task of opinion mining (OM) has been the topic of many publications. It has been approached with different goals in mind: Some research strived to perform subjectivity analysis at the document or sentence level, without focusing on what the individual opinions uttered in the document are about. Other approaches focused on extracting individual opinion words or phrases and what they are about. This aboutness has been referred to as the *opinion target* or opinion topic in the literature from the field. In this work our goal is to extract *opinion target* - *opinion word* pairs from sentences from movie reviews. A challenge which is frequently encountered in text mining tasks at this level of granularity is, that entities are being referred to by anaphora. In the task of OM, it can therefore also be necessary to analyze more than the content of one individual sentence when extracting opinion targets. Consider this example sentence: “*Simply*

put, it’s unfathomable that this movie cracks the Top 250. It is absolutely awful.”. If one wants to extract what the opinion in the second sentence is about, an algorithm which resolves the anaphoric reference to the opinion target is required.

The extraction of such anaphoric opinion targets has been noted as an open issue multiple times in the OM context (Zhuang et al., 2006; Hu and Liu, 2004; Nasukawa and Yi, 2003). It is not a marginal phenomenon, since Kessler and Nicolov (2009) report that in their data, 14% of the opinion targets are pronouns. However, the task of resolving anaphora to mine opinion targets has not been addressed and evaluated yet to the best of our knowledge.

In this work, we investigate whether anaphora resolution (AR) can be successfully integrated into an OM algorithm and whether we can achieve an improvement regarding the OM in doing so. This paper is structured as follows: Section 2 discusses the related work on opinion target identification and OM on movie reviews. Section 3 outlines the OM algorithm we employed by us, while in Section 4 we discuss two different algorithms for AR which we experiment with. Finally, in Section 5 we present our experimental work including error analysis and discussion, and we conclude in Section 6.

2 Related Work

We split the description of the related work in two parts: In Section 2.1 we discuss the related work on OM with a focus on approaches for opinion target identification. In Section 2.2 we elaborate on findings from related OM research which also worked with movie reviews as this is our target domain in the present paper.

2.1 Opinion Target Identification

The extraction of opinions and especially opinion targets has been performed with quite diverse

approaches. Initial approaches combined statistical information and basic linguistic features such as part-of-speech tags. The goal was to identify the opinion targets, here in form of products and their attributes, without a pre-built knowledge base which models the domain. For the target candidate identification, simple part-of-speech patterns were employed. The relevance ranking and extraction was then performed with different statistical measures: Pointwise Mutual Information (Popescu and Etzioni, 2005), the Likelihood Ratio Test (Yi et al., 2003) and Association Mining (Hu and Liu, 2004). A more linguistically motivated approach was taken by Kim and Hovy (2006) through identifying opinion holders and targets with semantic role labeling. This approach was promising, since their goal was to extract opinions from professionally edited content i.e. newswire.

Zhuang et al. (2006) present an algorithm for the extraction of *opinion target - opinion word* pairs. The opinion word and target candidates are identified in the annotated corpus and their extraction is then performed by applying possible paths connecting them in a dependency graph. These paths are combined with part-of-speech information and also learned from the annotated corpus.

To the best of our knowledge, there is currently only one system which integrates coreference information in OM. The algorithm by Stoyanov and Cardie (2008) identifies coreferring targets in newspaper articles. A candidate selection or extraction step for the opinion targets is not required, since they rely on manually annotated targets and focus solely on the coreference resolution. However they do not resolve pronominal anaphora in order to achieve that.

2.2 Opinion Mining on Movie Reviews

There is a huge body of work on OM in movie reviews which was sparked by the dataset from Pang and Lee (2005). This dataset consists of sentences which are annotated as expressing positive or negative opinions. An interesting insight was gained from the document level sentiment analysis on movie reviews in comparison to documents from other domains: Turney (2002) observes that the movie reviews are hardest to classify since the review authors tend to give information about the storyline of the movie which often contain characterizations, such as “*bad guy*” or “*violent scene*”. These statements however do not reflect any opin-

ions of the reviewers regarding the movie. Zhuang et al. (2006) also observe that movie reviews are different from e.g. customer reviews on Amazon.com. This is reflected in their experiments, in which their system outperforms the system by Hu and Liu (2004) which attributes an opinion target to the opinion word which is closest regarding word distance in a sentence. The sentences in the movie reviews tend to be more complex, which can also be explained by their origin. The reviews were taken from the Internet Movie Database¹, on which the users are given a set of guidelines on how to write a review. Due to these insights, we are confident that the overall textual quality of the movie reviews is high enough for linguistically more advanced technologies such as parsing or AR to be successfully applied.

3 Opinion Target Identification

3.1 Dataset

Currently the only freely available dataset annotated with opinions including annotated anaphoric opinion targets is a corpus of movie reviews by Zhuang et al. (2006). Kessler and Nicolov (2009) describe a collection of product reviews in which anaphoric opinion targets are also annotated, but it is not available to the public (yet). Zhuang et al. (2006) used a subset of the dataset they published (1829 documents), namely 1100 documents, however they do not state which documents comprise this subset used in their evaluation. In our experiments, we therefore use the complete dataset available, detailed in Table 1. As shown, roughly 9.5% of the opinion targets are referred to by pronouns. Table 2 outlines detailed statistics on which pronouns occur as opinion targets.

Table 1: Dataset Statistics

# Documents	1829
# Sentences	24918
# Tokens	273715
# Target + Opinion Pairs	5298
# Targets which are Pronouns	504
# Pronouns	> 11000

3.2 Baseline Opinion Mining

We reimplemented the algorithm presented by Zhuang et al. (2006) as the baseline for our

¹<http://www.imdb.com> (IMDB)

Table 2: Pronouns as Opinion Targets

it	274	he	58	she	22	they	22
this	77	his	26	her	10		
		him	15				

experiments. Their approach is a supervised one. The annotated dataset is split in five folds, of which four are used as the training data. In the first step, opinion target and opinion word candidates are extracted from the training data. Frequency counts of the annotated opinion targets and opinion words are extracted from four training folds. The most frequently occurring opinion targets and opinion words are selected as candidates. Then the annotated sentences are parsed and a graph containing the words of the sentence is created, which are connected by the dependency relations between them. For each *opinion target - opinion word* pair, the shortest path connecting them is extracted from the dependency graph. A path consists of the part-of-speech tags of the nodes and the dependency types of the edges.

In order to be able to identify rarely occurring opinion targets which are not in the candidate list, they expand it by crawling the cast and crew names of the movies from the IMDB. How this crawling and extraction is done is not explained.

4 Algorithms for Anaphora Resolution

As pointed out by Charniak and Elsner (2009) there are hardly any freely available systems for AR. Although Charniak and Elsner (2009) present a machine-learning based algorithm for AR, they evaluate its performance in comparison to three non machine-learning based algorithms, since those are the only ones available. They observe that the best performing baseline algorithm (OpenNLP) is hardly documented. The algorithm with the next-to-highest results in (Charniak and Elsner, 2009) is MARS (Mitkov, 1998) from the GuiTAR (Poesio and Kabadjov, 2004) toolkit. This algorithm is based on statistical analysis of the antecedent candidates. Another promising algorithm for AR employs a rule based approach for antecedent identification. The CogNIAC algorithm (Baldwin, 1997) was designed for high-precision AR. This approach seems like an adequate strategy for our OM task, since in the dataset used in our experiments only a small fraction of the total number of pronouns are ac-

tual opinion targets (see Table 1). We extended the CogNIAC implementation to also resolve “*it*” and “*this*” as anaphora candidates, since off-the-shelf it only resolves personal pronouns. We will refer to this extension with [id]. Both algorithms follow the common approach that noun phrases are antecedent candidates for the anaphora. In our experiments we employed both the MARS and the CogNIAC algorithm, for which we created three extensions which are detailed in the following.

4.1 Extensions of CogNIAC

We identified a few typical sources of errors in a preliminary error analysis. We therefore suggest three extensions to the algorithm which are on the one hand possible in the OM setting and on the other hand represent special features of the target discourse type: [1.] We observed that the Stanford Named Entity Recognizer (Finkel et al., 2005) is superior to the *Person* detection of the (MUC6 trained) CogNIAC implementation. We therefore filter out *Person* antecedent candidates which the Stanford NER detects for the impersonal and demonstrative pronouns and *Location & Organization* candidates for the personal pronouns. This way the input to the AR is optimized. [2.] The second extension exploits the fact that reviews from the IMDB exhibit certain contextual properties. They are gathered and to be presented in the context of one particular entity (=movie). The context or topic under which it occurs is therefore typically clear to the reader and is therefore not explicitly introduced in the discourse. This is equivalent to the situational context we often refer to in dialogue. In the reviews, the authors often refer to the movie or film as a whole by a pronoun. We exploit this by an additional rule which resolves an impersonal or demonstrative pronoun to “*movie*” or “*film*” if there is no other (matching) antecedent candidate in the previous two sentences. [3.] The rules by which CogNIAC resolves anaphora were designed so that anaphora which have ambiguous antecedents are left unresolved. This strategy should lead to a high precision AR, but at the same time it can have a negative impact on the recall. In the OM context, it happens quite frequently that the authors comment on the entity they want to criticize in a series of arguments. In such argument chains, we try to solve cases of antecedent ambiguity by analyzing the opinions: If there are ambiguous antecedent candidates for a

pronoun, we check whether there is an opinion uttered in the previous sentence. If this is the case and if the opinion target matches the pronoun regarding gender and number, we resolve the pronoun to the antecedent which was the previous opinion target.

In the results of our experiments in Section 5, we will refer to the configurations using these extensions with the numbers attributed to them above.

5 Experimental Work

To integrate AR in the OM algorithm, we add the antecedents of the pronouns annotated as opinion targets to the target candidate list. Then we extract the dependency paths connecting pronouns and opinion words and add them to the list of valid paths. When we run the algorithm, we extract anaphora which were resolved, if they occur with a valid dependency path to an opinion word. In such a case, the anaphor is substituted for its antecedent and thus extracted as part of an *opinion target - opinion word* pair.

To reproduce the system by Zhuang et al. (2006), we substitute the cast and crew list employed by them (see Section 3.2), with a NER component (Finkel et al., 2005). One aspect regarding the extraction of *opinion target - opinion word* pairs remains open in Zhuang et al. (2006): The dependency paths only identify connections between pairs of single words. However, almost 50% of the opinion target candidates are multiword expressions. Zhuang et al. (2006) do not explain how they extract multiword opinion targets with the dependency paths. In our experiments, we require a dependency path to be found to each word of a multiword target candidate for it to be extracted. Furthermore, Zhuang et al. (2006) do not state whether in their evaluation annotated multiword targets are treated as a single unit which needs to be extracted, or whether a partial matching is employed in such cases. We require all individual words of a multiword expression to be extracted by the algorithm. As mentioned above, the dependency path based approach will only identify connections between pairs of single words. We therefore employ a merging step, in which we combine adjacent opinion targets to a multiword expression. We have compiled two result sets: Table 3 shows the results of the overall OM in a five-fold cross-validation. Table 4 gives a detailed overview of the AR for opinion target identification summed

up over all folds. In Table 4, a true positive refers to an extracted pronoun which was annotated as an opinion target and is resolved to the correct antecedent. A false positive subsumes two error classes: A pronoun which was not annotated as an opinion target but extracted as such, or a pronoun which is resolved to an incorrect antecedent.

As shown in Table 3, the recall of our reimplementation is slightly higher than the recall reported in Zhuang et al. (2006). However, our precision and thus f-measure are lower. This can be attributed to the different document sets used in our experiments (see Section 3.1), or our substitution of the list of peoples' names with the NER component, or differences regarding the evaluation strategy as mentioned above.

We observe that the MARS algorithm yields an improvement regarding recall compared to the baseline system. However, it also extracts a high number of false positives for both the personal and impersonal / demonstrative pronouns. This is due to the fact that the MARS algorithm is designed for robustness and always resolves a pronoun to an antecedent.

CogNIAC in its off-the-shelf configuration already yields significant improvements over the baseline regarding f-measure². Our CogNIAC extension [id] improves recall slightly in comparison to the off-the-shelf system. As shown in Table 4, the algorithm extracts impersonal and demonstrative pronouns with lower precision than personal pronouns. Our error analysis shows that this is mostly due to the *Person / Location / Organization* classification of the CogNIAC implementation. The names of actors and movies are thus often misclassified. Extension [1] mitigates this problem, since it increases precision (Table 3 row 6), while not affecting recall. The overall improvement of our extensions [id] + [1] is however not statistically significant in comparison to off-the-shelf CogNIAC. Our extensions [2] and [3] in combination with [id] each increase recall at the expense of precision. The improvement in f-measure of CogNIAC [id] + [3] over the off-the-shelf system is statistically significant. The best overall results regarding f-measure are reached if we combine all our extensions of the CogNIAC algorithm. The results of this configuration show that the positive effects of extensions [2] and [3] are complemen-

²Significance of improvements was tested using a paired two-tailed t-test and $p \leq 0.05$ (*) and $p \leq 0.01$ (**)

Table 3: Op. Target - Op. Word Pair Extraction

Configuration	Reca.	Prec.	F-Meas.
Results in Zhuang et al.	0.548	0.654	0.596
Our Reimplementation	0.554	0.523	0.538
MARS off-the-shelf	0.595	0.467	0.523
CogNIAC off-the-shelf	0.586	0.534	0.559**
CogNIAC+[id]	0.594	0.516	0.552
CogNIAC+[id]+[1]	0.594	0.533	0.561
CogNIAC+[id]+[2]	0.603	0.501	0.547
CogNIAC+[id]+[3]	0.613	0.521	0.563*
CogNIAC+[id]+[1]+[2]+[3]	0.614	0.531	0.569*

Table 4: Results of AR for Opinion Targets

Algorithm	Pers. ¹		Imp. & Dem. ¹	
	TP ²	FP ²	TP	FP
MARS off-the-shelf	102	164	115	623
CogNIAC off-the-shelf	117	95	0	0
CogNIAC+[id]	117	95	105	180
CogNIAC+[id]+[1]	117	41	105	51
CogNIAC+[id]+[2]	117	95	153	410
CogNIAC+[id]+[3]	131	103	182	206
CogNIAC+[id]+[1]+[2]+[3]	124	64	194	132

¹ personal, impersonal & demonstrative pronouns

² true positives, false positives

tary regarding the extraction of impersonal and demonstrative pronouns. This configuration yields statistically significant improvements regarding f-measure over the off-the-shelf CogNIAC configuration, while also having the overall highest recall.

5.1 Error Analysis

When extracting opinions from movie reviews, we observe the same challenge as Turney (2002): The users often characterize events in the storyline or roles the characters play. These characterizations contain the same words which are also used to express opinions. Hence these combinations are frequently but falsely extracted as *opinion target* - *opinion word* pairs, negatively affecting the precision. The algorithm cannot distinguish them from opinions expressing the stance of the author. Overall, the recall of the baseline is rather low. This is due to the fact that the algorithm only learns a subset of the opinion words and opinion targets annotated in the training data. Currently, it cannot discover any new opinion words and targets. This could be addressed by integrating a component which identifies new opinion targets by calculating the relevance of a word in the corpus based on statistical measures.

The AR introduces new sources of errors regarding the extraction of opinion targets: Errors in

gender and number identification can lead to an incorrect selection of antecedent candidates. Even if the gender and number identification is correct, the algorithm might select an incorrect antecedent if there is more than one possible candidate. A non-robust algorithm as CogNIAC might leave a pronoun which is an actual opinion target unresolved, due to the ambiguity of its antecedent candidates.

The upper bound for the OM with perfect AR on top of the baseline would be recall: 0.649, precision: 0.562, f-measure: 0.602. Our best configuration reaches $\sim 50\%$ of the improvements which are theoretically possible with perfect AR.

6 Conclusions

We have shown that by extending an OM algorithm with AR for opinion target extraction significant improvements can be achieved. The rule based AR algorithm CogNIAC performs well regarding the extraction of opinion targets which are personal pronouns. The algorithm does not yield high precision when resolving impersonal and demonstrative pronouns. We present a set of extensions which address this challenge and in combination yield significant improvements over the off-the-shelf configuration. A robust AR algorithm does not yield any improvements regarding f-measure in the OM task. This type of algorithm creates many false positives, which are not filtered out by the dependency paths employed in the algorithm by Zhuang et al. (2006).

AR could also be employed in other OM algorithms which aim at identifying opinion targets by means of a statistical analysis. Vicedo and Ferrández (2000) successfully modified the relevance ranking of terms in their documents by replacing anaphora with their antecedents. The approach can be taken for OM algorithms which select the opinion target candidates with a relevance ranking (Hu and Liu, 2004; Yi et al., 2003).

Acknowledgments

The project was funded by means of the German Federal Ministry of Economy and Technology under the promotional reference "01MQ07012". The authors take the responsibility for the contents. This work has been supported by the Volkswagen Foundation as part of the Lichtenberg-Professorship Program under grant No. I/82806.

References

- Breck Baldwin. 1997. Cogniac: High precision coreference with limited knowledge and linguistic resources. In *Proceedings of a Workshop on Operational Factors in Practical, Robust Anaphora Resolution for Unrestricted Texts*, pages 38–45, Madrid, Spain, July.
- Eugene Charniak and Micha Elsner. 2009. EM works for pronoun anaphora resolution. In *Proceedings of the 12th Conference of the European Chapter of the ACL*, pages 148–156, Athens, Greece, March.
- Jenny Rose Finkel, Trond Grenager, and Christopher Manning. 2005. Incorporating non-local information into information extraction systems by gibbs sampling. In *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics*, pages 363–370, Michigan, USA, June.
- Minqing Hu and Bing Liu. 2004. Mining and summarizing customer reviews. In *Proceedings of the 10th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 168–177, Seattle, WA, USA, August.
- Jason Kessler and Nicolas Nicolov. 2009. Targeting sentiment expressions through supervised ranking of linguistic configurations. In *Proceedings of the Third International AAI Conference on Weblogs and Social Media*, San Jose, CA, USA, May.
- Soo-Min Kim and Eduard Hovy. 2006. Extracting opinions, opinion holders, and topics expressed in online news media text. In *Proceedings of the ACL Workshop on Sentiment and Subjectivity in Text*, pages 1–8, Sydney, Australia, July.
- Ruslan Mitkov. 1998. Robust pronoun resolution with limited knowledge. In *Proceedings of the 36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics*, pages 869–875, Montreal, Canada, August.
- Tetsuya Nasukawa and Jeonghee Yi. 2003. Sentiment analysis: Capturing favorability using natural language processing. In *Proceedings of the 2nd International Conference on Knowledge Capture*, pages 70–77, Sanibel Island, FL, USA, October.
- Bo Pang and Lillian Lee. 2005. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics*, pages 115–124, Michigan, USA, June.
- Massimo Poesio and Mijail A. Kabadjov. 2004. A general-purpose, off-the-shelf anaphora resolution module: Implementation and preliminary evaluation. In *Proceedings of the 4th International Conference on Language Resources and Evaluation*, pages 663–666, Lisboa, Portugal, May.
- Ana-Maria Popescu and Oren Etzioni. 2005. Extracting product features and opinions from reviews. In *Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing*, pages 339–346, Vancouver, Canada, October.
- Veselin Stoyanov and Claire Cardie. 2008. Topic identification for fine-grained opinion analysis. In *Proceedings of the 22nd International Conference on Computational Linguistics*, pages 817–824, Manchester, UK, August.
- Peter Turney. 2002. Thumbs up or thumbs down? semantic orientation applied to unsupervised classification of reviews. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 417–424, Philadelphia, Pennsylvania, USA, July.
- José L. Vicedo and Antonio Ferrández. 2000. Applying anaphora resolution to question answering and information retrieval systems. In *Proceedings of the First International Conference on Web-Age Information Management*, volume 1846 of *Lecture Notes In Computer Science*, pages 344–355. Springer, Shanghai, China.
- Jeonghee Yi, Tetsuya Nasukawa, Razvan Bunescu, and Wayne Niblack. 2003. Sentiment analyzer: Extracting sentiments about a given topic using natural language processing techniques. In *Proceedings of the 3rd IEEE International Conference on Data Mining*, pages 427–434, Melbourne, FL, USA, December.
- Li Zhuang, Feng Jing, and Xiao-Yan Zhu. 2006. Movie review mining and summarization. In *Proceedings of the ACM 15th Conference on Information and Knowledge Management*, pages 43–50, Arlington, VA, USA, November.

Hierarchical Sequential Learning for Extracting Opinions and their Attributes

Yejin Choi and Claire Cardie

Department of Computer Science

Cornell University

Ithaca, NY 14853

{ychoi, cardie}@cs.cornell.edu

Abstract

Automatic opinion recognition involves a number of related tasks, such as identifying the boundaries of opinion expression, determining their polarity, and determining their intensity. Although much progress has been made in this area, existing research typically treats each of the above tasks in isolation. In this paper, we apply a hierarchical parameter sharing technique using Conditional Random Fields for fine-grained opinion analysis, jointly detecting the boundaries of opinion expressions as well as determining two of their key attributes — polarity and intensity. Our experimental results show that our proposed approach improves the performance over a baseline that does not exploit hierarchical structure among the classes. In addition, we find that the joint approach outperforms a baseline that is based on cascading two separate components.

1 Introduction

Automatic opinion recognition involves a number of related tasks, such as identifying expressions of opinion (e.g. Kim and Hovy (2005), Popescu and Etzioni (2005), Breck et al. (2007)), determining their polarity (e.g. Hu and Liu (2004), Kim and Hovy (2004), Wilson et al. (2005)), and determining their strength, or intensity (e.g. Popescu and Etzioni (2005), Wilson et al. (2006)). Most previous work treats each subtask in isolation: opinion expression extraction (i.e. detecting the boundaries of opinion expressions) and opinion attribute classification (e.g. determining values for polarity and intensity) are tackled as separate steps in opinion recognition systems. Unfortunately, errors from individual components will propagate in

systems with cascaded component architectures, causing performance degradation in the end-to-end system (e.g. Finkel et al. (2006)) — in our case, in the end-to-end opinion recognition system.

In this paper, we apply a *hierarchical parameter sharing* technique (e.g., Cai and Hofmann (2004), Zhao et al. (2008)) using Conditional Random Fields (CRFs) (Lafferty et al., 2001) to fine-grained opinion analysis. In particular, we aim to jointly identify the boundaries of opinion expressions as well as to determine two of their key attributes — polarity and intensity.

Experimental results show that our proposed approach improves the performance over the baseline that does not exploit the hierarchical structure among the classes. In addition, we find that the joint approach outperforms a baseline that is based on cascading two separate systems.

2 Hierarchical Sequential Learning

We define the problem of joint extraction of opinion expressions and their attributes as a sequence tagging task as follows. Given a sequence of tokens, $x = x_1 \dots x_n$, we predict a sequence of labels, $y = y_1 \dots y_n$, where $y_i \in \{0, \dots, 9\}$ are defined as conjunctive values of polarity labels and intensity labels, as shown in Table 1. Then the conditional probability $p(y|x)$ for linear-chain CRFs is given as (Lafferty et al., 2001)

$$P(y|x) = \frac{1}{Z_x} \exp \sum_i \left(\lambda f(y_i, x, i) + \lambda' f'(y_{i-1}, y_i, x, i) \right)$$

where Z_x is the normalization factor.

In order to apply a hierarchical parameter sharing technique (e.g., Cai and Hofmann (2004), Zhao et al. (2008)), we extend parameters as follows.

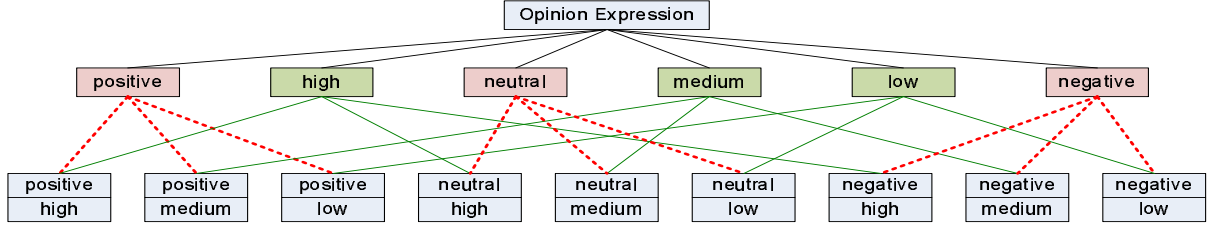


Figure 1: The hierarchical structure of classes for opinion expressions with polarity (positive, neutral, negative) and intensity (high, medium, low)

LABEL	0	1	2	3	4	5	6	7	8	9
POLARITY	none	positive	positive	positive	neutral	neutral	neutral	negative	negative	negative
INTENSITY	none	high	medium	low	high	medium	low	high	medium	low

Table 1: Labels for Opinion Extraction with Polarity and Intensity

$$\lambda f(y_i, x, i) = \lambda_\alpha g_o(\alpha, x, i) + \lambda_\beta g_p(\beta, x, i) + \lambda_\gamma g_s(\gamma, x, i) \quad (1)$$

$$\lambda' f'(y_{i-1}, y_i, x, i) = \lambda'_{\alpha, \hat{\alpha}} g'_o(\alpha, \hat{\alpha}, x, i) + \lambda'_{\beta, \hat{\beta}} g'_p(\beta, \hat{\beta}, x, i) + \lambda'_{\gamma, \hat{\gamma}} g'_s(\gamma, \hat{\gamma}, x, i)$$

where g_o and g'_o are feature vectors defined for **Opinion** extraction, g_p and g'_p are feature vectors defined for **Polarity** extraction, and g_s and g'_s are feature vectors defined for **Strength** extraction, and

$$\begin{aligned} \alpha, \hat{\alpha} &\in \{\text{OPINION, NO-OPINION}\} \\ \beta, \hat{\beta} &\in \{\text{POSITIVE, NEGATIVE, NEUTRAL, NO-POLARITY}\} \\ \gamma, \hat{\gamma} &\in \{\text{HIGH, MEDIUM, LOW, NO-INTENSITY}\} \end{aligned}$$

For instance, if $y_i = 1$, then

$$\begin{aligned} \lambda f(1, x, i) &= \lambda_{\text{OPINION}} g_o(\text{OPINION}, x, i) \\ &+ \lambda_{\text{POSITIVE}} g_p(\text{POSITIVE}, x, i) \\ &+ \lambda_{\text{HIGH}} g_s(\text{HIGH}, x, i) \end{aligned}$$

If $y_{i-1} = 0, y_i = 4$, then

$$\begin{aligned} \lambda' f'(0, 4, x, i) &= \lambda'_{\text{NO-OPINION, OPINION}} g'_o(\text{NO-OPINION, OPINION}, x, i) \\ &+ \lambda'_{\text{NO-POLARITY, NEUTRAL}} g'_p(\text{NO-POLARITY, NEUTRAL}, x, i) \\ &+ \lambda'_{\text{NO-INTENSITY, HIGH}} g'_s(\text{NO-INTENSITY, HIGH}, x, i) \end{aligned}$$

This hierarchical construction of feature and weight vectors allows similar labels to share the same subcomponents of feature and weight vectors. For instance, all $\lambda f(y_i, x, i)$ such that

$y_i \in \{1, 2, 3\}$ will share the same component $\lambda_{\text{POSITIVE}} g_p(\text{POSITIVE}, x, i)$. Note that there can be other variations of hierarchical construction. For instance, one can add $\lambda_\delta g_t(\delta, x, i)$ and $\lambda'_{\delta, \hat{\delta}} g'_t(\delta, \hat{\delta}, x, i)$ to Equation (1) for $\delta \in \{0, 1, \dots, 9\}$, in order to allow more individualized learning for each label.

Notice also that the number of sets of parameters constructed by Equation (1) is significantly smaller than the number of sets of parameters that are needed without the hierarchy. The former requires $(2 + 4 + 4) + (2 \times 2 + 4 \times 4 + 4 \times 4) = 46$ sets of parameters, but the latter requires $(10) + (10 \times 10) = 110$ sets of parameters. Because a combination of a polarity component and an intensity component can distinguish each label, it is not necessary to define a separate set of parameters for each label.

3 Features

We first introduce definitions of key terms that will be used to describe features.

- **PRIOR-POLARITY & PRIOR-INTENSITY:**

We obtain these prior-attributes from the *polarity lexicon* populated by Wilson et al. (2005).

- **EXP-POLARITY, EXP-INTENSITY & EXP-SPAN:**

Words in a given opinion expression often do not share the same prior-attributes. Such discontinuous distribution of features can make it harder to learn the desired opinion expression boundaries. Therefore, we try to obtain expression-level attributes (EXP-POLARITY and EXP-INTENSITY) using simple heuristics. In order to derive EXP-POLARITY, we perform simple

voting. If there is a word with a negation effect, such as “never”, “not”, “hardly”, “against”, then we flip the polarity. For EXP-INTENSITY, we use the highest PRIOR-INTENSITY in the span. The text span with the same expression-level attributes are referred to as EXP-SPAN.

3.1 Per-Token Features

Per-token features are defined in the form of $g_o(\alpha, x, i)$, $g_p(\beta, x, i)$ and $g_s(\gamma, x, i)$. The domains of α, β, γ are as given in Section 3.

Common Per-Token Features

Following features are common for all class labels. The notation \otimes indicates conjunctive operation of two values.

- PART-OF-SPEECH(x_i):
based on GATE (Cunningham et al., 2002).
- WORD(x_i), WORD(x_{i-1}), WORD(x_{i+1})
- WORDNET-HYPERNYM(x_i):
based on WordNet (Miller, 1995).
- OPINION-LEXICON(x_i):
based on *opinion lexicon* (Wiebe et al., 2002).
- SHALLOW-PARSER(x_i):
based on CASS partial parser (Abney, 1996).
- PRIOR-POLARITY(x_i) \otimes PRIOR-INTENSITY(x_i)
- EXP-POLARITY(x_i) \otimes EXP-INTENSITY(x_i)
- EXP-POLARITY(x_i) \otimes EXP-INTENSITY(x_i) \otimes STEM(x_i)
- EXP-SPAN(x_i):
boolean to indicate whether x_i is in an EXP-SPAN.
- DISTANCE-TO-EXP-SPAN(x_i): 0, 1, 2, 3+.
- EXP-POLARITY(x_i) \otimes EXP-INTENSITY(x_i) \otimes EXP-SPAN(x_i)

Polarity Per-Token Features

These features are included only for $g_o(\alpha, x, i)$ and $g_p(\beta, x, i)$, which are the feature functions corresponding to the polarity-based classes.

- PRIOR-POLARITY(x_i), EXP-POLARITY(x_i)
- STEM(x_i) \otimes EXP-POLARITY(x_i)
- COUNT-OF-Polarity:
where *Polarity* \in {positive, neutral, negative}.
This feature encodes the number of positive, neutral, and negative EXP-POLARITY words respectively, in the current sentence.
- STEM(x_i) \otimes COUNT-OF-Polarity
- EXP-POLARITY(x_i) \otimes COUNT-OF-Polarity
- EXP-SPAN(x_i) and EXP-POLARITY(x_i)
- DISTANCE-TO-EXP-SPAN(x_i) \otimes EXP-POLARITY(x_p)

Intensity Per-Token Features

These features are included only for $g_o(\alpha, x, i)$ and $g_s(\gamma, x, i)$, which are the feature functions corresponding to the intensity-based classes.

- PRIOR-INTENSITY(x_i), EXP-INTENSITY(x_i)
- STEM(x_i) \otimes EXP-INTENSITY(x_i)
- COUNT-OF-STRONG, COUNT-OF-WEAK:
the number of strong and weak EXP-INTENSITY words in the current sentence.
- INTENSIFIER(x_i): whether x_i is an intensifier, such as “extremely”, “highly”, “really”.
- STRONGMODAL(x_i): whether x_i is a strong modal verb, such as “must”, “can”, “will”.
- WEAKMODAL(x_i): whether x_i is a weak modal verb, such as “may”, “could”, “would”.
- DIMINISHER(x_i): whether x_i is a diminisher, such as “little”, “somewhat”, “less”.
- PRECEDED-BY- τ (x_i),
PRECEDED-BY- τ (x_i) \otimes EXP-INTENSITY(x_i):
where $\tau \in$ { INTENSIFIER, STRONGMODAL, WEAKMODAL, DIMINISHER }
- τ (x_i) \otimes EXP-INTENSITY(x_i),
 τ (x_i) \otimes EXP-INTENSITY(x_{i-1}),
 τ (x_{i-1}) \otimes EXP-INTENSITY(x_{i+1})
- EXP-SPAN(x_i) \otimes EXP-INTENSITY(x_i)
- DISTANCE-TO-EXP-SPAN(x_i) \otimes EXP-INTENSITY(x_p)

3.2 Transition Features

Transition features are employed to help with boundary extraction as follows:

Polarity Transition Features

Polarity transition features are features that are used only for $g'_o(\alpha, \hat{\alpha}, x, i)$ and $g'_p(\beta, \hat{\beta}, x, i)$.

- PART-OF-SPEECH(x_i) \otimes PART-OF-SPEECH(x_{i+1}) \otimes EXP-POLARITY(x_i)
- EXP-POLARITY(x_i) \otimes EXP-POLARITY(x_{i+1})

Intensity Transition Features

Intensity transition features are features that are used only for $g'_o(\alpha, \hat{\alpha}, x, i)$ and $g'_s(\gamma, \hat{\gamma}, x, i)$.

- PART-OF-SPEECH(x_i) \otimes PART-OF-SPEECH(x_{i+1}) \otimes EXP-INTENSITY(x_i)
- EXP-INTENSITY(x_i) \otimes EXP-INTENSITY(x_{i+1})

4 Evaluation

We evaluate our system using the Multi-Perspective Question Answering (MPQA) corpus¹. Our gold standard opinion expressions cor-

¹The MPQA corpus can be obtained at <http://nrrc.mitre.org/NRRC/publications.htm>.

Method Description	Positive			Neutral			Negative		
	r(%)	p(%)	f(%)	r(%)	p(%)	f(%)	r(%)	p(%)	f(%)
Polarity-Only \cap Intensity-Only (BASELINE1)	29.6	65.7	40.8	26.5	69.1	38.3	35.5	77.0	48.6
Joint without Hierarchy (BASELINE2)	30.7	65.7	41.9	29.9	66.5	41.2	37.3	77.1	50.3
Joint with Hierarchy	31.8	67.1	43.1	31.9	66.6	43.1	40.4	76.2	52.8

Table 2: Performance of Opinion Extraction with Correct Polarity Attribute

Method Description	High			Medium			Low		
	r(%)	p(%)	f(%)	r(%)	p(%)	f(%)	r(%)	p(%)	f(%)
Polarity-Only \cap Intensity-Only (BASELINE1)	26.4	58.3	36.3	29.7	59.0	39.6	15.4	60.3	24.5
Joint without Hierarchy (BASELINE2)	29.7	54.2	38.4	28.0	57.4	37.6	18.8	55.0	28.0
Joint with Hierarchy	27.1	55.2	36.3	32.0	56.5	40.9	21.1	56.3	30.7

Table 3: Performance of Opinion Extraction with Correct Intensity Attribute

Method Description	r(%)	p(%)	f(%)
Polar-Only \cap Intensity-Only	43.3	92.0	58.9
Joint without Hierarchy	46.0	88.4	60.5
Joint with Hierarchy	48.0	87.8	62.0

Table 4: Performance of Opinion Extraction

respond to *direct subjective expression* and *expressive subjective element* (Wiebe et al., 2005).²

Our implementation of hierarchical sequential learning is based on the Mallet (McCallum, 2002) code for CRFs. In all experiments, we use a Gaussian prior of 1.0 for regularization. We use 135 documents for development, and test on a different set of 400 documents using 10-fold cross-validation. We investigate three options for jointly extracting opinion expressions with their attributes as follows:

[Baseline-1] Polarity-Only \cap Intensity-Only:

For this baseline, we train two separate sequence tagging CRFs: one that extracts opinion expressions only with the polarity attribute (using common features and polarity extraction features in Section 3), and another that extracts opinion expressions only with the intensity attribute (using common features and intensity extraction features in Section 3). We then combine the results from two separate CRFs by collecting all opinion entities extracted by both sequence taggers.³ This

²Only 1.5% of the polarity annotations correspond to *both*; hence, we merge *both* into the *neutral*. Similarly, for gold standard intensity, we merge *extremely high* into *high*.

³We collect all entities whose portions of text spans are extracted by both models.

baseline effectively represents a cascaded component approach.

[Baseline-2] Joint without Hierarchy: Here we use simple linear-chain CRFs without exploiting the class hierarchy for the opinion recognition task. We use the tags shown in Table 1.

Joint with Hierarchy: Finally, we test the hierarchical sequential learning approach elaborated in Section 3.

4.1 Evaluation Results

We evaluate all experiments at the opinion entity level, i.e. at the level of each opinion expression rather than at the token level. We use three evaluation metrics: recall, precision, and F-measure with equally weighted recall and precision.

Table 4 shows the performance of opinion extraction without matching any attribute. That is, an extracted opinion entity is counted as correct if it overlaps⁴ with a gold standard opinion expression, without checking the correctness of its attributes. Table 2 and 3 show the performance of opinion extraction with the correct polarity and intensity respectively.

From all of these evaluation criteria, JOINT WITH

⁴Overlap matching is a reasonable choice as the annotator agreement study is also based on overlap matching (Wiebe et al., 2005). One might wonder whether the overlap matching scheme could allow a degenerative case where extracting the entire test dataset as one giant opinion expression would yield 100% recall and precision. Because each sentence corresponds to a different test instance in our model, and because some sentences do not contain any opinion expression in the dataset, such degenerative case is not possible in our experiments.

HIERARCHY performs the best, and the least effective one is BASELINE-1, which cascades two separately trained models. It is interesting that the simple sequential tagging approach even without exploiting the hierarchy (BASELINE-2) performs better than the cascaded approach (BASELINE-1).

When evaluating with respect to the polarity attribute, the performance of the negative class is substantially higher than that of other classes. This is not surprising as there is approximately twice as much data for the negative class. When evaluating with respect to the intensity attribute, the performance of the LOW class is substantially lower than that of other classes. This result reflects the fact that it is inherently harder to distinguish an opinion expression with low intensity from no opinion. In general, we observe that determining correct intensity attributes is a much harder task than determining correct polarity attributes.

In order to have a sense of upper bound, we also report the individual performance of two separately trained models used for BASELINE-1: for the Polarity-Only model that extracts opinion boundaries only with polarity attribute, the F-scores with respect to the positive, neutral, negative classes are 46.7, 47.5, 57.0, respectively. For the Intensity-Only model, the F-scores with respect to the high, medium, low classes are 37.1, 40.8, 26.6, respectively. Remind that neither of these models alone fully solve the joint task of extracting boundaries as well as determining two attributions simultaneously. As a result, when conjoining the results from the two models (BASELINE-1), the final performance drops substantially.

We conclude from our experiments that the simple joint sequential tagging approach even without exploiting the hierarchy brings a better performance than combining two separately developed systems. In addition, our hierarchical joint sequential learning approach brings a further performance gain over the simple joint sequential tagging method.

5 Related Work

Although there have been much research for fine-grained opinion analysis (e.g., Hu and Liu (2004), Wilson et al. (2005), Wilson et al. (2006), Choi and Claire (2008), Wilson et al. (2009)),⁵ none is

⁵For instance, the results of Wilson et al. (2005) is not comparable even for our Polarity-Only model used inside BASELINE-1, because Wilson et al. (2005) does not operate

directly comparable to our results; much of previous work studies only a subset of what we tackle in this paper. However, as shown in Section 4.1, when we train the learning models only for a subset of the tasks, we can achieve a better performance instantly by making the problem simpler. Our work differs from most of previous work in that we investigate how solving multiple related tasks affects performance on sub-tasks.

The hierarchical parameter sharing technique used in this paper has been previously used by Zhao et al. (2008) for opinion analysis. However, Zhao et al. (2008) employs this technique only to classify sentence-level attributes (polarity and intensity), without involving a much harder task of detecting boundaries of sub-sentential entities.

6 Conclusion

We applied a hierarchical parameter sharing technique using Conditional Random Fields for fine-grained opinion analysis. Our proposed approach jointly extract opinion expressions from unstructured text and determine their attributes — polarity and intensity. Empirical results indicate that the simple joint sequential tagging approach even without exploiting the hierarchy brings a better performance than combining two separately developed systems. In addition, we found that the hierarchical joint sequential learning approach improves the performance over the simple joint sequential tagging method.

Acknowledgments

This work was supported in part by National Science Foundation Grants BCS-0904822, BCS-0624277, IIS-0535099 and by the Department of Homeland Security under ONR Grant N0014-07-1-0152. We thank the reviewers and Ainur Yesenalina for many helpful comments.

References

- S. Abney. 1996. Partial parsing via finite-state cascades. In *Journal of Natural Language Engineering*, 2(4).
- E. Breck, Y. Choi and C. Cardie. 2007. Identifying Expressions of Opinion in Context. In *IJCAI*.

on the entire corpus as unstructured input. Instead, Wilson et al. (2005) evaluate only on known words that are in their opinion lexicon. Furthermore, Wilson et al. (2005) simplifies the problem by combining neutral opinions and no opinions into the same class, while our system distinguishes the two.

- L. Cai and T. Hofmann. 2004. Hierarchical document categorization with support vector machines. In *CIKM*.
- Y. Choi and C. Cardie. 2008. Learning with Compositional Semantics as Structural Inference for Subsentential Sentiment Analysis. In *EMNLP*.
- H. Cunningham, D. Maynard, K. Bontcheva and V. Tablan. 2002. GATE: A Framework and Graphical Development Environment for Robust NLP Tools and Applications. In *ACL*.
- J. R. Finkel, C. D. Manning and A. Y. Ng. 2006. Solving the Problem of Cascading Errors: Approximate Bayesian Inference for Linguistic Annotation Pipelines. In *EMNLP*.
- M. Hu and B. Liu. 2004. Mining and Summarizing Customer Reviews. In *KDD*.
- S. Kim and E. Hovy. 2004. Determining the sentiment of opinions. In *COLING*.
- S. Kim and E. Hovy. 2005. Automatic Detection of Opinion Bearing Words and Sentences. In Companion Volume to the *Proceedings of the Second International Joint Conference on Natural Language Processing (IJCNLP-05)*.
- J. Lafferty, A. McCallum and F. Pereira. 2001. Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data. In *ICML*.
- A. McCallum. 2002. MALLET: A Machine Learning for Language Toolkit. <http://mallet.cs.umass.edu>.
- G. A. Miller. 1995. WordNet: a lexical database for English. In *Communications of the ACM*, 38(11).
- Ana-Maria Popescu and O. Etzioni. 2005. Extracting Product Features and Opinions from Reviews. In *HLT-EMNLP*.
- J. Wiebe, E. Breck, C. Buckley, C. Cardie, P. Davis, B. Fraser, D. Litman, D. Pierce, E. Riloff and T. Wilson. 2002. Summer Workshop on Multiple-Perspective Question Answering: Final Report. In *NRRC*.
- J. Wiebe and T. Wilson and C. Cardie 2005. Annotating Expressions of Opinions and Emotions in Language. In *Language Resources and Evaluation, volume 39, issue 2-3*.
- T. Wilson, J. Wiebe and P. Hoffmann. 2005. Recognizing Contextual Polarity in Phrase-Level Sentiment Analysis. In *HLT-EMNLP*.
- T. Wilson, J. Wiebe and R. Hwa. 2006. Recognizing strong and weak opinion clauses. In *Computational Intelligence*. 22 (2): 73-99.
- T. Wilson, J. Wiebe and P. Hoffmann. 2009. Recognizing Contextual Polarity: an exploration of features for phrase-level sentiment analysis. *Computational Linguistics* 35(3).
- J. Zhao, K. Liu and G. Wang. 2008. Adding Redundant Features for CRFs-based Sentence Sentiment Classification. In *EMNLP*.

Jointly optimizing a two-step conditional random field model for machine transliteration and its fast decoding algorithm

Dong Yang, Paul Dixon and Sadaoki Furui

Department of Computer Science

Tokyo Institute of Technology

Tokyo 152-8552 Japan

{raymond,dixonp,furui}@furui.cs.titech.ac.jp

Abstract

This paper presents a joint optimization method of a two-step conditional random field (CRF) model for machine transliteration and a fast decoding algorithm for the proposed method. Our method lies in the category of direct orthographical mapping (DOM) between two languages without using any intermediate phonemic mapping. In the two-step CRF model, the first CRF segments an input word into chunks and the second one converts each chunk into one unit in the target language. In this paper, we propose a method to jointly optimize the two-step CRFs and also a fast algorithm to realize it. Our experiments show that the proposed method outperforms the well-known joint source channel model (JSCM) and our proposed fast algorithm decreases the decoding time significantly. Furthermore, combination of the proposed method and the JSCM gives further improvement, which outperforms state-of-the-art results in terms of top-1 accuracy.

1 Introduction

There are more than 6000 languages in the world and 10 languages of them have more than 100 million native speakers. With the information revolution and globalization, systems that support multiple language processing and spoken language translation become urgent demands. The translation of named entities from alphabetic to syllabary language is usually performed through transliteration, which tries to preserve the pronunciation in the original language.

For example, in Chinese, foreign words are written with Chinese characters; in Japanese, foreign words are usually written with special char-

Source Name	Target Name	Note
Google	谷歌 gu ge	English-to-Chinese Chinese Romanized writing
Google	グーグル guu gu ru	English-to-Japanese Japanese Romanized writing

Figure 1: Transliteration examples

acters called Katakana; examples are given in Figure 1.

An intuitive transliteration method (Knight and Graehl, 1998; Oh et al., 2006) is to firstly convert a source word into phonemes, then find the corresponding phonemes in the target language, and finally convert them to the target language's written system. There are two reasons why this method does not work well: first, the named entities have diverse origins and this makes the grapheme-to-phoneme conversion very difficult; second, the transliteration is usually not only determined by the pronunciation, but also affected by how they are written in the original language.

Direct orthographical mapping (DOM), which performs the transliteration between two languages directly without using any intermediate phonemic mapping, is recently gaining more attention in the transliteration research community, and it is also the "Standard Run" of the "NEWS 2009 Machine Transliteration Shared Task" (Li et al., 2009). In this paper, we try to make our system satisfy the standard evaluation condition, which requires that the system uses the provided parallel corpus (without pronunciation) only, and cannot use any other bilingual or monolingual resources.

The source channel and joint source channel models (JSCMs) (Li et al., 2004) have been proposed for DOM, which try to model $P(T|S)$ and $P(T, S)$ respectively, where T and S denote the words in the target and source languages. Ekbal et al. (2006) modified the JSCM to incorporate different context information into the model for

Indian languages. In the “NEWS 2009 Machine Transliteration Shared Task”, a new two-step CRF model for transliteration task has been proposed (Yang et al., 2009), in which the first step is to segment a word in the source language into character chunks and the second step is to perform a context-dependent mapping from each chunk into one written unit in the target language.

In this paper, we propose to jointly optimize a two-step CRF model. We also propose a fast decoding algorithm to speed up the joint search. The rest of this paper is organized as follows: Section 2 explains the two-step CRF method, followed by Section 3 which describes our joint optimization method and its fast decoding algorithm; Section 4 introduces a rapid implementation of a JSCM system in the weighted finite state transducer (WFST) framework; and the last section reports the experimental results and conclusions. Although our method is language independent, we use an English-to-Chinese transliteration task in all the explanations and experiments.

2 Two-step CRF method

2.1 CRF introduction

A chain-CRF (Lafferty et al., 2001) is an undirected graphical model which assigns a probability to a label sequence $L = l_1 l_2 \dots l_T$, given an input sequence $C = c_1 c_2 \dots c_T$. CRF training is usually performed through the L-BFGS algorithm (Walach, 2002) and decoding is performed by the Viterbi algorithm. We formalize machine transliteration as a CRF tagging problem, as shown in Figure 2.

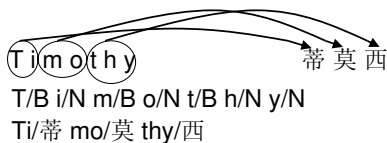


Figure 2: An pictorial description of a CRF segmenter and a CRF converter

2.2 CRF segmenter

In the CRF, a feature function describes a co-occurrence relation, and it is usually a binary function, taking the value 1 when both an observation and a label transition are observed. Yang et al. (2009) used the following features in the segmentation tool:

- Single unit features: $C_{-2}, C_{-1}, C_0, C_1, C_2$
- Combination features: $C_{-1}C_0, C_0C_1$

Here, C_0 is the current character, C_{-1} and C_1 denote the previous and next characters, and C_{-2} and C_2 are the characters located two positions to the left and right of C_0 .

One limitation of their work is that only top-1 segmentation is output to the following CRF converter.

2.3 CRF converter

Similar to the CRF segmenter, the CRF converter has the format shown in Figure 2.

For this CRF, Yang et al. (2009) used the following features:

- Single unit features: CK_{-1}, CK_0, CK_1
- Combination features: $CK_{-1}CK_0, CK_0CK_1$

where CK represents the source language chunk, and the subscript notation is the same as the CRF segmenter.

3 Joint optimization and its fast decoding algorithm

3.1 Joint optimization

We denote a word in the source language by S , a segmentation of S by A , and a word in the target language by T . Our goal is to find the best word \hat{T} in the target language which maximizes the probability $P(T|S)$.

Yang et al. (2009) used only the best segmentation in the first CRF and the best output in the second CRF, which is equivalent to

$$\begin{aligned} \hat{A} &= \arg \max_A P(A|S) \\ \hat{T} &= \arg \max_T P(T|S, \hat{A}), \end{aligned} \quad (1)$$

where $P(A|S)$ and $P(T|S, A)$ represent two CRFs respectively. This method considers the segmentation and the conversion as two independent steps. A major limitation is that, if the segmentation from the first step is wrong, the error propagates to the second step, and the error is very difficult to recover.

In this paper, we propose a new method to jointly optimize the two-step CRF, which can be

written as:

$$\begin{aligned}
\hat{T} &= \arg \max_T P(T|S) \\
&= \arg \max_T \sum_A P(T, A|S) \\
&= \arg \max_T \sum_A P(A|S)P(T|S, A)
\end{aligned} \tag{2}$$

The joint optimization considers all the segmentation possibilities and sums the probability over all the alternative segmentations which generate the same output. It considers the segmentation and conversion in a unified framework and is robust to segmentation errors.

3.2 N-best approximation

In the process of finding the best output using Equation 2, a dynamic programming algorithm for joint decoding of the segmentation and conversion is possible, but the implementation becomes very complicated. Another direction is to divide the decoding into two steps of segmentation and conversion, which is this paper’s method. However, exact inference by listing all possible candidates explicitly and summing over all possible segmentations is intractable, because of the exponential computation complexity with the source word’s increasing length.

In the segmentation step, the number of possible segmentations is 2^N , where N is the length of the source word and 2 is the size of the tagging set. In the conversion step, the number of possible candidates is $M^{N'}$, where N' is the number of chunks from the 1st step and M is the size of the tagging set. M is usually large, e.g., about 400 in Chinese and 50 in Japanese, and it is impossible to list all the candidates.

Our analysis shows that beyond the 10th candidate, almost all the probabilities of the candidates in both steps drop below 0.01. Therefore we decided to generate top-10 results for both steps to approximate the Equation 2.

3.3 Fast decoding algorithm

As introduced in the previous subsection, in the whole decoding process we have to perform n-best CRF decoding in the segmentation step and 10 n-best CRF decoding in the second CRF. Is it really necessary to perform the second CRF for all the segmentations? The answer is “No” for candidates

with low probabilities. Here we propose a no-loss fast decoding algorithm for deciding when to stop performing the second CRF decoding.

Suppose we have a list of segmentation candidates which are generated by the 1st CRF, ranked by probabilities $P(A|S)$ in descending order $A : A_1, A_2, \dots, A_N$ and we are performing the 2nd CRF decoding starting from A_1 . Up to A_k , we get a list of candidates $T : T_1, T_2, \dots, T_L$, ranked by probabilities in descending order. If we can guarantee that, even performing the 2nd CRF decoding for all the remaining segmentations $A_{k+1}, A_{k+2}, \dots, A_N$, the top 1 candidate does not change, then we can stop decoding.

We can show that the following formula is the stop condition:

$$P_k(T_1|S) - P_k(T_2|S) > 1 - \sum_{j=1}^k P(A_j|S). \tag{3}$$

The meaning of this formula is that the probability of all the remaining candidates is smaller than the probability difference between the best and the second best candidates; on the other hand, even if all the remaining probabilities are added to the second best candidate, it still cannot overturn the top candidate. The mathematical proof is provided in Appendix A.

The stop condition here has no approximation nor pre-defined assumption, and it is a no-loss fast decoding algorithm.

4 Rapid development of a JSCM system

The JSCM represents how the source words and target names are generated simultaneously (Li et al., 2004):

$$\begin{aligned}
P(S, T) &= P(s_1, s_2, \dots, s_k, t_1, t_2, \dots, t_k) \\
&= P(\langle s, t \rangle_1, \langle s, t \rangle_2, \dots, \langle s, t \rangle_k) \\
&= \prod_{k=1}^K P(\langle s, t \rangle_k | \langle s, t \rangle_1^{k-1})
\end{aligned} \tag{4}$$

where $S = (s_1, s_2, \dots, s_k)$ is a word in the source language and $T = (t_1, t_2, \dots, t_k)$ is a word in the target language.

The training parallel data without alignment is first aligned by a Viterbi version EM algorithm (Li et al., 2004).

The decoding problem in JSCM can be written as:

$$\hat{T} = \arg \max_T P(S, T). \tag{5}$$

After the alignments are generated, we use the MITLM toolkit (Hsu and Glass, 2008) to build a trigram model with modified Kneser-Ney smoothing. We then convert the n-gram to a WFST M (Sproat et al., 2000; Caseiro et al., 2002). To allow transliteration from a sequence of characters, a second WFST T is constructed. The input word is converted to an acceptor I , and it is then combined with T and M according to $O = I \circ T \circ M$ where \circ denotes the composition operator. The n-best paths are extracted by projecting the output, removing the epsilon labels and applying the n-shortest paths algorithm with determinization in the OpenFst Toolkit (Allauzen et al., 2007).

5 Experiments

We use several metrics from (Li et al., 2009) to measure the performance of our system.

1. Top-1 ACC: word accuracy of the top-1 candidate
2. Mean F-score: fuzziness in the top-1 candidate, how close the top-1 candidate is to the reference
3. MRR: mean reciprocal rank, $1/\text{MRR}$ tells approximately the average rank of the correct result

5.1 Comparison with the baseline and JSCM

We use the training, development and test sets of NEWS 2009 data for English-to-Chinese in our experiments as detailed in Table 1. This is a parallel corpus without alignment.

Training data	Development data	Test data
31961	2896	2896

Table 1: Corpus size (number of word pairs)

We compare the proposed decoding method with the baseline which uses only the best candidates in both CRF steps, and also with the well known JSCM. As we can see in Table 2, the proposed method improves the baseline top-1 ACC from 0.670 to 0.708, and it works as well as, or even better than the well known JSCM in all the three measurements.

Our experiments also show that the decoding time can be reduced significantly via using our fast decoding algorithm. As we have explained, without fast decoding, we need 11 CRF n-best decoding for each word; the number can be reduced to 3.53 (1 “the first CRF”+2.53 “the second CRF”) via the fast decoding algorithm.

We should notice that the decoding time is significantly shorter than the training time. While

testing takes minutes on a normal PC, the training of the CRF converter takes up to 13 hours on an 8-core (8*3G Hz) server.

Measure	Top-1 ACC	Mean F-score	MRR
Baseline	0.670	0.869	0.750
Joint optimization	0.708	0.885	0.789
JSCM	0.706	0.882	0.789

Table 2: Comparison of the proposed decoding method with the previous method and the JSCM

5.2 Further improvement

We tried to combine the two-step CRF model and the JSCM. From the two-step CRF model we get the conditional probability $P_{CRF}(T|S)$ and from the JSCM we get the joint probability $P(S,T)$. The conditional probability of $P_{JSCM}(T|S)$ can be calculated as follows:

$$P_{JSCM}(T|S) = \frac{P(T,S)}{P(S)} = \frac{P(T,S)}{\sum_T P(T,S)}. \quad (6)$$

They are used in our combination method as:

$$P(T|S) = \lambda P_{CRF}(T|S) + (1 - \lambda) P_{JSCM}(T|S) \quad (7)$$

where λ denotes the interpolation weight (λ is set by development data in this paper).

As we can see in Table 3, the linear combination of two systems further improves the top-1 ACC to 0.720, and it has outperformed the best reported “Standard Run” (Li et al., 2009) result 0.717. (The reported best “Standard Run” result 0.731 used target language phoneme information, which requires a monolingual dictionary; as a result it is not a standard run.)

Measure	Top-1 ACC	Mean F-score	MRR
Baseline+JSCM	0.713	0.883	0.794
Joint optimization + JSCM	0.720	0.888	0.797
state-of-the-art (Li et al., 2009)	0.717	0.890	0.785

Table 3: Model combination results

6 Conclusions and future work

In this paper we have presented our new joint optimization method for a two-step CRF model and its fast decoding algorithm. The proposed

method improved the system significantly and outperformed the JSCM. Combining the proposed method with JSCM, the performance was further improved.

In future work we are planning to combine our system with multilingual systems. Also we want to make use of acoustic information in machine transliteration. We are currently investigating discriminative training as a method to further improve the JSCM. Another issue of our two-step CRF method is that the training complexity increases quadratically according to the size of the label set, and how to reduce the training time needs more research.

Appendix A. Proof of Equation 3

The CRF segmentation provides a list of segmentations: $A : A_1, A_2, \dots, A_N$, with conditional probabilities $P(A_1|S), P(A_2|S), \dots, P(A_N|S)$.

$$\sum_{j=1}^N P(A_j|S) = 1.$$

The CRF conversion, given a segmentation A_i , provides a list of transliteration output T_1, T_2, \dots, T_M , with conditional probabilities $P(T_1|S, A_i), P(T_2|S, A_i), \dots, P(T_M|S, A_i)$.

In our fast decoding algorithm, we start performing the CRF conversion from A_1 , then A_2 , and then A_3 , etc. Up to A_k , we get a list of candidates $T : T_1, T_2, \dots, T_L$, ranked by probabilities $P_k(T|S)$ in descending order. The probability $P_k(T_l|S) (l = 1, 2, \dots, L)$ is accumulated probability of $P(T_l|S)$ over A_1, A_2, \dots, A_k , calculated by:

$$P_k(T_l|S) = \sum_{j=1}^k P(A_j|S)P(T_l|S, A_j)$$

If we continue performing the CRF conversion to cover all $N (N \geq k)$ segmentations, eventually we will get:

$$\begin{aligned} P(T_l|S) &= \sum_{j=1}^N P(A_j|S)P(T_l|S, A_j) \\ &\geq \sum_{j=1}^k P(A_j|S)P(T_l|S, A_j) \\ &= P_k(T_l|S) \end{aligned} \quad (8)$$

If Equation 3 holds, then for $\forall i \neq 1$,

$$\begin{aligned} P_k(T_1|S) &> P_k(T_2|S) + (1 - \sum_{j=1}^k P(A_j|S)) \\ &\geq P_k(T_i|S) + (1 - \sum_{j=1}^k P(A_j|S)) \\ &= P_k(T_i|S) + \sum_{j=k+1}^N P(A_j|S) \\ &\geq P_k(T_i|S) \\ &\quad + \sum_{j=k+1}^N P(A_j|S)P(T_i|S, A_j) \\ &= P(T_i|S) \end{aligned} \quad (9)$$

Therefore, $P(T_1|S) > P(T_i|S) (i \neq 1)$, and T_1 maximizes the probability $P(T|S)$.

References

- Cyril Allauzen, Michael Riley, Johan Schalkwyk, Wojciech Skut and Mehryar Mohri 2007. *OpenFst: A General and Efficient Weighted Finite-State Transducer Library*. Proceedings of the Ninth International Conference on Implementation and Application of Automata, (CIAA), pages 11-23.
- Diamantino Caseiro, Isabel Trancosoo, Luis Oliveira and Ceu Viana 2002. *Grapheme-to-phone using finite state transducers*. Proceedings IEEE Workshop on Speech Synthesis.
- Asif Ekbal, Sudip Kumar Naskar and Sivaji Bandyopadhyay. 2006. *A modified joint source-channel model for transliteration*, Proceedings of the COLING/ACL, pages 191-198.
- Bo-June Hsu and James Glass 2008. *Iterative Language Model Estimation: Efficient Data Structure & Algorithms*. Proceedings Interspeech, pages 841-844.
- Kevin Knight and Jonathan Graehl. 1998. *Machine Transliteration*, Association for Computational Linguistics.
- John Lafferty, Andrew McCallum, and Fernando Pereira 2001. *Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data.*, Proceedings of International Conference on Machine Learning, pages 282-289.
- Haizhou Li, Min Zhang and Jian Su. 2004. *A joint source-channel model for machine transliteration*, Proceedings of the 42nd Annual Meeting on Association for Computational Linguistics.
- Haizhou Li, A. Kumaran, Vladimir Pervouchine and Min Zhang 2009. *Report of NEWS 2009 Machine Transliteration Shared Task*, Proceedings of the 2009 Named Entities Workshop: Shared Task on Transliteration (NEWS 2009), pages 1-18
- Jong-Hoon Oh, Key-Sun Choi and Hitoshi Isahara. 2006. *A comparison of different machine transliteration models*, Journal of Artificial Intelligence Research, 27, pages 119-151.
- Richard Sproat 2000. *Corpus-Based Methods and Hand-Built Methods*. Proceedings of International Conference on Spoken Language Processing, pages 426-428.
- Andrew J. Viterbi 1967. *Error Bounds for Convolutional Codes and an Asymptotically Optimum Decoding Algorithm*. IEEE Transactions on Information Theory, Volume IT-13, pages 260-269.
- Hanna Wallach 2002. *Efficient Training of Conditional Random Fields*. M. Thesis, University of Edinburgh.
- Dong Yang, Paul Dixon, Yi-Cheng Pan, Tasuku Oonishi, Masanobu Nakamura and Sadaoki Furui 2009. *Combining a Two-step Conditional Random Field Model and a Joint Source Channel Model for Machine Transliteration*, Proceedings of the 2009 Named Entities Workshop: Shared Task on Transliteration (NEWS 2009), pages 72-75

Arabic Named Entity Recognition: Using Features Extracted from Noisy Data

Yassine Benajiba¹ Imed Zitouni² Mona Diab¹ Paolo Rosso³

¹ Center for Computational Learning Systems, Columbia University

² IBM T.J. Watson Research Center, Yorktown Heights

³ Natural Language Engineering Lab. - ELiRF, Universidad Politécnic de Valencia

{ybenajiba, mdiab}@ccls.columbia.edu, izitouni@us.ibm.com, proso@dsic.upv.es

Abstract

Building an accurate Named Entity Recognition (NER) system for languages with complex morphology is a challenging task. In this paper, we present research that explores the feature space using both gold and bootstrapped noisy features to build an improved highly accurate Arabic NER system. We bootstrap noisy features by projection from an Arabic-English parallel corpus that is automatically tagged with a baseline NER system. The feature space covers lexical, morphological, and syntactic features. The proposed approach yields an improvement of up to 1.64 F-measure (absolute).

1 Introduction

Named Entity Recognition (NER) has earned an important place in Natural Language Processing (NLP) as an enabling process for other tasks. When explicitly taken into account, research shows that it helps such applications achieve better performance levels (Babych and Hartley, 2003; Thompson and Dozier, 1997). NER is defined as the computational identification and classification of Named Entities (NEs) in running text. For instance, consider the following text:

Barack Obama is visiting the Middle East.

A NER system should be able to identify *Barack Obama* and *Middle East* as NEs and classify them as *Person* (PER) and *Geo-Political Entity* (GPE), respectively. The class-set used to tag NEs may vary according to user needs. In this research, we adopt the Automatic Content Extraction (ACE) 2007 nomenclature¹.

According to (Nadeau and Sekine, 2007), optimization of the feature set is the key component in enhancing the performance of a global NER system. In this paper we investigate the possibility of building a high performance Arabic NER system by using a large space of available feature sets that go beyond the explored shallow feature sets used to date in the literature for Arabic NER.

¹<http://www.nist.gov/speech/tests/ace/index.htm>

Given current state-of-the-art syntactic processing of Arabic text and the relative small size of manually annotated Arabic NER data, we set out to explore a main concrete research goal: to fully exploit the level of advancement in Arabic lexical and syntactic processing to explore deeper linguistic features for the NER task. Realizing that the gold data available for NER is quite limited in size especially given the diverse genres in the set, we devise a method to bootstrap additional instances for the new features of interest from noisily NER tagged Arabic data.

2 Our Approach

We use our state-of-the-art NER system described in (Benajiba et al., 2008) as our baseline system (BASE) since it yields, to our knowledge, the best performance for Arabic NER. BASE employs Support Vector Machines (SVMs) and Conditional Random Fields (CRFs) as Machine Learning (ML) approaches. BASE uses lexical, syntactic and morphological features extracted using highly accurate automatic Arabic POS-taggers. BASE employs a multi-classifier approach where each classifier is tagging a NE class separately. The feature selection is performed by using an incremental approach selecting the top n features (the features are ranked according to their individual impact) at each iteration and keeping the set that yields the best results. In case of conflict - a word is classified with more than one class/tag simultaneously - the global NER system selects the output of the classifier with the highest precision.

The following is the feature set used in (Benajiba et al., 2008) and accordingly in the BASE system. **1. Context:** a $-/+1$ token window; **2. Lexical:** character $n - grams$ where n ranges from $1 - 3$; **3. Gazetteers:** automatically harvested and manually cleaned Person NE class (PER), Geopolitical Entity NE class (GPE), and Organization NE class (ORG) lexica; **4. POS-tag and Base Phrase Chunk (BPC):** automatically tagged using AMIRA (Diab et al., 2007) which yields F-measures for both tasks in the high 90's; **5. Morphological features:** automatically tagged using the Morphological Analysis and Disambiguation for Arabic (MADA) tool to extract information about gender, number, person, definiteness and as-

pect for each word (Habash and Rambow, 2005);

6. Capitalization: derived as a side effect from running MADA. MADA chooses a specific morphological analysis given the context of a given word. As part of the morphological information available in the underlying lexicon that MADA exploits. As part of the information present, the underlying lexicon has an English gloss associated with each entry. More often than not, if the word is a NE in Arabic then the gloss will also be a NE in English and hence capitalized.

We devise an extended Arabic NER system (EXTENDED) that uses the same architecture as BASE but employs additional features to those in BASE. EXTENDED defines new additional syntagmatic features.

We specifically investigate the space of the surrounding context for the NEs. We explore generalizations over the kinds of words that occur with NEs and the syntactic relations NEs engage in. We use an off-the-shelf Arabic syntactic parser. State-of-the-art for Arabic syntactic parsing for the most common genre (with the most training data) of Arabic data, newswire, is in the low 80%s. Hence, we acknowledge that some of the derived syntactic features will be noisy.

Similar to all supervised ML problems, it is desirable to have sufficient training data for the relevant phenomena. The size of the manually annotated gold data typically used for training Arabic NER systems poses a significant challenge for robustly exploring deeper syntactic and lexical features. Accordingly, we bootstrap more NE tagged data via projection over Arabic-English parallel data. The role of this data is simply to give us more instances of the newly defined features (namely the syntagmatic features) in the EXTENDED system as well as more instances for the Gazetteers and Context features defined in BASE. It is worth noting that we do not use the bootstrapped NE tagged data directly as training data with the gold data.

2.1 Syntagmatic Features

For deriving our deeper linguistic features, we parse the Arabic sentences that contain an NE. For each of the NEs, we extract a number of features described as follows:

- Syntactic head-word (SHW): The idea here is to look for a broader **relevant** context. Whereas the feature lexical n-gram context feature used in BASE, and hence here for EXTENDED, considers the linearly adjacent neighboring words of a NE, SHW uses a parse tree to look at farther, yet related, words. For instance, in the Arabic phrase “SrH Ams An

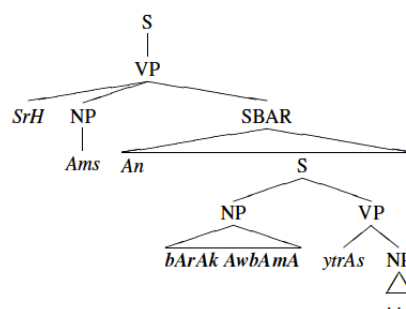


Figure 1: Example for the head word and syntactic environment feature

bArAk AwbAmA ytrAs”, which means “declared yesterday that Barack Obama governs ...”, glossed “SrH/declared Ams/yesterday An/that bArAk/Barack AwbAmA/Obama ytrAs/governs ...”, is parsed in Figure 1. According to the phrase structure parse, the first parent sub-tree headword of the NE “bArAk AwbAmA” is the verb ‘ytrAs’ (governs), the second one is ‘An’ (that) and the third one is the verb ‘SrH’ (declared). This example illustrates that the word “Ams” is ignored for this feature set since it is not a syntactic head. This is a lexicalized feature.

- Syntactic Environment (SE): This follows in the same *spirit* as SHW, but expands the idea in that it looks at the parent non-terminal instead of the parent head word, hence it is not a lexicalized feature. The goal being to use a more abstract representation level of the context in which a NE appears. For instance, for the same example presented in Figure 1, the first, second, and third non-terminal parents of the NE “bArAk AwbAmA” are ‘S’, ‘SBAR’ and ‘VP’, respectively.

In our experiments we use the Bikel implementation (Bikel, 2004) of the Collins parser (Collins, 1999) which is freely available on the web². It is a head-driven CFG-style parser trained to parse English, Arabic, and Chinese.

2.2 Bootstrapping Noisy Arabic NER Data

Extracting the syntagmatic features from the training data yields relatively small number of instances. Hence the need for additional tagged data. The new Arabic NER tagged data is derived via projection exploiting parallel Arabic English data. The process depends on the availability of two key components: a large Arabic English parallel corpus that is sentence and word aligned, and a robust high performing English NER system. The process is as follows. We NE tag the

²<http://www.cis.upenn.edu/~dbikel/software.html#stat-parser>

English side of the parallel corpus. We project the automatically tagged NER tags from the English side to the Arabic side of the parallel corpus. In our case, we have access to a large manually aligned parallel corpus, therefore the NER projection is direct. However, the English side of the parallel corpus is not NER tagged, hence we use an off-the-shelf competitive robust automatic English NER system which has a published performance of 92% (Zitouni and Florian, 2009). The result of these two processes is a large Arabic NER, albeit noisy, tagged data set. As mentioned earlier this data is used only for deriving additional instances for training for the syntagmatic features and for the context and gazetteer features.³ Given this additional source of data, we changed the lexical features extracted from the BASE to the EXTENDED. We added two other lexical features: CBG and NGC, described as follows: - Class Based Gazetteers (CBG): This feature focuses on the surface form of the NEs. We group the NEs encountered on the Arabic side of the parallel corpus by class as they are found in different dictionaries. The difference between this feature and that in BASE is that the Gazetteers are not restricted to Wikipedia sources.

- N-gram context (NGC): Here we disregard the surface form of the NE, instead we focus on its lexical context. For each n , where n varies from 1 to 3, we compile a list of the $-n$, $+n$, and $-/+n$ words surrounding the NE. Similar to the CBG feature, these lists are also separated by NE class. It is worth highlighting that the NCG feature is different from the Context feature in BASE in that the window size is different $+/-1-3$ for EXTENDED versus $+/-1$ for BASE.

3 Experiments and Results

3.1 Gold Data for training and evaluation

We use the standard sets of ACE 2003, ACE 2004 and ACE 2005.⁴ The ACE data is annotated for many tasks: Entity Detection and Tracking (EDT), Relation Detection and Recognition (RDR), Event Detection and Recognition (EDR). All the data sets comprise *Broadcast News* (BN) and *NewsWire* (NW) genres. ACE 2004 includes an additional NW data set from the Arabic TreeBank (ATB). ACE 2005 includes a different genre of *Weblogs* (WL). The NE classes adopted in the annotation of the ACE 2003 data are: Person (PER), Geo Political Entity (GPE), Organization (ORG) and Facility (FAC).

³Therefore, we did not do the full feature extraction for the other features described in BASE for this data.

⁴<http://www.nist.gov/speech/tests/ace/>

Additionally for the ACE 2004 and 2005 data, two NE classes are added to the ACE 2003 tag-set: Vehicles (e.g. Rotterdam Ship) and Weapons (e.g. Kalashnikof). We use the same split for train, development, and test used in (Benajiba et al., 2008).

3.2 Parallel Data

Most of the hand-aligned Arabic-English parallel data used in our experiments is from the Language Data Consortium (LDC).⁵ Another set of the parallel data is annotated in-house by professional annotators. The corpus has texts of five different genres, namely: newswire, news groups, broadcast news, broadcast conversation and weblogs corresponding to the data genres in the ACE gold data. The Arabic side of the parallel corpus contains 941,282 tokens. After projecting the NE tags from the English side to the Arabic side of the parallel corpus, we obtain a total of 57,290 Arabic NE instances. Table 1 shows the number of NEs for each class.

Class	Number of NEs	Class	Number of NEs
FAC	998	PER	17,964
LOC	27,651	VEH	85
ORG	10,572	WEA	20

Table 1: Number of NEs per class in the Arabic side of the parallel corpus

3.3 Individual Feature Impact

Across the board, all the features yield improved performance. The highest obtained result is observed where the first non-terminal parent is used as a feature, a Syntactic Environment (SE) feature, yielding an improvement of up to 4 points over the baseline. We experiment with different sizes for the SE, i.e. taking the first parent versus adding neighboring non-terminal parents. We note that even though we observe an overall increase in performance, considering both the {first, second} or the {first, second, and third} non-terminal parents decreases performance by 0.5 and 1.5 F-measure points, respectively, compared to considering the first parent information alone. The head word features, SHW, show a higher positive impact than the lexical context feature, NGC. Finally, the Gazetteer feature, CBG, impact is comparable to the obtained improvement of the lexical context feature.

3.4 Feature Combination Experiments

Table 2 illustrates the final results. It shows for each data set and each genre the F-measure obtained using the best feature set and ML approach. It shows results for both the dev and test data using the optimal number of features selected from

⁵All the LDC data are publicly available

		ACE 2003		ACE 2004			ACE 2005		
		BN	NW	BN	NW	ATB	BN	NW	WL
	<i>FreqBaseline</i>	73.74	67.61	62.17	51.67	62.94	70.18	57.17	27.66
dev	<i>All-Synt.</i>	83.41	79.11	76.90	72.90	74.82	81.42	76.07	54.49
	<i>All</i>	83.93	79.72	78.54	72.80	74.97	81.82	75.92	55.65
test	<i>All-Synt.</i>	83.50	78.90	76.70	72.40	73.50	81.31	75.30	57.30
	<i>All</i>	84.32	79.4	78.12	72.13	74.54	81.73	75.67	58.11

Table 2: Final Results obtained with selected features contrasted against all features combined

the all the features except the syntagmatic ones (*All-Synt.*) contrasted against the system including the semantic features, i.e. All the features, per class *All*. The baseline results, *FreqBaseline*, assigns a test token the most frequent tag observed for it in the gold training data, if a test token is not observed in the training data, it is assigned the most frequent tag which is the O tag.

4 Results Discussion

Individual feature impact results show that the syntagmatic features are helpful for most of the data sets. The highest improvements are obtained for the 2003 BN and 2005 WL data-sets. The improvement varies significantly from one data-set to another because it highly depends on the number of NEs which the model has not been able to capture using the contextual, lexical, syntactic and morphological features.

Impact of the features extracted from the parallel corpus per class: The syntagmatic features have varied in their influence on the different NE classes. Generally, the LOC and PER classes benefitted more from the head word features, SHW), than the other classes. On the other hand for the syntactic environment feature (SE), the PER class seemed not to benefit much from the presence of this feature. *Weblogs:* Our results show that the random contexts in which the NEs tend to appear in the WL documents stand against obtaining a significant improvement. Consequently, the features which use a more global context (syntactic environment, SE, and head word, SHW, features) have helped obtain better results than the ones which we have obtained using local context namely CBG and NGC.

5 Related Work

Projecting explicit linguistic tags from another language via parallel corpora has been widely used in the NLP tasks and has proved to contribute significantly to achieving better performance. Different research works report positive results when using this technique to enhance WSD (Diab and Resnik, 2002; Ng et al., 2003). In the latter two

works, they augment training data from parallel data for training supervised systems. In (Diab, 2004), the author uses projections from English into Arabic to bootstrap a sense tagging system for Arabic as well as a seed Arabic WordNet through projection. In (Hwa et al., 2002), the authors report promising results of inducing Chinese dependency trees from English. The obtained model outperformed the baseline. More recently, in (Chen and Ji, 2009), the authors report their comparative study between monolingual and cross-lingual bootstrapping. Finally, in Mention Detection (MD), a task which includes NER and adds the identification and classification of nominal and pronominal mentions, (Zitouni and Florian, 2008) show the impact of using a MT system to enhance the performance of an Arabic MD model. The authors report an improvement of up to 1.6F when the baseline system uses lexical features only. Unlike the work we present here, their approach requires the availability of an accurate MT system which is a more expensive process.

6 Conclusion and Future Directions

In this paper we investigate the possibility of building a high performance Arabic NER system by using lexical, syntactic and morphological features and augmenting the model with deeper lexical features and more syntagmatic features. These extra features are extracted from noisy data obtained via projection from an Arabic-English parallel corpus. Our results show that we achieve a significantly high performance for almost all the data-sets. The greatest impact of the syntagmatic features (1.64 points of F-measure) is obtained for the ACE 2004, BN genre. Also, the WL genre yields an improvement of 1.16 F1 points absolute.

Acknowledgments

This work has been partially funded by DARPA GALE project. The research of the last author was funded by MICINN research project TEXT-ENTERPRISE 2.0 TIN2009-13391-C04-03 (Plan I+D+i).

References

- B. Babych and A. Hartley. 2003. Improving Machine Translation Quality with Automatic Named Entity Recognition. In *Proc. of EACL-EAMT*.
- Y. Benajiba, M. Diab, and P. Rosso. 2008. Arabic named entity recognition using optimized feature sets. In *Proceedings of EMNLP'08*, pages 284–293.
- Daniel M. Bikel. 2004. *On the parameter space of generative lexicalized statistical parsing models*. University of Pennsylvania, Philadelphia, PA, USA. Supervisor-Marcus, Mitchell P.
- Z. Chen and H. Ji. 2009. Can one language bootstrap the other: A case study of event extraction. In *Proceedings of NAACL'09*.
- M. Collins. 1999. *Head-Driven Statistical Models for Natural Language Parsing*. University of Pennsylvania, Philadelphia, PA, USA.
- Mona Diab and Philip Resnik. 2002. An unsupervised method for word sense tagging using parallel corpora. In *Proceedings of 40th Annual Meeting of the Association for Computational Linguistics*, pages 255–262, Philadelphia, Pennsylvania, USA, July. Association for Computational Linguistics.
- M. Diab, K. Hacioglu, and D. Jurafsky, 2007. *Arabic Computational Morphology: Knowledge-based and Empirical Methods*, chapter 9. Springer.
- Mona Diab. 2004. Bootstrapping a wordnet taxonomy for arabic. In *Proceedings of First Arabic Language Technology Conference (NEMLAR), Cairo Egypt*.
- N. Habash and O. Rambow. 2005. Arabic Tokenization, Part-of-Speech Tagging and Morphological Disambiguation in One Fell Swoop. In *Proc. of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL'05)*, pages 573–580, Ann Arbor, Michigan, June. Association for Computational Linguistics.
- R. Hwa, P. Resnik, and A. Weinberg. 2002. Breaking the resource bottleneck for multilingual parsing. In *Proceedings of the Workshop on Linguistic Knowledge Acquisition and Representation: Bootstrapping Annotated Language Data*.
- D. Nadeau and S. Sekine. 2007. A Survey of Named Entity Recognition and Classification. *Linguisticae Investigationes*, 30(7).
- H.-T. Ng, B. Wang, and Y.-S. Chan. 2003. Exploiting parallel texts for word sense disambiguation: An empirical study. In *ACL'03*, pages 455–462, Sapporo, Japan.
- P. Thompson and C. Dozier. 1997. Name Searching and Information Retrieval. In *In Proc. of Second Conference on Empirical Methods in Natural Language Processing*, Providence, Rhode Island.
- I. Zitouni and R. Florian. 2008. Mention detection crossing the language barrier. In *Proceedings of EMNLP'08*, Honolulu, Hawaii, October.
- Imed Zitouni and Radu Florian. 2009. Cross language information propagation for arabic mention detection. *Journal of ACM Transactions on Asian Language Information Processing*, December.

Extracting Sequences from the Web

Anthony Fader, Stephen Soderland, and Oren Etzioni

University of Washington, Seattle
{afader,soderlan,etzioni}@cs.washington.edu

Abstract

Classical Information Extraction (IE) systems fill slots in domain-specific frames. This paper reports on SEQ, a novel open IE system that leverages a *domain-independent* frame to extract ordered sequences such as presidents of the United States or the most common causes of death in the U.S. SEQ leverages regularities about sequences to extract a coherent set of sequences from Web text. SEQ nearly doubles the area under the precision-recall curve compared to an extractor that does not exploit these regularities.

1 Introduction

Classical IE systems fill slots in domain-specific frames such as the time and location slots in seminar announcements (Freitag, 2000) or the terrorist organization slot in news stories (Chieu et al., 2003). In contrast, open IE systems are domain-independent, but extract “flat” sets of assertions that are not organized into frames and slots (Sekine, 2006; Banko et al., 2007). This paper reports on SEQ—an open IE system that leverages a domain-independent frame to extract ordered sequences of objects from Web text. We show that the novel, domain-independent sequence frame in SEQ substantially boosts the precision and recall of the system and yields coherent sequences filtered from low-precision extractions (Table 1).

Sequence extraction is distinct from set expansion (Etzioni et al., 2004; Wang and Cohen, 2007) because sequences are *ordered* and because the extraction process does not require seeds or HTML lists as input.

The domain-independent *sequence frame* consists of a sequence name s (e.g., presidents of the United States), and a set of ordered pairs (x, k) where x is a string naming a member of the sequence with name s , and k is an integer indicating

Most common cause of death in the United States: 1. heart disease, 2. cancer, 3. stroke, 4. COPD, 5. pneumonia, 6. cirrhosis, 7. AIDS, 8. chronic liver disease, 9. sepsis, 10. suicide, 11. septic shock.
Largest tobacco company in the world: 1. Philip Morris, 2. BAT, 3. Japan Tobacco, 4. Imperial Tobacco, 5. Altadis.
Largest rodent in the world: 1. Capybara, 2. Beaver, 3. Patagonian Cavies. 4. Maras.
Sign of the zodiac: 1. Aries, 2. Taurus, 3. Gemini, 4. Cancer, 5. Leo, 6. Virgo, 7. Libra, 8. Scorpio, 9. Sagittarius, 10. Capricorn, 11. Aquarius, 12. Pisces, 13. Ophiuchus.

Table 1: Examples of sequences extracted by SEQ from unstructured Web text.

its position (e.g., (Washington, 1) and (JFK, 35)). The task of *sequence extraction* is to automatically instantiate sequence frames given a corpus of unstructured text.

By definition, sequences have two properties that we can leverage in creating a sequence extractor: *functionality* and *density*. Functionality means position k in a sequence is occupied by a single real-world entity x . Density means that if a value has been observed at position k then there must exist values for all $i < k$, and possibly more after it.

2 The SEQ System

Sequence extraction has two parts: identifying possible extractions (x, k, s) from text, and then classifying those extractions as either correct or incorrect. In the following section, we describe a way to identify candidate extractions from text using a set of lexico-syntactic patterns. We then show that classifying extractions based on sentence-level features and redundancy alone yields low precision, which is improved by leveraging the functionality and density properties of sequences as done in our SEQ system.

Pattern	Example
the ORD	the fifth
the RB ORD	the very first
the JJS	the best
the RB JJS	the very best
the ORD JJS	the third biggest
the RBS JJ	the most popular
the ORD RBS JJ	the second least likely

Table 2: The patterns used by SEQ to detect ordinal phrases are noun phrases that begin with one of the part-of-speech patterns listed above.

2.1 Generating Sequence Extractions

To obtain candidate sequence extractions (x, k, s) from text, the SEQ system finds sentences in its input corpus that contain an ordinal phrase (OP). Table 2 lists the lexico-syntactic patterns SEQ uses to detect ordinal phrases. The value of k is set to the integer corresponding to the ordinal number in the OP.¹

Next, SEQ takes each sentence that contains an ordinal phrase o , and finds candidate items of the form (x, k) for the sequence with name s . SEQ constrains x to be an NP that is disjoint from o , and s to be an NP (which may have post-modifying PPs or clauses) following the ordinal number in o .

For example, given the sentence “With help from his father, JFK was elected as the 35th President of the United States in 1960”, SEQ finds the candidate sequences with names “President”, “President of the United States”, and “President of the United States in 1960”, each of which has candidate extractions (JFK, 35), (his father, 35), and (help, 35). We use heuristics to filter out many of the candidate values (*e.g.*, no value should cross a sentence-like boundary, and x should be at most some distance from the OP).

This process of generating candidate extractions has high coverage, but low precision. The first step in identifying correct extractions is to compute a confidence measure $localConf(x, k, s|sentence)$, which measures how likely (x, k, s) is given the sentence it came from. We do this using domain-independent syntactic features based on POS tags and the pattern-based features “ x {is,are,was,were} the k th s ” and “the k th s {is,are,was,were} x ”. The features are then combined using a Naive Bayes classifier.

In addition to the local, sentence-based features,

¹Sequences often use a superlative for the first item ($k = 1$) such as “the deepest lake in Africa”, “the second deepest lake in Africa” (or “the 2nd deepest ...”), etc.

we define the measure $totalConf$ that takes into account redundancy in an input corpus \mathcal{C} . As Downey *et al.* observed (2005), extractions that occur more frequently in multiple distinct sentences are more likely to be correct.

$$totalConf(x, k, s|\mathcal{C}) = \sum_{sentence \in \mathcal{C}} localConf(x, k, s|sentence) \quad (1)$$

2.2 Challenges

The scores $localConf$ and $totalConf$ are not sufficient to identify valid sequence extractions. They tend to give high scores to extractions where the sequence scope is too general or too specific. In our running example, the sequence name “President” is too general – many countries and organizations have a president. The sequence name “President of the United States in 1960” is too specific – there were not multiple U.S. presidents in 1960.

These errors can be explained as violations of functionality and density. The sequence with name “President” will have many distinct candidate extractions in its positions, which is a violation of functionality. The sequence with name “President of the United States in 1960” will not satisfy density, since it will have extractions for only one position.

In the next section, we present the details of how SEQ incorporates functionality and density into its assessment of a candidate extraction.

Given an extraction (x, k, s) , SEQ must classify it as either correct or incorrect. SEQ breaks this problem down into two parts: (1) determining whether s is a correct sequence name, and (2) determining whether (x, k) is an item in s , assuming s is correct.

A joint probabilistic model of these two decisions would require a significant amount of labeled data. To get around this problem, we represent each (x, k, s) as a vector of features and train two Naive Bayes classifiers: one for classifying s and one for classifying (x, k) . We then rank extractions by taking the product of the two classifiers’ confidence scores.

We now describe the features used in the two classifiers and how the classifiers are trained.

Classifying Sequences To classify a sequence name s , SEQ uses features to measure the functionality and density of s . Functionality means

that a correct sequence with name s has one correct value x at each position k , possibly with additional noise due to extraction errors and synonymous values of x . For a fixed sequence name s and position k , we can weight each of the candidate x values in that position by their normalized total confidence:

$$w(x|k, s, \mathcal{C}) = \frac{\text{totalConf}(x, k, s|\mathcal{C})}{\sum_{x'} \text{totalConf}(x', k, s|\mathcal{C})}$$

For overly general sequences, the distribution of weights for a position will tend to be more flat, since there are many equally-likely candidate x values. To measure this property, we use a function analogous to information entropy:

$$H(k, s|\mathcal{C}) = - \sum_x w(x|k, s, \mathcal{C}) \log_2 w(x|k, s, \mathcal{C})$$

Sequences s that are too general will tend to have high values of $H(k, s|\mathcal{C})$ for many values of k . We found that a good measure of the overall non-functionality of s is the average value of $H(k, s|\mathcal{C})$ for $k = 1, 2, 3, 4$.

For a sequence name s that is too specific, we would expect that there are only a few filled-in positions. We model the density of s with two metrics. The first is $\text{numFilledPos}(s|\mathcal{C})$, the number of distinct values of k such that there is some extraction (x, k) for s in the corpus. The second is $\text{totalSeqConf}(s|\mathcal{C})$, which is the sum of the scores of most confident x in each position:

$$\text{totalSeqConf}(s|\mathcal{C}) = \sum_k \max_x \text{totalConf}(x, k, s|\mathcal{C}) \quad (2)$$

The functionality and density features are combined using a Naive Bayes classifier. To train the classifier, we use a set of sequence names s labeled as either correct or incorrect, which we describe in Section 3.

Classifying Sequence Items To classify (x, k) given s , SEQ uses two features: the total confidence $\text{totalConf}(x, k, s|\mathcal{C})$ and the same total confidence normalized to sum to 1 over all x , holding k and s constant. To train the classifier, we use a set of extractions (x, k, s) where s is known to be a correct sequence name.

3 Experimental Results

This section reports on two experiments. First, we measured how the density and functionality features improve performance on the sequence name

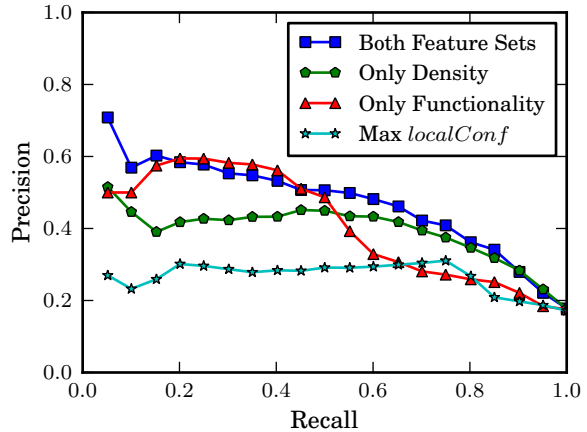


Figure 1: Using density or functionality features alone is effective in identifying correct sequence names. Combining both types of features outperforms either by a statistically significant margin (paired t-test, $p < 0.05$).

classification sub-task (Figure 1). Second, we report on SEQ’s performance on the sequence-extraction task (Figure 2).

To create a test set, we selected all sentences containing ordinal phrases from Banko’s 500M Web page corpus (2008). To enrich this set O , we obtained additional sentences from Bing.com as follows. For each sequence name s satisfying $\text{localConf}(x, k, s|\text{sentence}) \geq 0.5$ for some sentence in O , we queried Bing.com for “the k th s ” for $k = 1, 2, \dots$ until no more hits were returned.² For each query, we downloaded the search snippets and added them to our corpus. This procedure resulted in making 95,611 search engine queries. The final corpus contained 3,716,745 distinct sentences containing an OP.

Generating candidate extractions using the method from Section 2.1 resulted in a set of over 40 million distinct extractions, the vast majority of which are incorrect. To get a sample with a significant number of correct extractions, we filtered this set to include only extractions with $\text{totalConf}(x, k, s|\mathcal{C}) \geq 0.8$ for some sentence, resulting in a set of 2,409,211 extractions.

We then randomly sampled and manually labeled 2,000 of these extractions for evaluation. We did a Web search to verify the correctness of the sequence name s and that x is the k th item in the sequence. In some cases, the ordering relation of the sequence name was ambiguous (e.g.,

²We queried for both the numeric form of the ordinal and the number spelled out (e.g. “the 2nd ...” and “the second ...”). We took up to 100 results per query.

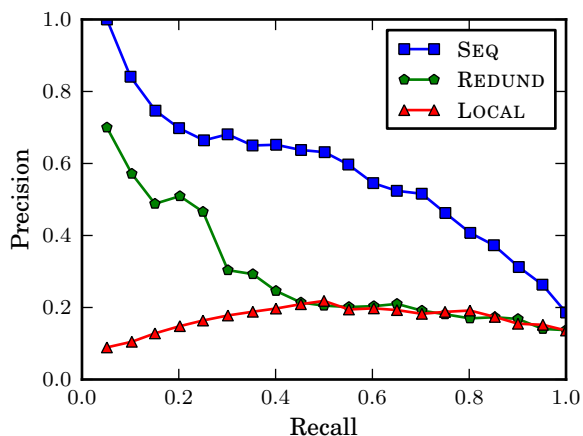


Figure 2: SEQ outperforms the baseline systems, increasing the area under the curve by 247% relative to LOCAL and by 90% relative to REDUND.

“largest state in the US” could refer to land area or population), which could lead to merging two distinct sequences. In practice, we found that most ordering relations were used in a consistent way (e.g., “largest city in” always means largest by population) and only about 5% of the sequence names in our sample have an ambiguous ordering relation.

We compute precision-recall curves relative to this random sample by changing a confidence threshold. Precision is the percentage of correct extractions above a threshold, while recall is the percentage correct above a threshold divided by the total number of correct extractions. Because SEQ requires training data, we used 15-fold cross validation on the labeled sample.

The functionality and density features boost SEQ’s ability to correctly identify sequence names. Figure 1 shows how well SEQ can identify correct sequence names using only functionality, only density, and using functionality and density in concert. The baseline used is the maximum value of $localConf(x, k, s)$ over all (x, k) . Both the density features and the functionality features are effective at this task, but using both types of features resulted in a statistically significant improvement over using either type of feature individually (paired t-test of area under the curve, $p < 0.05$).

We measure SEQ’s efficacy on the complete sequence-extraction task by contrasting it with two baseline systems. The first is LOCAL, which ranks extractions by $localConf$.³ The second is

³If an extraction arises from multiple sentences, we use

REDUND, which ranks extractions by $totalConf$. Figure 2 shows the precision-recall curves for each system on the test data. The area under the curves for SEQ, REDUND, and LOCAL are 0.59, 0.31, and 0.17, respectively. The low precision and flat curve for LOCAL suggests that $localConf$ is not informative for classifying extractions on its own.

REDUND outperformed LOCAL, especially at the high-precision part of the curve. On the subset of extractions with correct s , REDUND can identify x as the k th item with precision of 0.85 at recall 0.80. This is consistent with previous work on redundancy-based extractors on the Web. However, REDUND still suffered from the problems of over-specification and over-generalization described in Section 2. SEQ reduces the negative effects of these problems by decreasing the scores of sequence names that appear too general or too specific.

4 Related Work

There has been extensive work in extracting lists or sets of entities from the Web. These extractors rely on either (1) HTML features (Cohen et al., 2002; Wang and Cohen, 2007) to extract from structured text or (2) lexico-syntactic patterns (Hearst, 1992; Etzioni et al., 2005) to extract from unstructured text. SEQ is most similar to this second type of extractor, but additionally leverages the sequence regularities of functionality and density. These regularities allow the system to overcome the poor performance of the purely syntactic extractor LOCAL and the redundancy-based extractor REDUND.

5 Conclusions

We have demonstrated that an extractor leveraging sequence regularities can greatly outperform extractors without this knowledge. Identifying likely sequence names and *then* filling in sequence items proved to be an effective approach to sequence extraction.

One line of future research is to investigate other types of domain-independent frames that exhibit useful regularities. Other examples include *events* (with regularities about actor, location, and time) and a generic *organization-role* frame (with regularities about person, organization, and role played).

the maximal $localConf$.

6 Acknowledgements

This research was supported in part by NSF grant IIS-0803481, ONR grant N00014-08-1-0431, DARPA contract FA8750-09-C-0179, and an NSF Graduate Research Fellowship, and was carried out at the University of Washington's Turing Center.

References

- Michele Banko and Oren Etzioni. 2008. The tradeoffs between open and traditional relation extraction. In *Proceedings of ACL-08: HLT*, pages 28–36.
- Michele Banko, Michael J. Cafarella, Stephen Soderland, Matthew Broadhead, and Oren Etzioni. 2007. Open information extraction from the web. In *IJCAI*, pages 2670–2676.
- H. Chieu, H. Ng, and Y. Lee. 2003. Closing the gap: Learning-based information extraction rivaling knowledge-engineering methods. In *ACL*, pages 216–223.
- William W. Cohen, Matthew Hurst, and Lee S. Jensen. 2002. A flexible learning system for wrapping tables and lists in html documents. In *In International World Wide Web Conference*, pages 232–241.
- Doug Downey, Oren Etzioni, and Stephen Soderland. 2005. A probabilistic model of redundancy in information extraction. In *IJCAI*, pages 1034–1041.
- O. Etzioni, M. Cafarella, D. Downey, A. Popescu, T. Shaked, S. Soderland, D. Weld, and A. Yates. 2004. Methods for domain-independent information extraction from the Web: An experimental comparison. In *Proceedings of the Nineteenth National Conference on Artificial Intelligence (AAAI-2004)*, pages 391–398.
- Oren Etzioni, Michael Cafarella, Doug Downey, Ana maria Popescu, Tal Shaked, Stephen Soderl, Daniel S. Weld, and Er Yates. 2005. Unsupervised named-entity extraction from the web: An experimental study. *Artificial Intelligence*, 165:91–134.
- D. Freitag. 2000. Machine learning for information extraction in informal domains. *Machine Learning*, 39(2-3):169–202.
- Marti A. Hearst. 1992. Automatic acquisition of hyponyms from large text corpora. In *COLING*, pages 539–545.
- Satoshi Sekine. 2006. On-demand information extraction. In *Proceedings of the COLING/ACL on Main conference poster sessions*, pages 731–738, Morristown, NJ, USA. Association for Computational Linguistics.
- Richard C. Wang and William W. Cohen. 2007. Language-independent set expansion of named entities using the web. In *ICDM*, pages 342–350. IEEE Computer Society.

An Entity-Level Approach to Information Extraction

Aria Haghighi

UC Berkeley, CS Division
aria42@cs.berkeley.edu

Dan Klein

UC Berkeley, CS Division
klein@cs.berkeley.edu

Abstract

We present a generative model of template-filling in which coreference resolution and role assignment are jointly determined. Underlying template roles first generate abstract entities, which in turn generate concrete textual mentions. On the standard corporate acquisitions dataset, joint resolution in our entity-level model reduces error over a mention-level discriminative approach by up to 20%.

1 Introduction

Template-filling information extraction (IE) systems must merge information across multiple sentences to identify all role fillers of interest. For instance, in the MUC4 terrorism event extraction task, the entity filling the *individual perpetrator* role often occurs multiple times, variously as proper, nominal, or pronominal mentions. However, most template-filling systems (Freitag and McCallum, 2000; Patwardhan and Riloff, 2007) assign roles to individual textual mentions using only local context as evidence, leaving aggregation for post-processing. While prior work has acknowledged that coreference resolution and discourse analysis are integral to accurate role identification, to our knowledge no model has been proposed which jointly models these phenomena.

In this work, we describe an entity-centered approach to template-filling IE problems. Our model jointly merges surface mentions into underlying entities (coreference resolution) and assigns roles to those discovered entities. In the generative process proposed here, document entities are generated for each template role, along with a set of non-template entities. These entities then generate mentions in a process sensitive to both lexical and structural properties of the mention. Our model outperforms a discriminative mention-level baseline. Moreover, since our model is generative, it

Template			
SELLER	BUSINESS	ACQUIRED	PURCHASER
CSR Limited	Oil and Gas	Delhi Fund	Esso Inc.

(a)

Document					
[S CSR]	has said that	[S it]	has sold	[S its]	[B oil
interests]	held in	[A Delhi Fund].	[P Esso Inc.]	did not	disclose how much
[P they]	paid for	[A Dehli].			

(b)

Figure 1: Example of the corporate acquisitions role-filling task. In (a), an example template specifying the entities playing each domain role. In (b), an example document with coreferent mentions sharing the same role label. Note that pronoun mentions provide direct clues to entity roles.

can naturally incorporate unannotated data, which further increases accuracy.

2 Problem Setting

Figure 1(a) shows an example *template-filling* task from the corporate acquisitions domain (Freitag, 1998).¹ We have a template of K roles (PURCHASER, AMOUNT, etc.) and we must identify which entity (if any) fills each role (*CSR Limited*, etc.). Often such problems are modeled at the mention level, directly labeling individual mentions as in Figure 1(b). Indeed, in this data set, the mention-level perspective is evident in the gold annotations, which ignore pronominal references. However, roles in this domain appear in several locations throughout the document, with pronominal mentions often carrying the critical information for template filling. Therefore, Section 3 presents a model in which entities are explicitly modeled, naturally merging information across all mention types and explicitly representing latent structure very much like the entity-level template structure from Figure 1(a).

¹In Freitag (1998), some of these fields are split in two to distinguish a full versus abbreviated name, but we ignore this distinction. Also we ignore the *status* field as it doesn't apply to entities and its meaning is not consistent.

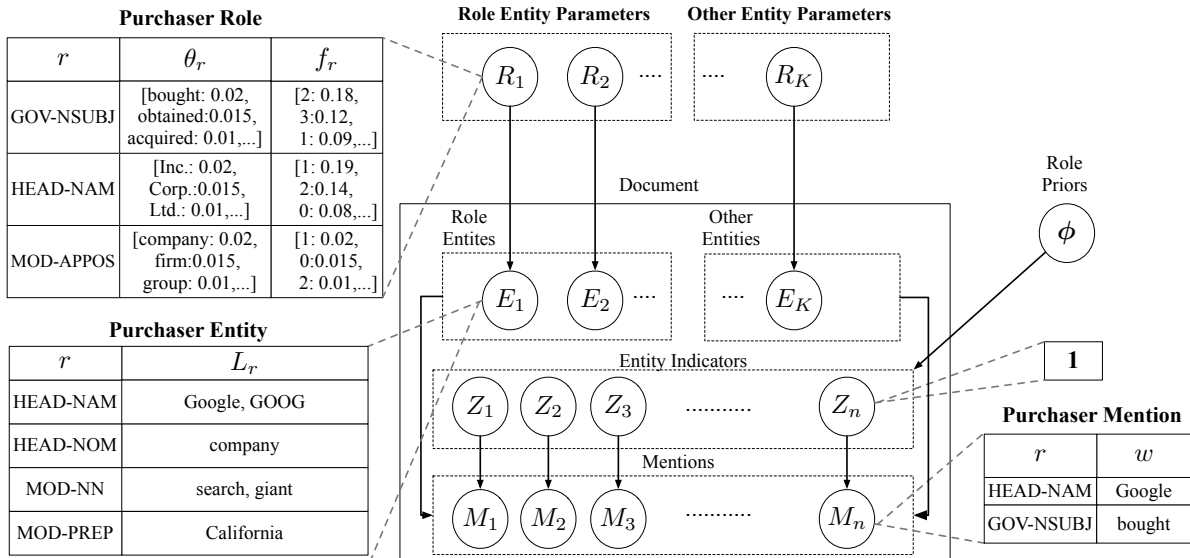


Figure 2: Graphical model depiction of our generative model described in Section 3. Sample values are illustrated for key parameters and latent variables.

3 Model

We describe our generative model for a document, which has many similarities to the coreference-only model of Haghighi and Klein (2010), but which integrally models template role-fillers. We briefly describe the key abstractions of our model.

Mentions: A mention is an observed textual reference to a latent real-world entity. Mentions are associated with nodes in a parse tree and are typically realized as NPs. There are three basic forms of mentions: proper (NAM), nominal (NOM), and pronominal (PRO). Each mention M is represented as collection of key-value pairs. The keys are called *properties* and the values are words. The set of properties utilized here, denoted \mathcal{R} , are the same as in Haghighi and Klein (2010) and consist of the mention head, its dependencies, and its governor. See Figure 2 for a concrete example. Mention types are trivially determined from mention head POS tag. All mention properties and their values are observed.

Entities: An entity is a specific individual or object in the world. Entities are always latent in text. Where a mention has a single word for each property, an entity has a *list* of signature words. Formally, entities are mappings from properties $r \in \mathcal{R}$ to lists L_r of “canonical” words which that entity uses for that property.

Roles: The elements we have described so far are standard in many coreference systems. Our model performs role-filling by assuming that each entity is drawn from an underlying role. These

roles include the K template roles as well as ‘junk’ roles to represent entities which do not fill a template role (see Section 5.2). Each role R is represented as a mapping between properties r and pairs of multinomials (θ_r, f_r) . θ_r is a unigram distribution of words for property r that are semantically licensed for the role (e.g., being the subject of “acquired” for the ACQUIRED role). f_r is a “fertility” distribution over the integers that characterizes entity list lengths. Together, these distributions control the lists L_r for entities which instantiate the role.

We first present a broad sketch of our model’s components and then detail each in a subsequent section. We temporarily assume that all mentions belong to a template role-filling entity; we lift this restriction in Section 5.2. First, a semantic component generates a sequence of entities $\mathbf{E} = (E_1, \dots, E_K)$, where each E_i is generated from a corresponding role R_i . We use $\mathbf{R} = (R_1, \dots, R_K)$ to denote the vector of template role parameters. Note that this work assumes that there is a one-to-one mapping between entities and roles; in particular, at most one entity can fill each role. This assumption is appropriate for the domain considered here.

Once entities have been generated, a discourse component generates which entities will be evoked in each of the n mention positions. We represent these choices using *entity indicators* denoted by $\mathbf{Z} = (Z_1, \dots, Z_n)$. This component utilizes a learned global prior ϕ over roles. The Z_i in-

dicators take values in $1, \dots, K$ indicating the entity number (and thereby the role) underlying the i th mention position. Finally, a mention generation component renders each mention conditioned on the underlying entity and role. Formally:

$$P(\mathbf{E}, \mathbf{Z}, \mathbf{M} | \mathbf{R}, \phi) =$$

$$\left(\prod_{i=1}^K P(E_i | R_i) \right) \quad \text{[Semantic, Sec. 3.1]}$$

$$\left(\prod_{j=1}^n P(Z_j | \mathbf{Z}_{<j}, \phi) \right) \quad \text{[Discourse, Sec. 3.2]}$$

$$\left(\prod_{j=1}^n P(M_j | E_{Z_j}, R_{Z_j}) \right) \quad \text{[Mention, Sec. 3.3]}$$

3.1 Semantic Component

Each role R generates an entity E as follows: for each mention property r , a word list, L_r , is drawn by first generating a list length from the corresponding f_r distribution in R .² This list is then populated by an independent draw from R 's unigram distribution θ_r . Formally, for each $r \in \mathcal{R}$, an entity word list is drawn according to,³

$$P(L_r | R) = P(\text{len}(L_r) | f_r) \prod_{w \in L_r} P(w | \theta_r)$$

3.2 Discourse Component

The discourse component draws the entity indicator Z_j for the j th mention according to,

$$P(Z_j | \mathbf{Z}_{<j}, \phi) = \begin{cases} P(Z_j | \phi), & \text{if non-pronominal} \\ \sum_{j'} \mathbf{1}[Z_j = Z_{j'}] P(j' | j), & \text{o.w.} \end{cases}$$

When the j th mention is non-pronominal, we draw Z_j from ϕ , a global prior over the K roles. When M_j is a pronoun, we first draw an antecedent mention position j' , such that $j' < j$, and then we set $Z_j = Z_{j'}$. The antecedent position is selected according to the distribution,

$$P(j' | j) \propto \exp\{-\gamma \text{TREEDIST}(j', j)\}$$

where $\text{TREEDIST}(j', j)$ represents the tree distance between the parse nodes for M_j and $M_{j'}$.⁴ Mass is

²There is one exception: the sizes of the proper and nominal head property lists are jointly generated, but their word lists are still independently populated.

³While, in principle, this process can yield word lists with duplicate words, we constrain the model during inference to not allow that to occur.

⁴Sentence parse trees are merged into a right-branching document parse tree. This allows us to extend tree distance to inter-sentence nodes.

restricted to antecedent mention positions j' which occur earlier in the same sentence or in the previous sentence.⁵

3.3 Mention Generation

Once the entity indicator has been drawn, we generate words associated with mention conditioned on the underlying entity E and role R . For each mention property r associated with the mention, a word w is drawn utilizing E 's word list L_r as well as the multinomials (f_r, θ_r) from role R . The word w is drawn according to,

$$P(w | E, R) = (1 - \alpha_r) \frac{\mathbf{1}[w \in L_r]}{\text{len}(L_r)} + \alpha_r P(w | \theta_r)$$

For each property r , there is a hyper-parameter α_r which interpolates between selecting a word uniformly from the entity list L_r and drawing from the underlying role distribution θ_r . Intuitively, a small α_r indicates that an entity prefers to re-use a small number of words for property r . This is typically the case for proper and nominal heads as well as modifiers. At the other extreme, setting α_r to 1 indicates the property isn't particular to the entity itself, but rather always drawn from the underlying role distribution. We set α_r to 1 for pronoun heads as well as for the governor properties.

4 Learning and Inference

Since we will make use of unannotated data (see Section 5), we utilize a variational EM algorithm to learn parameters \mathbf{R} and ϕ . The E-Step requires the posterior $P(\mathbf{E}, \mathbf{Z} | \mathbf{R}, \mathbf{M}, \phi)$, which is intractable to compute exactly. We approximate it using a surrogate variational distribution of the following factored form:

$$Q(\mathbf{E}, \mathbf{Z}) = \left(\prod_{i=1}^K q_i(E_i) \right) \left(\prod_{j=1}^n r_j(Z_j) \right)$$

Each $r_j(Z_j)$ is a distribution over the entity indicator for mention M_j , which approximates the true posterior of Z_j . Similarly, $q_i(E_i)$ approximates the posterior over entity E_i which is associated with role R_i . As is standard, we iteratively update each component distribution to minimize KL-divergence, fixing all other distributions:

$$q_i \leftarrow \underset{q_i}{\text{argmin}} KL(Q(\mathbf{E}, \mathbf{Z}) | P(\mathbf{E}, \mathbf{Z} | \mathbf{M}, \mathbf{R}, \phi))$$

$$\propto \exp\{\mathbb{E}_{Q/q_i} \ln P(\mathbf{E}, \mathbf{Z} | \mathbf{M}, \mathbf{R}, \phi)\}$$

⁵The sole parameter γ is fixed at 0.1.

	Ment Acc.	Ent. Acc.
INDEP	60.0	43.7
JOINT	64.6	54.2
JOINT+PRO	68.2	57.8

Table 1: Results on corporate acquisition tasks with given role mention boundaries. We report mention role accuracy and entity role accuracy (correctly labeling all entity mentions).

For example, the update for a non-pronominal entity indicator component $r_j(\cdot)$ is given by:⁶

$$\begin{aligned} \ln r_j(z) &\propto \mathbb{E}_{Q/r_j} \ln P(\mathbf{E}, \mathbf{Z}, \mathbf{M} | \mathbf{R}, \phi) \\ &\propto \mathbb{E}_{q_z} \ln (P(z|\phi) P(M_j | E_z, R_z)) \\ &= \ln P(z|\phi) + \mathbb{E}_{q_z} \ln P(M_j | E_z, R_z) \end{aligned}$$

A similar update is performed on pronominal entity indicator distributions, which we omit here for space. The update for variational entity distribution is given by:

$$\begin{aligned} \ln q_i(e_i) &\propto \mathbb{E}_{Q/q_i} \ln P(\mathbf{E}, \mathbf{Z}, \mathbf{M} | \mathbf{R}, \phi) \\ &\propto \mathbb{E}_{\{r_j\}} \ln \left(P(e_i | R_i) \prod_{j: Z_j=i} P(M_j | e_i, R_i) \right) \\ &= \ln P(e_i | R_i) + \sum_j r_j(i) \ln P(M_j | e_i, R_i) \end{aligned}$$

It is intractable to enumerate all possible entities e_i (each consisting of several sets of words). We instead limit the support of $q_i(e_i)$ to several sampled entities. We obtain entity samples by sampling mention entity indicators according to r_j . For a given sample, we assume that E_i consists of the non-pronominal head words and modifiers of mentions such that Z_j has sampled value i .

During the E-Step, we perform 5 iterations of updating each variational factor, which results in an approximate posterior distribution. Using expectations from this approximate posterior, our M-Step is relatively straightforward. The role parameters R_i are computed from the $q_i(e_i)$ and $r_j(z)$ distributions, and the global role prior ϕ from the non-pronominal components of $r_j(z)$.

5 Experiments

We present results on the corporate acquisitions task, which consists of 600 annotated documents split into a 300/300 train/test split. We use 50 training documents as a development set. In all

⁶For simplicity of exposition, we omit terms where M_j is an antecedent to a pronoun.

documents, proper and (usually) nominal mentions are annotated with roles, while pronouns are not. We preprocess each document identically to Haghighi and Klein (2010): we sentence-segment using the OpenNLP toolkit, parse sentences with the Berkeley Parser (Petrov et al., 2006), and extract mention properties from parse trees and the Stanford Dependency Extractor (de Marneffe et al., 2006).

5.1 Gold Role Boundaries

We first consider the simplified task where role mention boundaries are given. We map each labeled token span in training and test data to a parse tree node that shares the same head. In this setting, the role-filling task is a collective classification problem, since we know each mention is filling some role.

As our baseline, INDEP, we built a maximum entropy model which independently classifies each mention’s role. It uses features as similar as possible to the generative model (and more), including the head word, typed dependencies of the head, various tree features, governing word, and several conjunctions of these features as well as coarser versions of lexicalized features. This system yields 60.0 mention labeling accuracy (see Table 1). The primary difficulty in classification is the disambiguation amongst the acquired, seller, and purchaser roles, which have similar internal structure, and differ primarily in their semantic contexts. Our entity-centered model, JOINT in Table 1, has no latent variables at training time in this setting, since each role maps to a unique entity. This model yields 64.6, outperforming INDEP.⁷

During development, we noted that often the most direct evidence of the role of an entity was associated with pronoun usage (see the first “it” in Figure 1). Training our model with pronominal mentions, whose roles are latent variables at training time, improves accuracy to 68.2.⁸

5.2 Full Task

We now consider the more difficult setting where role mention boundaries are not provided at test time. In this setting, we automatically extract mentions from a parse tree using a heuristic ap-

⁷We use the mode of the variational posteriors $r_j(Z_j)$ to make predictions (see Section 4).

⁸While this approach incorrectly assumes that all pronouns have antecedents amongst our given mentions, this did not appear to degrade performance.

	ROLE ID			OVERALL		
	P	R	F ₁	P	R	F ₁
INDEP	79.0	65.5	71.6	48.6	40.3	44.0
JOINT+PRO	80.3	69.2	74.3	53.4	46.4	49.7
BEST	80.1	70.1	74.8	57.3	49.2	52.9

Table 2: Results on corporate acquisitions data where mention boundaries are not provided. Systems must determine which mentions are template role-fillers as well as label them. ROLE ID only evaluates the binary decision of whether a mention is a template role-filler or not. OVERALL includes correctly labeling mentions. Our BEST system, see Section 5, adds extra unannotated data to our JOINT+PRO system.

proach. Our mention extraction procedure yields 95% recall over annotated role mentions and 45% precision.⁹ Using extracted mentions as input, our task is to label some subset of the mentions with template roles. Since systems can label mentions as non-role bearing, only recall is critical to mention extraction. To adapt INDEP to this setting, we first use a binary classifier trained to distinguish role-bearing mentions. The baseline then classifies mentions which pass this first phase as before. We add ‘junk’ roles to our model to flexibly model entities that do not correspond to annotated template roles. During training, extracted mentions which are not matched in the labeled data have posteriors which are constrained to be amongst the ‘junk’ roles.

We first evaluate role identification (ROLE ID in Table 2), the task of identifying mentions which play some role in the template. The binary classifier for INDEP yields 71.6 F₁. Our JOINT+PRO system yields 74.3. On the task of identifying and correctly labeling role mentions, our model outperforms INDEP as well (OVERALL in Table 2). As our model is generative, it is straightforward to utilize totally unannotated data. We added 700 fully unannotated documents from the mergers and acquisitions portion of the Reuters 21857 corpus. Training JOINT+PRO on this data as well as our original training data yields the best performance (BEST in Table 2).¹⁰

To our knowledge, the best previously published results on this dataset are from Siefkes (2008), who report 45.9 weighted F₁. Our BEST system evaluated in their slightly stricter way yields 51.1.

⁹Following Patwardhan and Riloff (2009), we match extracted mentions to labeled spans if the head of the mention matches the labeled span.

¹⁰We scaled expected counts from the unlabeled data so that they did not overwhelm those from our (partially) labeled data.

6 Conclusion

We have presented a joint generative model of coreference resolution and role-filling information extraction. This model makes role decisions at the entity, rather than at the mention level. This approach naturally aggregates information across multiple mentions, incorporates unannotated data, and yields strong performance.

Acknowledgements: This project is funded in part by the Office of Naval Research under MURI Grant No. N000140911081.

References

- M. C. de Marneffe, B. Maccartney, and C. D. Manning. 2006. Generating typed dependency parses from phrase structure parses. In *LREC*.
- Dayne Freitag and Andrew McCallum. 2000. Information extraction with hmm structures learned by stochastic optimization. In *Association for the Advancement of Artificial Intelligence (AAAI)*.
- Dayne Freitag. 1998. Machine learning for information extraction in informal domains.
- A. Haghighi and D. Klein. 2010. Coreference resolution in a modular, entity-centered model. In *North American Association of Computational Linguistics (NAACL)*.
- P. Liang and D. Klein. 2007. Structured Bayesian non-parametric models with variational inference (tutorial). In *Association for Computational Linguistics (ACL)*.
- S. Patwardhan and E. Riloff. 2007. Effective information extraction with semantic affinity patterns and relevant regions. In *Joint Conference on Empirical Methods in Natural Language Processing*.
- S. Patwardhan and E. Riloff. 2009. A unified model of phrasal and sentential evidence for information extraction. In *Empirical Methods in Natural Language Processing (EMNLP)*.
- Slav Petrov, Leon Barrett, Romain Thibaux, and Dan Klein. 2006. Learning accurate, compact, and interpretable tree annotation. In *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics*, pages 433–440, Sydney, Australia, July. Association for Computational Linguistics.
- Christian Siefkes. 2008. *An Incrementally Trainable Statistical Approach to Information Extraction: Based on Token Classification and Rich Context Model*. VDM Verlag, Saarbrücken, Germany, Germany.

A Semi-Supervised Key Phrase Extraction Approach: Learning from Title Phrases through a Document Semantic Network

Decong Li¹, Sujian Li¹, Wenjie Li², Wei Wang¹, Weiguang Qu³

¹Key Laboratory of Computational Linguistics, Peking University

²Department of Computing, The Hong Kong Polytechnic University

³School of Computer Science and Technology, Nanjing Normal University

{lidecong,lisujian, wwei }@pku.edu.cn cswjli@comp.polyu.edu.hk wgqu@njnu.edu.cn

Abstract

It is a fundamental and important task to extract key phrases from documents. Generally, phrases in a document are not independent in delivering the content of the document. In order to capture and make better use of their relationships in key phrase extraction, we suggest exploring the Wikipedia knowledge to model a document as a semantic network, where both n -ary and binary relationships among phrases are formulated. Based on a commonly accepted assumption that the title of a document is always elaborated to reflect the content of a document and consequently key phrases tend to have close semantics to the title, we propose a novel semi-supervised key phrase extraction approach in this paper by computing the phrase importance in the semantic network, through which the influence of title phrases is propagated to the other phrases iteratively. Experimental results demonstrate the remarkable performance of this approach.

1 Introduction

Key phrases are defined as the phrases that express the main content of a document. Guided by the given key phrases, people can easily understand what a document describes, saving a great amount of time reading the whole text. Consequently, automatic key phrase extraction is in high demand. Meanwhile, it is also fundamental to many other natural language processing applications, such as information retrieval, text clustering and so on.

Key phrase extraction can be normally cast as a ranking problem solved by either supervised or unsupervised methods. Supervised learning requires a large amount of expensive training data, whereas unsupervised learning totally ignores human knowledge. To overcome the deficiencies

of these two kinds of methods, we propose a novel semi-supervised key phrase extraction approach in this paper, which explores title phrases as the source of knowledge.

It is well agreed that the title has a similar role to the key phrases. They are both elaborated to reflect the content of a document. Therefore, phrases in the titles are often appropriate to be key phrases. That is why position has been a quite effective feature in the feature-based key phrase extraction methods (Witten, 1999), i.e., if a phrase is located in the title, it is ranked higher.

However, one can only include a couple of most important phrases in the title prudently due to the limitation of the title length, even though many other key phrases are all pivotal to the understanding of the document. For example, when we read the title “China Tightens Grip on the Web”, we can only have a glimpse of what the document says. On the other hand, the key phrases, such as “China”, “Censorship”, “Web”, “Domain name”, “Internet”, and “CNNIC”, etc. can tell more details about the main topics of the document. In this regard, title phrases are often good key phrases but they are far from enough.

If we review the above example again, we will find that the key phrase “Internet” can be inferred from the title phrase “Web”. As a matter of fact, key phrases often have close semantics to title phrases. Then a question comes to our minds: can we make use of these title phrases to infer the other key phrases?

To provide a foundation of inference, a semantic network that captures the relationships among phrases is required. In the previous works (Turdakov and Velikhov, 2008), semantic networks are constructed based on the binary relations, and the semantic relatedness between a pair of phrases is formulated by the weighted edges that connects them. The deficiency of these approaches is the incapability to capture the n -ary relations among multiple phrases. For example, a group of

phrases may collectively describe an entity or an event.

In this study, we propose to model a semantic network as a hyper-graph, where vertices represent phrases and weighted hyper-edges measure the semantic relatedness of both binary relations and n -ary relations among phrases. We explore a universal knowledge base – Wikipedia – to compute the semantic relatedness. Yet our major contribution is to develop a novel semi-supervised key phrase extraction approach by computing the phrase importance in the semantic network, through which the influence of title phrases is propagated to the other phrases iteratively.

The goal of the semi-supervised learning is to design a function that is sufficiently smooth with respect to the intrinsic structure revealed by title phrases and other phrases. Based on the assumption that semantically related phrases are likely to have similar scores, the function to be estimated is required to assign title phrases a higher score and meanwhile locally smooth on the constructed hyper-graph. Zhou et al.’s work (Zhou 2005) lays down a foundation for our semi-supervised phrase ranking algorithm introduced in Section 3. Experimental results presented in Section 4 demonstrate the effectiveness of this approach.

2 Wikipedia-based Semantic Network Construction

Wikipedia¹ is a free online encyclopedia, which has unarguably become the world’s largest collection of encyclopedic knowledge. *Articles* are the basic entries in the Wikipedia, with each article explaining one Wikipedia term. Articles contain *links* pointing from one article to another. Currently, there are over 3 million articles and 90 million links in English Wikipedia. In addition to providing a large vocabulary, Wikipedia articles also contain a rich body of lexical semantic information expressed via the extensive number of links. During recent years, Wikipedia has been used as a powerful tool to compute semantic relatedness between terms in a good few of works (Turdakov 2008).

We consider a document composed of the phrases that describe various aspects of entities or events with different semantic relationships. We then model a document as a semantic network formulated by a weighted hyper-graph

$G=(V, E, W)$, where each vertex $v_i \in V$ ($1 \leq i \leq n$) represents a phrase, each hyper-edge $e_j \in E$ ($1 \leq j \leq m$) is a subset of V , representing binary relations or n -ary relations among phrases, and the weight $w(e_j)$ measures the semantic relatedness of e_j .

By applying the WSD technique proposed by (Turdakov and Velikhov, 2008), each phrase is assigned with a single Wikipedia article that describes its meaning. Intuitively, if the fraction of the links that the two articles have in common to the total number of the links in both articles is high, the two phrases corresponding to the two articles are more semantically related. Also, an article contains different types of links, which are relevant to the computation of semantic relatedness to different extent. Hence we adopt the weighted Dice metric proposed by (Turdakov 2008) to compute the semantic relatedness of each binary relation, resulting in the edge weight $w(e_{ij})$, where e_{ij} is an edge connecting the phrases v_i and v_j .

To define the n -ary relations in the semantic network, a proper graph clustering technique is needed. We adopt the weighted Girvan-Newman algorithm (Newman 2004) to cluster phrases (including title phrases) by computing their betweenness centrality. The advantage of this algorithm is that it need not specify a pre-defined number of clusters. Then the phrases, within each cluster, are connected by a n -ary relation. n -ary relations among the phrases in the same cluster are then measured based on binary relations. The weight of a hyper-edge e is defined as:

$$w(e) = \frac{\alpha}{|e|} \sum_{e_{ij} \subseteq e} w(e_{ij}) \quad (1)$$

where $|e|$ is the number of the vertices in e , e_{ij} is an edge with two vertices included in e and $\alpha \geq 0$ is a parameter balancing the relative importance of n -ary hyper-edges compared with binary ones.

3 Semi-supervised Learning from Title

Given the document semantic network represented as a phrase hyper-graph, one way to make better use of the semantic information is to rank phrases with a semi-supervised learning strategy, where the title phrases are regarded as labeled samples, while the other phrases as unlabeled ones. That is, the information we have at the beginning about how to rank phrases is that the title phrases are the most important phrases. Initially, the title phrases are assigned with a positive score of 1 indicating its importance and oth-

¹ www.wikipedia.org

er phrases are assigned zero. Then the importance scores of the phrases are learned iteratively from the title phrases through the hyper-graph. The key idea behind hyper-graph based semi-supervised ranking is that the vertices which usually belong to the same hyper-edges should be assigned with similar scores. Then, we have the following two constraints:

1. The phrases which have many incident hyper-edges in common should be assigned similar scores.

2. The given initial scores of the title phrases should be changed as little as possible.

Given a weighted hyper-graph G , assume a ranking function f over V , which assigns each vertex v an importance score $f(v)$. f can be thought as a vector in Euclid space $R^{|V|}$. For the convenience of computation, we use an incidence matrix H to represent the hypergraph, defined as:

$$h(v, e) = \begin{cases} 0, & \text{if } v \notin e \\ 1, & \text{if } v \in e \end{cases} \quad (2)$$

Based on the incidence matrix, we define the degrees of the vertex v and the hyper-edge e as

$$d(v) = \sum_{e \in E} w(e) h(v, e) \quad (3)$$

and

$$\delta(e) = \sum_{v \in V} h(v, e) \quad (4)$$

Then, to formulate the above-mentioned constraints, let y denote the initial score vector, then the importance scores of the phrases are learned iteratively by solving the following optimization problem:

$$\arg \min_{f \in R^{|V|}} \{ \Omega(f) + \mu \|f - y\|^2 \} \quad (5)$$

$$\Omega(f) = \frac{1}{2} \sum_{e \in E} \frac{1}{\delta(e)} \sum_{\{u, v\} \subseteq e} w(e) \left(\frac{f(u)}{\sqrt{d(u)}} - \frac{f(v)}{\sqrt{d(v)}} \right)^2 \quad (6)$$

where $\mu > 0$ is the parameter specifying the tradeoff between the two competitive items. Let D_v and D_e denote the diagonal matrices containing the vertex and the hyper-edge degrees respectively, W denote the diagonal matrix containing the hyper-edge weights, f^* denote the solution of (6). Zhou has given the solution (Zhou, 2005) as.

$$f^* = \beta \Theta f^* + (1 - \beta) y \quad (7)$$

where $\Theta = D_v^{-1/2} H W D_e^{-1} H^T D_v^{-1/2}$ and $\beta = 1 / (\mu + 1)$. Using an approximation algorithm (e.g. Algorithm 1), we can finally get a vector f representing the approximate phrase scores.

Algorithm 1: PhraseRank(V, T, a, b)

Input: Title phrase set = $\{v_1, v_2, \dots, v_t\}$, the set of other phrases = $\{v_{t+1}, v_{t+2}, \dots, v_n\}$, parameters α and β , con-

vergence threshold ζ

Output: The approximate phrase scores f

Construct a document semantic network for all the phrases $\{v_1, v_2, \dots, v_n\}$ using the method described in section 2.

Let $\Theta = D_v^{-1/2} H W D_e^{-1} H^T D_v^{-1/2}$;

Initialize the score vector y as $y_i = 1, 1 \leq i \leq t$, and

$y_j = 0, t < j \leq n$;

Let $f^0 = y, k = 0$;

REPEAT

$f^{k+1} = \beta \Theta f^k + (1 - \beta) y$;

$\nabla \leftarrow \max_i |f_i^{k+1} - f_i^k|$, for $0 \leq i \leq n$;

$k \leftarrow k + 1$;

UNTIL $\nabla < \zeta$

END

Finally we rank phrases in descending order of the calculated importance scores and select those highest ranked phrases as key phrases. According to the number of all the candidate phrases, we choose an appropriate proportion, i.e. 10%, of all the phrases as key phrases.

4 Evaluation

4.1 Experiment Set-up

We first collect all the Wikipedia terms to compose of a dictionary. The word sequences that occur in the dictionary are identified as phrases. Here we use a finite-state automaton to accomplish this task to avoid the imprecision of pre-processing by POS tagging or chunking. Then, we adopt the WSD technique proposed by (Turakov and Velikhov 2008) to find the corresponding Wikipedia article for each phrase. As mentioned in Section 2, a document semantic network in the form of a hyper-graph is constructed, on which Algorithm 1 is applied to rank the phrases.

To evaluate our proposed approach, we select 200 pieces of news from well-known English media. 5 to 10 key phrases are manually labeled in each news document and the average number of the key phrases is 7.2 per document. Due to the abbreviation and synonymy phenomena, we construct a thesaurus and convert all manual and automatic phrases into their canonical forms when evaluated. The traditional Recall, Precision and F1-measure metrics are adopted for evaluation. This section conducts two sets of experiment: (1) to examine the influence of two parameters: α and β , on the key phrase extraction performance; (2) to compare with other well known state-of-art key phrase extraction approaches.

4.2 Parameter tuning

The approach involves two parameters: α ($\alpha \geq 0$) is a relation factor balancing the influence of n -ary relations and binary relations; β ($0 \leq \beta \leq 1$) is a learning factor tuning the influence from the title phrases. It is hard to find a global optimized solution for the combination of these two factors. So we apply a gradient search strategy. At first, the learning factor is set to $\beta=0.8$. Different values of α ranging from 0 to 3 are examined. Then, given that α is set to the value with the best performance, we conduct experiments to find an appropriate value for β .

4.2.1 α : Relation Factor

First, we fix the learning factor β as 0.8 randomly and evaluate the performance by varying α value from 0 to 3. When $\alpha=0$, it means that the weight of n -ary relations is zero and only binary relations are considered. As we can see from Figure 1, the performance is improved in most cases in terms of F1-measure and reaches a peak at $\alpha=1.8$. This justifies the rationale to incorporate n -ary relations with binary relations in the document semantic network.

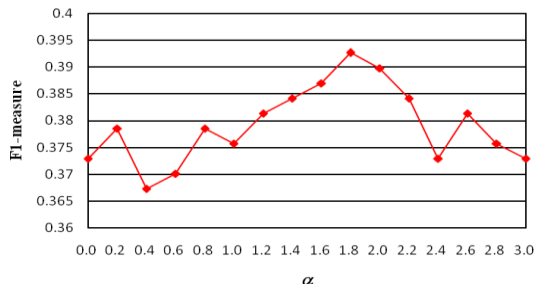


Figure 1. F1-measures with α in [0, 3]

4.2.2 β : Learning factor

Next, we set the relation factor $\alpha=1.8$, we inspect the performance with the learning factor β ranging from 0 to 1. $\beta=1$ means that the ranking scores learn from the semantic network without any consideration of title phrases. As shown in Figure 2, we find that the performance almost keep a smooth fluctuation as β increases from 0 to 0.9, and then a diving when $\beta=1$. This proves that title phrases indeed provide valuable information for learning.

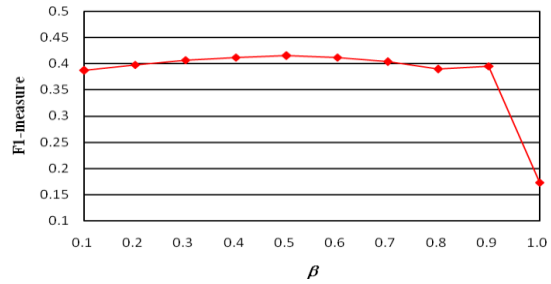


Figure 2. F1-measure with β in [0, 1]

4.3 Comparison with Other Approaches

Our approach aims at inferring important key phrases from title phrases through a semantic network. Here we take a method of synonym expansion as the baseline, called WordNet expansion here. The WordNet² expansion approach selects all the synonyms of the title phrases in the document as key phrases. Afterwards, our approach is evaluated against two existing approaches, which rely on the conventional semantic network and are able to capture binary relations only. One approach combines the title information into the Grineva's community-based method (Grineva *et al.*, 2009), called title-community approach. The title-community approach uses the Girvan-Newman algorithm to cluster phrases into communities and selects those phrases in the communities containing the title phrases as key phrases. We do not limit the number of key phrases selected. The other one is based on topic-sensitive LexRank (Otterbacher *et al.*, 2005), called title-sensitive PageRank here. The title-sensitive PageRank approach makes use of title phrases to re-weight the transitions between vertices and picks up 10% top-ranked phrases as key phrases.

Approach	Precision	Recall	F1
Title-sensitive PageRank ($d=0.15$)	34.8%	39.5%	37.0%
Title-community	29.8%	56.9%	39.1%
Our approach ($\alpha=1.8, \beta=0.5$)	39.4%	44.6%	41.8%
WordNet expansion (baseline)	7.9%	32.9%	12.5%

Table 1. Comparison with other approaches

Table 1 summarizes the performance on the test data. The results presented in the table show that our approach exhibits the best performance among all the four approaches. It follows that the key phrases inferred from a document semantic network are not limited to the synonyms of title phrases. As the title-sensitive PageRank ap-

² <http://wordnet.princeton.edu>

proach totally ignores the n -ary relations, its performance is the worst. Based on binary relations, the title-community approach clusters phrases into communities and each community can be considered as an n -ary relation. However, this approach lacks of an importance propagation process. Consequently, it has the highest recall value but the lowest precision. In contrast, our approach achieves the highest precision, due to its ability to infer many correct key phrases using importance propagation among n -ary relations.

5 Conclusion

This work is based on the belief that key phrases tend to have close semantics to the title phrases. In order to make better use of phrase relations in key phrase extraction, we explore the Wikipedia knowledge to model one document as a semantic network in the form of hyper-graph, through which the other phrases learned their importance scores from the title phrases iteratively. Experimental results demonstrate the effectiveness and robustness of our approach.

Acknowledgments

The work described in this paper was partially supported by NSFC programs (No: 60773173, 60875042 and 90920011), and Hong Kong RGC Projects (No: PolyU5217/07E). We thank the anonymous reviewers for their insightful comments.

References

- David Milne, Ian H. Witten. 2008. *An Effective, Low-Cost Measure of Semantic Relatedness Obtained from Wikipedia Links*. In Wikipedia and AI workshop at the AAAI-08 Conference, Chicago, US.
- Dengyong Zhou, Jiayuan Huang and Bernhard Schölkopf. 2005. *Beyond Pairwise Classification and Clustering Using Hypergraphs*. MPI Technical Report, Tübingen, Germany.
- Denis Turdakov and Pavel Velikhov. 2008. *Semantic relatedness metric for wikipedia concepts based on link analysis and its application to word sense disambiguation*. In Colloquium on Databases and Information Systems (SYRCODIS).
- Ian H. Witten, Gordon W. Paynter, Eibe Frank, Carl Gutwin, Craig G. Nevill-Manning. 1999. *KEA: practical automatic keyphrase extraction*, In Proceedings of the fourth ACM conference on Digital libraries, pp.254-255, California, USA.

Jahna Otterbacher, Gunes Erkan and Dragomir R. Radev. 2005. *Using Random Walks for Question-focused Sentence Retrieval*. In Proceedings of HLT/EMNLP 2005, pp. 915-922, Vancouver, Canada.

Maria Grineva, Maxim Grinev and Dmitry Lizorkin. 2009. *Extracting key terms from noisy and multitheme documents*, In Proceedings of the 18th international conference on World wide web, pp. 661-670, Madrid, Spain.

Michael Strube and Simone Paolo Ponzetto. 2006. *WikiRelate! Computing Semantic Relatedness using Wikipedia*. In Proceedings of the 21st National Conference on Artificial Intelligence, pp. 1419-1424, Boston, MA.

M. E. J. Newman. 2004. *Analysis of Weighted Networks*. Physical Review E 70, 056131.

Domain Adaptation of Maximum Entropy Language Models

Tanel Alumäe*

Adaptive Informatics Research Centre
School of Science and Technology
Aalto University
Helsinki, Finland
tanel@cis.hut.fi

Mikko Kurimo

Adaptive Informatics Research Centre
School of Science and Technology
Aalto University
Helsinki, Finland
Mikko.Kurimo@tkk.fi

Abstract

We investigate a recently proposed Bayesian adaptation method for building style-adapted maximum entropy language models for speech recognition, given a large corpus of written language data and a small corpus of speech transcripts. Experiments show that the method consistently outperforms linear interpolation which is typically used in such cases.

1 Introduction

In large vocabulary speech recognition, a language model (LM) is typically estimated from large amounts of written text data. However, recognition is typically applied to speech that is stylistically different from written language. For example, in an often-tried setting, speech recognition is applied to broadcast news, that includes introductory segments, conversations and spontaneous interviews. To decrease the mismatch between training and test data, often a small amount of speech data is human-transcribed. A LM is then built by interpolating the models estimated from large corpus of written language and the small corpus of transcribed data. However, in practice, different models might be of different importance depending on the word context. Global interpolation doesn't take such variability into account and all predictions are weighted across models identically, regardless of the context.

In this paper we investigate a recently proposed Bayesian adaptation approach (Daume III, 2007; Finkel and Manning, 2009) for adapting a conditional maximum entropy (ME) LM (Rosenfeld, 1996) to a new domain, given a large corpus of out-of-domain training data and a small corpus of in-domain data. The main contribution of this

* Currently with Tallinn University of Technology, Estonia

paper is that we show how the suggested hierarchical adaptation can be used with suitable priors and combined with the class-based speedup technique (Goodman, 2001) to adapt ME LMs in large-vocabulary speech recognition when the amount of target data is small. The results outperform the conventional linear interpolation of background and target models in both N -grams and ME models. It seems that with the adapted ME models, the same recognition accuracy for the target evaluation data can be obtained with 50% less adaptation data than in interpolated ME models.

2 Review of Conditional Maximum Entropy Language Models

Maximum entropy (ME) modeling is a framework that has been used in a wide area of natural language processing (NLP) tasks. A conditional ME model has the following form:

$$P(x|h) = \frac{e^{\sum_i \lambda_i f_i(x,h)}}{\sum_{x'} e^{\sum_j \lambda_j f_j(x',h)}} \quad (1)$$

where x is an outcome (in case of a LM, a word), h is a context (the word history), and x' a set of all possible outcomes (words). The functions f_i are (typically binary) feature functions. During ME training, the optimal weights λ_i corresponding to features $f_i(x, h)$ are learned. More precisely, finding the ME model is equal to finding weights that maximize the log-likelihood $L(X; \Lambda)$ of the training data X . The weights are learned via improved iterative scaling algorithm or some of its modern fast counterparts (i.e., conjugate gradient descent).

Since LMs typically have a vocabulary of tens of thousands of words, the use of a normalization factor over all possible outcomes makes estimating a ME LM very memory and time consuming. Goodman (2001) proposed a class-based method that drastically reduces the resource requirements for training such models. The idea is to cluster

words in the vocabulary into classes (e.g., based on their distributional similarity). Then, we can decompose the prediction of a word given its history into prediction of its class given the history, and prediction of the word given the history and its class :

$$P(w|h) = P(C(w)|h) \cdot P(w|h, C(w)) \quad (2)$$

Using such decomposition, we can create two ME models: one corresponding to $P(C(w)|h)$ and the other corresponding to $P(w|h, C(w))$. It is easy to see that computing the normalization factor of the first component model now requires only looping over all classes. It turns out that normalizing the second model is also easier: for a context $h, C(w)$, we only need to normalize over words that belong to class $C(w)$, since other words cannot occur in this context. This decomposition can be further extended by using hierarchical classes.

To avoid overfitting, ME models are usually smoothed (regularized). The most widely used smoothing method for ME LMs is Gaussian priors (Chen and Rosenfeld, 2000): a zero-mean prior with a given variance is added to all feature weights, and the model optimization criteria becomes:

$$L'(X; \Lambda) = L(X; \Lambda) - \sum_{i=1}^F \frac{\lambda_i^2}{2\sigma_i^2} \quad (3)$$

where F is the number of feature functions. Typically, a fixed hyperparameter $\sigma_i = \sigma$ is used for all parameters. The optimal variance is usually estimated on a development set. Intuitively, this method encourages feature weights to be smaller, by penalizing weights with big absolute values.

3 Domain Adaptation of Maximum Entropy Models

Recently, a hierarchical Bayesian adaptation method was proposed that can be applied to a large family of discriminative learning tasks (such as ME models, SVMs) (Daume III, 2007; Finkel and Manning, 2009). In NLP problems, data often comes from different sources (e.g., newspapers, web, textbooks, speech transcriptions). There are three classic approaches for building models from multiple sources. We can pool all training data and estimate a single model, and apply it for all tasks. The second approach is to “unpool” the data, i.e. only use training data from the test domain. The

third and often the best performing approach is to train separate models for each data source, apply them to test data and interpolate the results.

The hierarchical Bayesian adaptation method is a generalization of the three approaches described above. The hierarchical model jointly optimizes global and domain-specific parameters, using parameters built from pooled data as priors for domain-specific parameters. In other words, instead of using smoothing to encourage parameters to be closer to zero, it encourages domain-specific model parameters to be closer to the corresponding global parameters, while a zero mean Gaussian prior is still applied for global parameters. For processing test data during runtime, the domain-specific model is applied. Intuitively, this approach can be described as follows: the domain-specific parameters are largely determined by global data, unless there is good domain-specific evidence that they should be different. The key to this approach is that the global and domain-specific parameters are learned jointly, not hierarchically. This allows domain-specific parameters to influence the global parameters, and vice versa. Formally, the joint optimization criteria becomes:

$$L_{hier}(X; \Lambda) = \sum_d \left(L_{orig}(X_d, \Lambda_d) - \sum_{i=1}^F \frac{(\lambda_{d,i} - \lambda_{*,i})^2}{2\sigma_d^2} \right) - \sum_{i=1}^F \frac{\lambda_{*,i}^2}{2\sigma_*^2} \quad (4)$$

where X_d is data for domain d , $\lambda_{*,i}$ the global parameters, $\lambda_{d,i}$ the domain-specific parameters, σ_*^2 the global variance and σ_d^2 the domain-specific variances. The global and domain-specific variances are optimized on the heldout data. Usually, larger values are used for global parameters and for domains with more data, while for domains with less data, the variance is typically set to be smaller, encouraging the domain-specific parameters to be closer to global values.

This adaptation scheme is very similar to the approaches proposed by (Chelba and Acero, 2006) and (Chen, 2009b): both use a model estimated from background data as a prior when learning a model from in-domain data. The main difference is the fact that in this method, the models are estimated jointly while in the other works, back-

ground model has to be estimated before learning the in-domain model.

4 Experiments

In this section, we look at experimental results over two speech recognition tasks.

4.1 Tasks

Task 1: English Broadcast News. This recognition task consists of the English broadcast news section of the 2003 NIST Rich Transcription Evaluation Data. The data includes six news recordings from six different sources with a total length of 176 minutes.

As acoustic models, the CMU Sphinx open source triphone HUB4 models for wideband (16kHz) speech¹ were used. The models have been trained using 140 hours of speech.

For training the LMs, two sources were used: first 5M sentences from the Gigaword (2nd ed.) corpus (99.5M words), and broadcast news transcriptions from the TDT4 corpus (1.19M words). The latter was treated as in-domain data in the adaptation experiments. A vocabulary of 26K words was used. It is a subset of a bigger 60K vocabulary, and only includes words that occurred in the training data. The OOV rate against the test set was 2.4%.

The audio used for testing was segmented into parts of up to 20 seconds in length. Speaker diarization was applied using the LIUM.SpKDiArization toolkit (Deléglise et al., 2005). The CMU Sphinx 3.7 was used for decoding. A three-pass recognition strategy was applied: the first pass recognition hypotheses were used for calculating MLLR-adapted models for each speaker. In the second pass, the adapted acoustic models were used for generating a 5000-best list of hypotheses for each segment. In the third pass, the ME LM was used to re-rank the hypotheses and select the best one. During decoding, a trigram LM model was used. The trigram model was an interpolation of source-specific models which were estimated using Kneser-Ney discounting.

Task 2: Estonian Broadcast Conversations. The second recognition task consists of four recordings from different live talk programs from

three Estonian radio stations. Their format consists of hosts and invited guests, spontaneously discussing current affairs. There are 40 minutes of transcriptions, with 11 different speakers.

The acoustic models were trained on various wideband Estonian speech corpora: the BABEL speech database (9h), transcriptions of Estonian broadcast news (7.5h) and transcriptions of radio live talk programs (10h). The models are triphone HMMs, using MFCC features.

For training the LMs, two sources were used: about 10M sentences from various Estonian newspapers, and manual transcriptions of 10 hours of live talk programs from three Estonian radio stations. The latter is identical in style to the test data, although it originates from a different time period and covers a wider variety of programs, and was treated as in-domain data.

As Estonian is a highly inflective language, morphemes are used as basic units in the LM. We use a morphological analyzer (Kaalep and Vaino, 2001) for splitting the words into morphemes. After such processing, the newspaper corpus includes of 185M tokens, and the transcribed data 104K tokens. A vocabulary of 30K tokens was used for this task, with an OOV rate of 1.7% against the test data. After recognition, morphemes were concatenated back to words.

As with English data, a three-pass recognition strategy involving MLLR adaptation was applied.

4.2 Results

For both tasks, we rescored the N-best lists in two different ways: (1) using linear interpolation of source-specific ME models and (2) using hierarchically domain-adapted ME model (as described in previous chapter). The English ME models had a three-level and Estonian models a four-level class hierarchy. The classes were derived using the word exchange algorithm (Kneser and Ney, 1993). The number of classes at each level was determined experimentally so as to optimize the resource requirements for training ME models (specifically, the number of classes was 150, 1000 and 5000 for the English models and 20, 150, 1000 and 6000 for the Estonian models). We used unigram, bigram and trigram features that occurred at least twice in the training data. The feature cut-off was applied in order to accommodate the memory requirements. The feature set was identical for interpolated and adapted models.

¹<http://www.speech.cs.cmu.edu/sphinx/models/>

Adaptation data (No of words)	Interp. models		Adapted models		
	σ_{OD}^2	σ_{ID}^2	σ_*^2	σ_{OD}^2	σ_{ID}^2
English Broadcast News					
147K	2e8	3e5	5e7	2e7	2e6
292K	2e8	5e5	5e7	2e7	2e6
591K	2e8	1e6	5e7	2e7	2e6
1119K	2e8	2e6	5e7	2e7	5e6
Estonian Broadcast Conversations					
104K	5e8	3e5	5e7	1e7	2e6

Table 1: The unnormalized values of Gaussian prior variances for interpolated out-of-domain (OD) and in-domain (ID) ME models, and hierarchically adapted global (*), out-of-odomain (OD) and in-domain (ID) models that were used in the experiments.

For the English task, we also explored the efficiency of these two approaches with varying size of adaptation data: we repeated the experiments when using one eighth, one quarter, half and all of the TDT4 transcription data for interpolation/adaptation. The amount of used Gigaword data was not changed. In all cases, interpolation weights were re-optimized and new Gaussian variance values were heuristically determined.

The TADM toolkit² was used for estimating ME models, utilizing its implementation of the conjugate gradient algorithm.

The models were regularized using Gaussian priors. The variance parameters were chosen heuristically based on light tuning on development set perplexity. For the source-specific ME models, the variance was fixed on per-model basis. For the adapted model, that jointly models global and domain-specific data, the Gaussian priors were fixed for each hierarchy node (i.e., the variance was fixed across global, out-of-domain, and in-domain parameters). Table 1 lists values for the variances of Gaussian priors (as in equations 3 and 4) that we used in the experiments. In other publications, the variance values are often normalized to the size of the data. We chose not to normalize the values, since in the hierarchical adaptation scheme, also data from other domains have impact on the learned model parameters, thus

²<http://tadm.sourceforge.net/>

it’s not possible to simply normalize the variances.

The experimental results are presented in Table 2. Perplexity and word error rate (WER) results of the interpolated and adapted models are compared. For the Estonian task, letter error rate (LER) is also reported, since it tends to be a more indicative measure of speech recognition quality for highly inflected languages. In all experiments, using the adapted models resulted in lower perplexity and lower error rate. Improvements in the English experiment were less evident than in the Estonian system, with under 10% improvement in perplexity and 1-3% in WER, against 15% and 4% for the Estonian experiment. In most cases, there was a significant improvement in WER when using the adapted ME model (according to the Wilcoxon test), with and exception of the English experiments on the 292K and 591K data sets.

The comparison between N -gram models and ME models is not entirely fair since ME models are actually class-based. Such transformation introduces additional smoothing into the model and can improve model perplexity, as also noticed by Goodman (2001).

5 Discussion

In this paper we have tested a hierarchical adaptation method (Daume III, 2007; Finkel and Manning, 2009) on building style-adapted LMs for speech recognition. We showed that the method achieves consistently lower error rates than when using linear interpolation which is typically used in such scenarios.

The tested method is ideally suited for language modeling in speech recognition: we almost always have access to large amounts of data from written sources but commonly the speech to be recognized is stylistically noticeably different. The hierarchical adaptation method enables to use even a small amount of in-domain data to modify the parameters estimated from out-of-domain data, if there is enough evidence.

As Finkel and Manning (2009) point out, the hierarchical nature of the method makes it possible to estimate highly specific models: we could draw style-specific models from general high-level priors, and topic-and-style specific models from style-specific priors. Furthermore, the models don’t have to be hierarchical: it is easy to generalize the method to general multilevel approach where a model is drawn from multiple priors. For

Adaptation data (No. of words)	Perplexity				WER			LER		
	Pooled N-gram	Interp. N-gram	Interp. ME	Adapted ME	Interp. N-gram	Interp. ME	Adapted ME	Interp. N-gram	Interp. ME	Adapted ME
English Broadcast News										
147K	290	255	243	230	27.2	26.3	25.9			
292K	286	250	236	223	26.7	25.8	25.6			
591K	280	243	228	215	26.6	25.9	25.6			
1119K	272	232	217	204	26.2	25.6	24.9			
Estonian Broadcast Conversations										
104K	237	197	200	169	40.5	38.9	37.4	17.7	17.3	16.6

Table 2: Perplexity, WER and LER results comparing pooled and interpolated N -gram models and interpolated and adapted ME models, with changing amount of available in-domain data.

instance, we could build a model for recognizing computer science lectures, given data from textbooks, including those about computer science, and transcripts of lectures on various topics (which don't even need to include lectures about computer science).

The method has some considerable shortcomings from the practical perspective. First, training ME LMs in general has much higher resource requirements than training N -gram models which are typically used in speech recognition. Moreover, training hierarchical ME models requires even more memory than training simple ME models, proportional to the number of nodes in the hierarchy. However, it should be possible to alleviate this problem by profiting from the hierarchical nature of n -gram features, as proposed in (Wu and Khudanpur, 2002). It is also difficult to determine good variance values σ_i^2 for the global and domain-specific priors. While good variance values for simple ME models can be chosen quite reliably based on the size of the training data (Chen, 2009a), we have found that it is more demanding to find good hyperparameters for hierarchical models since weights for the same feature in different nodes in the hierarchy are all related to each other. We plan to investigate this problem in the future since the choice of hyperparameters has a strong impact on the performance of the model.

Acknowledgments

This research was partly funded by the Academy of Finland in the project Adaptive Informatics, by the target-financed theme No. 0322709s06 of the Estonian Ministry of Education and Research and by the National Programme for Estonian Lan-

guage Technology.

References

- Ciprian Chelba and Alex Acero. 2006. Adaptation of maximum entropy capitalizer: Little data can help a lot. *Computer Speech & Language*, 20(4):382–399, October.
- S. F. Chen and R. Rosenfeld. 2000. A survey of smoothing techniques for ME models. *IEEE Transactions on Speech and Audio Processing*, 8(1):37–50.
- S. F. Chen. 2009a. Performance prediction for exponential language models. In *Proceedings of HLT-NAACL*, pages 450–458, Boulder, Colorado.
- Stanley F. Chen. 2009b. Shrinking exponential language models. In *Proceedings of HLT-NAACL*, pages 468–476, Boulder, Colorado.
- H. Daume III. 2007. Frustratingly easy domain adaptation. In *Proceedings of ACL*, pages 256–263.
- P. Deléglise, Y. Estève, S. Meignier, and T. Merlin. 2005. The LIUM speech transcription system: a CMU Sphinx III-based system for French broadcast news. In *Proceedings of Interspeech*, Lisboa, Portugal.
- J. R. Finkel and Ch. Manning. 2009. Hierarchical Bayesian domain adaptation. In *Proceedings of HLT-NAACL*, pages 602–610, Boulder, Colorado.
- J. Goodman. 2001. Classes for fast maximum entropy training. In *Proceedings of ICASSP*, Utah, USA.
- H.-J. Kaalep and T. Vaino. 2001. Complete morphological analysis in the linguist's toolbox. In *Congressus Nonus Internationalis Fenno-Ugristarum Pars V*, pages 9–16, Tartu, Estonia.
- R. Kneser and H. Ney. 1993. Improved clustering techniques for class-based statistical language modelling. In *Proceedings of the European Conference*

on *Speech Communication and Technology*, pages 973–976.

R. Rosenfeld. 1996. A maximum entropy approach to adaptive statistical language modeling. *Computer, Speech and Language*, 10:187–228.

J. Wu and S. Khudanpur. 2002. Building a topic-dependent maximum entropy model for very large corpora. In *Proceedings of ICASSP*, Orlando, Florida, USA.

Decision detection using hierarchical graphical models

Trung H. Bui

CSLI

Stanford University
Stanford, CA 94305, USA
thbui@stanford.edu

Stanley Peters

CSLI

Stanford University
Stanford, CA 94305, USA
peters@csli.stanford.edu

Abstract

We investigate hierarchical graphical models (HGMs) for automatically detecting decisions in multi-party discussions. Several types of dialogue act (DA) are distinguished on the basis of their roles in formulating decisions. HGMs enable us to model dependencies between observed features of discussions, decision DAs, and subdialogues that result in a decision. For the task of detecting decision regions, an HGM classifier was found to outperform non-hierarchical graphical models and support vector machines, raising the F1-score to 0.80 from 0.55.

1 Introduction

In work environments, people share information and make decisions in multi-party conversations known as meetings. The demand for systems that can automatically process information contained in audio and video recordings of meetings is growing rapidly. Our own research, and that of other contemporary projects (Janin et al., 2004) aim at meeting this demand.

We are currently investigating the automatic detection of decision discussions. Our approach involves distinguishing between different dialogue act (DA) types based on their role in the decision-making process. These DA types are called Decision Dialogue Acts (DDAs). Groups of DDAs combine to form a decision region.

Recent work (Bui et al., 2009) showed that Directed Graphical Models (DGMs) outperform other machine learning techniques such as Support Vector Machines (SVMs) for detecting individual DDAs. However, the proposed models, which were non-hierarchical, did not significantly improve identification of decision regions. This paper tests whether giving DGMs hierarchical structure (making them HGMs) can improve

their performance at this task compared with non-hierarchical DGMs.

We proceed as follows. Section 2 discusses related work, and section 3 our data set and annotation scheme for decision discussions. Section 4 summarizes previous decision detection experiments using DGMs. Section 5 presents the HGM approach, and section 6 describes our HGM experiments. Finally, section 7 draws conclusions and presents ideas for future work.

2 Related work

User studies (Banerjee et al., 2005) have confirmed that meeting participants consider decisions to be one of the most important meeting outputs, and Whittaker et al. (2006) found that the development of an automatic decision detection component is critical for re-using meeting archives. With the new availability of substantial meeting corpora such as the AMI corpus (McCowan et al., 2005), recent years have seen an increasing amount of research on decision-making dialogue. This research has tackled issues such as the automatic detection of agreement and disagreement (Galley et al., 2004), and of the level of involvement of conversational participants (Gatica-Perez et al., 2005). Recent work on automatic detection of decisions has been conducted by Hsueh and Moore (2007), Fernández et al. (2008), and Bui et al. (2009).

Fernández et al. (2008) proposed an approach to modeling the structure of decision-making dialogue. These authors designed an annotation scheme that takes account of the different roles that utterances can play in the decision-making process—for example it distinguishes between DDAs that initiate a decision discussion by raising an issue, those that propose a resolution of the issue, and those that express agreement to a proposed resolution. The authors annotated a portion of the AMI corpus, and then applied what

they refer to as “hierarchical classification.” Here, one *sub-classifier* per DDA class hypothesizes occurrences of that type of DDA and then, based on these hypotheses, a *super-classifier* determines which regions of dialogue are decision discussions. All of the classifiers, (sub and super), were linear kernel binary SVMs. Results were better than those obtained with (Hsueh and Moore, 2007)’s approach—the F1-score for detecting decision discussions in manual transcripts was 0.58 vs. 0.50. Purver et al. (2007) had earlier detected action items with the approach Fernández et al. (2008) extended to decisions.

Bui et al. (2009) built on the promising results of (Fernández et al., 2008), by employing DGMs in place of SVMs. DGMs are attractive because they provide a natural framework for modeling sequence and dependencies between variables, including the DDAs. Bui et al. (2009) were especially interested in whether DGMs better exploit non-lexical features. Fernández et al. (2008) obtained much more value from lexical than non-lexical features (and indeed no value at all from prosodic features), but lexical features have limitations. In particular, they can be domain specific, increase the size of the feature space dramatically, and deteriorate more in quality than other features when automatic speech recognition (ASR) is poor. More detail about decision detection using DGMs will be presented in section 4.

Beyond decision detection, DGMs are used for labeling and segmenting sequences of observations in many different fields—including bioinformatics, ASR, Natural Language Processing (NLP), and information extraction. In particular, Dynamic Bayesian Networks (DBNs) are a popular model for probabilistic sequence modeling because they exploit structure in the problem to compactly represent distributions over multi-state and observation variables. Hidden Markov Models (HMMs), a special case of DBNs, are a classical method for important NLP applications such as unsupervised part-of-speech tagging (Gael et al., 2009) and grammar induction (Johnson et al., 2007) as well as for ASR. More complex DBNs have been used for applications such as DA recognition (Crook et al., 2009) and activity recognition (Bui et al., 2002).

Undirected graphical models (UGMs) are also valuable for building probabilistic models for segmenting and labeling sequence data. Conditional

Random Fields (CRFs), a simple UGM case, can avoid the label bias problem (Lafferty et al., 2001) and outperform maximum entropy Markov models and HMMs.

However, the graphical models used in these applications are mainly non-hierarchical, including those in Bui et al. (2009). Only Sutton et al. (2007) proposed a three-level HGM (in the form of a dynamic CRF) for the joint noun phrase chunking and part of speech labeling problem; they showed that this model performs better than a non-hierarchical counterpart.

3 Data

For the experiments reported in this study, we used 17 meetings from the AMI Meeting Corpus¹, a freely available corpus of multi-party meetings with both audio and video recordings, and a wide range of annotated information including DAs and topic segmentation. The meetings last around 30 minutes each, and are scenario-driven, wherein four participants play different roles in a company’s design team: *project manager*, *marketing expert*, *interface designer* and *industrial designer*.

We use the same annotation scheme as Fernández et al. (2008) to model decision-making dialogue. As stated in section 2, this scheme distinguishes between a small number of DA types based on the role which they perform in the formulation of a decision. Besides improving the detection of decision discussions (Fernández et al., 2008), such a scheme also aids in summarization of them, because it indicates which utterances provide particular types of information.

The annotation scheme is based on the observation that a decision discussion typically contains the following main structural components: (a) A topic or issue requiring resolution is raised; (b) One or more possible resolutions are considered; (c) A particular resolution is agreed upon, and so adopted as the decision. Hence the scheme distinguishes between three main DDA classes: *issue* (*I*), *resolution* (*R*), and *agreement* (*A*). Class *R* is further subdivided into *resolution proposal* (*RP*) and *resolution restatement* (*RR*). *I* utterances introduce the topic of the decision discussion, examples being “*Are we going to have a backup?*” and “*But would a backup really be necessary?*” in Table 1. In comparison, *R* utterances specify the resolution which is ultimately adopted as the deci-

¹<http://corpus.amiproject.org/>

- (1) A: Are we going to have a backup? Or we do just–
 B: But would a backup really be necessary?
 A: I think maybe we could just go for the kinetic energy and be bold and innovative.
 C: Yeah.
 B: I think– yeah.
 A: It could even be one of our selling points.
 C: Yeah –*laugh*–.
 D: Environmentally conscious or something.
 A: Yeah.
 B: Okay, fully kinetic energy.
 D: Good.

Table 1: An excerpt from the AMI dialogue ES2015c. It has been modified slightly for presentation purposes.

sion. *RP* utterances propose this resolution (e.g. “*I think maybe we could just go for the kinetic energy ...*”), while *RR* utterances close the discussion by confirming/summarizing the decision (e.g. “*Okay, fully kinetic energy*”). Finally, *A* utterances agree with the proposed resolution, signaling that it is adopted as the decision, (e.g. “*Yeah*”, “*Good*” and “*Okay*”). Unsurprisingly, an utterance may be assigned to more than one DDA class; and within a decision discussion, more than one utterance can be assigned to the same DDA class.

We use manual transcripts in the experiments described here. Inter-annotator agreement was satisfactory, with kappa values ranging from .63 to .73 for the four DDA classes. The manual transcripts contain a total of 15,680 utterances, and on average 40 DDAs per meeting. DDAs are sparse in the transcripts: for all DDAs, 6.7% of the totality of utterances; for *I*, 1.6%; for *RP*, 2%; for *RR*, 0.5%; and for *A*, 2.6%. In all, 3753 utterances (i.e., 23.9%) are tagged as decision-related utterances, and on average there are 221 decision-related utterances per meeting.

4 Prior Work on Decision Detection using Graphical Models

To detect each individual DDA class, Bui et al. (2009) examined the four simple DGMs shown in Fig. 1. The DDA node is binary valued, with value 1 indicating the presence of a DDA and 0 its absence. The evidence node (*E*) is a multi-dimensional vector of observed values of non-lexical features. These include utterance features

(UTT) such as length in words², duration in milliseconds, position within the meeting (as percentage of elapsed time), manually annotated dialogue act (DA) features³ such as *inform*, *assess*, *suggest*, and prosodic features (PROS) such as energy and pitch. These features are the same as the non-lexical features used by Fernández et al. (2008). The hidden component node (*C*) in the *-mix* models represents the distribution of observable evidence *E* as a mixture of Gaussian distributions. The number of Gaussian components was hand-tuned during the training phase.

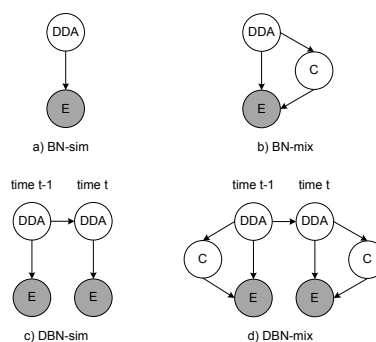


Figure 1: Simple DGMs for individual decision dialogue act detection. The clear nodes are hidden, and the shaded nodes are observable.

More complex models were constructed from the four simple models in Fig. 1 to allow for dependencies between different DDAs. For example, the model in Fig. 2 generalizes Fig. 1c with arcs connecting the DDA classes based on analysis of the annotated AMI data.

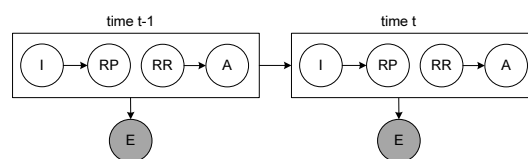


Figure 2: A DGM that takes the dependencies between decision dialogue acts into account.

Decision discussion regions were identified using the DGM output and the following two simple rules: (1) A decision discussion region begins with an *Issue* DDA; (2) A decision discussion region contains at least one *Issue* DDA and one *Resolution* DDA.

²This feature is a manual count of lexical tokens; but word count was extracted automatically from ASR output by Bui et al. (2009). We plan experiments to determine how much using ASR output degrades detection of decision regions.

³The authors used the AMI DA annotations.

The authors conducted experiments using the AMI corpus and found that when using non-lexical features, the DGMs outperform the hierarchical SVM classification method of (Fernández et al., 2008). The F1-score for the four DDA classes increased between 0.04 and 0.19 ($p < 0.005$), and for identifying decision discussion regions, by 0.05 ($p > 0.05$).

5 Hierarchical graphical models

Although the results just discussed showed graphical models are better than SVMs for detecting decision dialogue acts (Bui et al., 2009), two-level graphical models like those shown in Figs. 1 and 2 cannot exploit dependencies between high-level discourse items such as decision discussions and DDAs; and the “superclassifier” rule (Bui et al., 2009) used for detecting decision regions did not significantly improve the F1-score for decisions.

We thus investigate whether HGMs (structured as three or more levels) are superior for discovering the structure and learning the parameters of decision recognition. Our approach composes graphical models to increase hierarchy with an additional level above or below previous ones, or inserts a new level such as for discourse topics into the interior of a given model.

Fig. 3 shows a simple structure for three-level HGMs. The top level corresponds to high-level discourse regions such as decision discussions. The segmentation into these regions is represented in terms of a random variable (at each DR node) that takes on discrete values: {positive, negative} (the utterance belongs to a decision region or not) or {begin, middle, end, outside} (indicating the position of the utterance relative to a decision discussion region). The middle level corresponds to mid-level discourse items such as issues, resolution proposals, resolution restatements, and agreements. These classes (C_1, C_2, \dots, C_n nodes) are represented as a collection of random variables, each corresponding to an individual mid-level utterance class. For example, the middle level of the three-level HGM Fig. 3 could be the top-level of the two-level DGM in Fig. 2, each middle level node containing random variables for the DDA classes I, RP, RR, and A. The bottom level corresponds to vectors of observed features as before, e.g. lexical, utterance, and prosodic features.

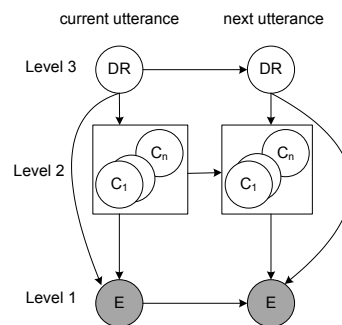


Figure 3: A simple structure of a three-level HGM: DRs are high-level discourse regions; C_1, C_2, \dots, C_n are mid-level utterance classes; and Es are vectors of observed features.

6 Experiments

The HGM classifier in Figure 3 was implemented in Matlab using the BNT software⁴. The classifier hypothesizes that an utterance belongs to a decision region if the marginal probability of the utterance’s DR node is above a hand-tuned threshold. The threshold is selected using the ROC curve analysis⁵ to obtain the highest F1-score. To evaluate the accuracy of hypothesized decision regions, we divided the dialogue into 30-second windows and evaluated on a per window basis.

The best model structure was selected by comparing the performance of various handcrafted structures. For example, the model in Fig. 4b outperforms the one in Fig. 4a. Fig. 4b explicitly models the dependency between the decision regions and the observed features.

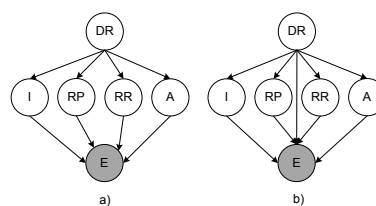


Figure 4: Three-level HGMs for recognition of decisions. This illustrates the choice of the structure for each time slice of the HGM sequence models.

Table 2 shows the results of 17-fold cross-validation for the hierarchical SVM classification (Fernández et al., 2008), rule-based classification with DGM output (Bui et al., 2009), and our HGM classification using the best combination of non-lexical features. All three methods

⁴<http://www.cs.ubc.ca/~murphyk/Software/BNT/bnt.html>

⁵http://en.wikipedia.org/wiki/Receiver_operating_characteristic

were implemented by us using exactly the same data and 17-fold cross-validation. The features were selected based on the best combination of non-lexical features for each method. The HGM classifier outperforms both its SVM and DGM counterparts ($p < 0.0001$)⁶. In fact, even when the SVM uses lexical as well as non-lexical features, its F1-score is still lower than the HGM classifier.

Classifier	Pr	Re	F1
SVM	0.35	0.88	0.50
DGM	0.39	0.93	0.55
HGM	0.69	0.96	0.80

Table 2: Results for detection of decision discussion regions by the SVM super-classifier, rule-based DGM classifier, and HGM classifier, each using its best combination of non-lexical features: SVM (UTT+DA), DGM (UTT+DA+PROS), HGM (UTT+DA).

In contrast with the hierarchical SVM and rule-based DGM methods, the HGM method identifies decision-related utterances by exploiting not just DDAs but also direct dependencies between decision regions and UTT, DA, and PROS features. As mentioned in the second paragraph of this section, explicitly modeling the dependency between decision regions and observable features helps to improve detection of decision regions. Furthermore, a three-level HGM can straightforwardly model the composition of each high-level decision region as a sequence of mid-level DDA utterances. While the hierarchical SVM method can also take dependency between successive utterances into account, it has no principled way to associate this dependency with more extended decision regions. In addition, this dependency is only meaningful for lexical features (Fernández et al., 2008).

The HGM result presented in Table 2 was computed using the three-level DBN model (see Fig. 4b) using the combination of UTT and DA features. Without DA features, the F1-score degrades from 0.8 to 0.78. However, this difference is not statistically significant (i.e., $p > 0.5$).

7 Conclusions and Future Work

To detect decision discussions in multi-party dialogue, we investigated HGMs as an extension of

⁶We used the paired t test for computing statistical significance. <http://www.graphpad.com/quickcalcs/ttest1.cfm>

the DGMs studied in (Bui et al., 2009). When using non-lexical features, HGMs outperform the non-hierarchical DGMs of (Bui et al., 2009) and also the hierarchical SVM classification method of Fernández et al. (2008). The F1-score for identifying decision discussion regions increased to 0.80 from 0.55 and 0.50 respectively ($p < 0.0001$).

In future work, we plan to (a) investigate cascaded learning methods (Sutton et al., 2007) to improve the detection of DDAs further by using detected decision regions and (b) extend HGMs beyond three levels in order to integrate useful semantic information such as topic structure.

Acknowledgments

The research reported in this paper was sponsored by the Department of the Navy, Office of Naval Research, under grants number N00014-09-1-0106 and N00014-09-1-0122. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the Office of Naval Research.

References

- Satanjeev Banerjee, Carolyn Rosé, and Alex Rudnicky. 2005. The necessity of a meeting recording and playback system, and the benefit of topic-level annotations to meeting browsing. In *Proceedings of the 10th International Conference on Human-Computer Interaction*.
- H. H. Bui, S. Venkatesh, and G. West. 2002. Policy recognition in the abstract hidden markov model. *Journal of Artificial Intelligence Research*, 17:451–499.
- Trung Huu Bui, Matthew Frampton, John Dowding, and Stanley Peters. 2009. Extracting decisions from multi-party dialogue using directed graphical models and semantic similarity. In *Proceedings of the 10th Annual SIGDIAL Meeting on Discourse and Dialogue (SIGdial09)*.
- Nigel Crook, Ramon Granel, and Stephen Pulman. 2009. Unsupervised classification of dialogue acts using a dirichlet process mixture model. In *Proceedings of SIGDIAL 2009: the 10th Annual Meeting of the Special Interest Group in Discourse and Dialogue*, pages 341–348.
- Raquel Fernández, Matthew Frampton, Patrick Ehlen, Matthew Purver, and Stanley Peters. 2008. Modelling and detecting decisions in multi-party dialogue. In *Proceedings of the 9th SIGdial Workshop on Discourse and Dialogue*.

- Jurgen Van Gael, Andreas Vlachos, and Zoubin Ghahramani. 2009. The infinite HMM for unsupervised PoS tagging. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing*, pages 678–687.
- Michel Galley, Kathleen McKeown, Julia Hirschberg, and Elizabeth Shriberg. 2004. Identifying agreement and disagreement in conversational speech: Use of Bayesian networks to model pragmatic dependencies. In *Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Daniel Gatica-Perez, Ian McCowan, Dong Zhang, and Samy Bengio. 2005. Detecting group interest level in meetings. In *Proceedings of ICASSP*.
- Pey-Yun Hsueh and Johanna Moore. 2007. Automatic decision detection in meeting speech. In *Proceedings of MLMI 2007*, Lecture Notes in Computer Science. Springer-Verlag.
- Adam Janin, Jeremy Ang, Sonali Bhagat, Rajdip Dhillon, Jane Edwards, Javier Marciás-Guarasa, Nelson Morgan, Barbara Peskin, Elizabeth Shriberg, Andreas Stolcke, Chuck Wooters, and Britta Wrede. 2004. The ICSI meeting project: Resources and research. In *Proceedings of the 2004 ICASSP NIST Meeting Recognition Workshop*.
- Mark Johnson, Thomas Griffiths, and Sharon Goldwater. 2007. Bayesian inference for PCFGs via Markov chain Monte Carlo. In *Proceedings of Human Language Technologies 2007: The Conference of the North American Chapter of the Association for Computational Linguistics*, pages 139–146, Rochester, New York, April. Association for Computational Linguistics.
- John Lafferty, Andrew McCallum, and Fernando Pereira. 2001. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of the 18th International Conference on Machine Learning*, pages 282–289. Morgan Kaufmann.
- Iain McCowan, Jean Carletta, W. Kraaij, S. Ashby, S. Bourban, M. Flynn, M. Guillemot, T. Hain, J. Kadlec, V. Karaiskos, M. Kronenthal, G. Lathoud, M. Lincoln, A. Lisowska, W. Post, D. Reidsma, and P. Wellner. 2005. The AMI Meeting Corpus. In *Proceedings of Measuring Behavior, the 5th International Conference on Methods and Techniques in Behavioral Research*, Wageningen, Netherlands.
- Matthew Purver, John Dowding, John Niekrasz, Patrick Ehlen, Sharareh Noorbaloochi, and Stanley Peters. 2007. Detecting and summarizing action items in multi-party dialogue. In *Proceedings of the 8th SIGdial Workshop on Discourse and Dialogue*, Antwerp, Belgium.
- Charles Sutton, Andrew McCallum, and Khashayar Rohanimanesh. 2007. Dynamic conditional random fields: Factorized probabilistic models for labeling and segmenting sequence data. *Journal of Machine Learning Research*, 8:693–723.
- Steve Whittaker, Rachel Laban, and Simon Tucker. 2006. Analysing meeting records: An ethnographic study and technological implications. In S. Renals and S. Bengio, editors, *Machine Learning for Multimodal Interaction: Second International Workshop, MLMI 2005, Revised Selected Papers*, volume 3869 of *Lecture Notes in Computer Science*, pages 101–113. Springer.

Using Speech to Reply to SMS Messages While Driving: An In-Car Simulator User Study

Yun-Cheng Ju, Tim Paek

Microsoft Research

Redmond, WA USA

{yuncj|timpaek}@microsoft.com

Abstract

Speech recognition affords automobile drivers a hands-free, eyes-free method of replying to Short Message Service (SMS) text messages. Although a voice search approach based on template matching has been shown to be more robust to the challenging acoustic environment of automobiles than using dictation, users may have difficulties verifying whether SMS response templates match their intended meaning, especially while driving. Using a high-fidelity driving simulator, we compared dictation for SMS replies versus voice search in increasingly difficult driving conditions. Although the two approaches did not differ in terms of driving performance measures, users made about six times more errors on average using dictation than voice search.

1 Introduction

Users love Short Message Service (SMS) text messaging; so much so that 3 trillion SMS messages are expected to have been sent in 2009 alone (Stross, 2008). Because research has shown that SMS messaging while driving results in 35% slower reaction time than being intoxicated (Reed & Robbins, 2008), campaigns have been launched by states, governments and even cell phone carriers to discourage and ban SMS messaging while driving (DOT, 2009). Yet, automobile manufacturers have started to offer infotainment systems, such as the Ford Sync, which feature the ability to listen to incoming SMS messages using text-to-speech (TTS). Automatic speech recognition (ASR) affords users a hands-free, eyes-free method of replying to SMS messages. However, to date, manufacturers have not established a safe and reliable method of leveraging ASR, though some researchers have

begun to explore techniques. In previous research (Ju & Paek, 2009), we examined three ASR approaches to replying to SMS messages: dictation using a language model trained on SMS responses, canned responses using a probabilistic context-free grammar (PCFG), and a “voice search” approach based on template matching. Voice search proceeds in two steps (Natarajan et al., 2002): an utterance is first converted into text, which is then used as a search query to match the most similar items of an index using IR techniques (Yu et al., 2007). For SMS replies, we created an index of SMS response templates, with slots for semantic concepts such as time and place, from a large SMS corpus. After convolving recorded SMS replies so that the audio would exhibit the acoustic characteristics of in-car recognition, they compared how the three approaches handled the convolved audio with respect to the top n-best reply candidates. The voice search approach consistently outperformed dictation and canned responses, achieving as high as 89.7% task completion with respect to the top 5 reply candidates.

Even if the voice search approach may be more robust to in-car noise, this does not guarantee that it will be more usable. Indeed, because voice search can only match semantic concepts contained in the templates (which may or may not utilize the same wording as the reply), users must verify that a retrieved template matches the semantics of their intended reply. For example, suppose a user replies to the SMS message “*how about lunch*” with “*can’t right now running errands*”. Voice search may find “*nope, got errands to run*” as the closest template match, in which case, users will have to decide whether this response has the same meaning as their reply. This of course entails cognitive effort, which is very limited in the context of driving. On the other hand, a dictation approach to replying to SMS messages may be far worse due to misrecognitions. For example, dictation may interpret “*can’t right now running errands*” as “*can right*

now fun in errands”. We posited that voice search has the advantage because it always generates intelligible SMS replies (since response templates are manually filtered), as opposed to dictation, which can sometimes result in unpredictable and nonsensical misrecognitions. However, this advantage has not been empirically demonstrated in a user study. This paper presents a user study investigating how the two approaches compare when users are actually driving – that is, when usability matters most.

2 Driving Simulator Study

Although ASR affords users hands-free, eyes-free interaction, the benefits of leveraging speech can be forfeit if users are expending cognitive effort judging whether the speech interface correctly interpreted their utterances. Indeed, research has shown that the cognitive demands of dialogue seem to play a more important role in distracting drivers than physically handling cell phones (Nunes & Recarte, 2002; Strayer & Johnston, 2001). Furthermore, Kun et al. (2007) have found that when in-car speech interfaces encounter recognition problems, users tend to drive more dangerously as they attempt to figure out why their utterances are failing. Hence, any approach to replying to SMS messages in automobiles must avoid distracting drivers with errors and be highly usable while users are engaged in their primary task, driving.

2.1 Method

To assess the usability and performance of both the voice search approach and dictation, we conducted a controlled experiment using the STISIM Drive™ simulator. Our simulation setup consisted of a central console with a steering wheel and two turn signals, surrounded by three 47” flat panels placed at a 45° angle to immerse the driver. Figure 1 displays the setup.

We recruited 16 participants (9 males, 7 females) through an email sent to employees of our organization. The mean age was 38.8. All participants had a driver’s license and were compensated for their time.

We examined two independent variables: *SMS Reply Approach*, consisting of *voice search* and *dictation*, and *Driving Condition*, consisting of *no driving*, *easy driving* and *difficult driving*. We included *Driving Condition* as a way of increasing cognitive demand (see next section). Overall, we conducted a 2 (*SMS Reply Approach*) × 3 (*Driving Condition*) repeated measures, within-



Figure 1. Driving simulator setup.

subjects design experiment in which the order of *SMS Reply* for each *Driving Condition* was counter-balanced. Because our primary variable of interest was *SMS Reply*, we had users experience both *voice search* and *dictation* with *no driving* first, then *easy driving*, followed by *difficult driving*. This gave users a chance to adjust themselves to increasingly difficult road conditions.

Driving Task: As the primary task, users were asked to drive two courses we developed with *easy driving* and *difficult driving* conditions while obeying all rules of the road, as they would in real driving and not in a videogame. With speed limits ranging from 25 mph to 55 mph, both courses contained five sequential sections which took about 15-20 minutes to complete: a residential area, a country highway, and a small city with a downtown area as well as a business/industrial park. Although both courses were almost identical in the number of turns, curves, stops, and traffic lights, the easy course consisted mostly of simple road segments with relatively no traffic, whereas the difficult course had four times as many vehicles, cyclists, and pedestrians. The difficult course also included a foggy road section, a few busy construction sites, and many unexpected events, such as a car in front suddenly breaking, a parked car merging into traffic, and a pedestrian jaywalking. In short, the difficult course was designed to fully engage the attention and cognitive resources of drivers.

SMS Reply Task: As the secondary task, we asked users to listen to an incoming SMS message together with a formulated reply, such as:

- (1) *Message Received:* “Are you lost?” *Your Reply:* “No, never with my GPS”

The users were asked to repeat the reply back to the system. For Example (1) above, users would have to utter “No, never with my GPS”. Users

could also say “Repeat” if they had any difficulties understanding the TTS rendering or if they experienced lapses in attention. For each course, users engaged in 10 SMS reply tasks. SMS messages were cued every 3000 feet, roughly every 90 seconds, which provided enough time to complete each SMS dialogue. Once users uttered the formulated reply, they received a list of 4 possible reply candidates (each labeled as “One”, “Two”, etc.), from which they were asked to either pick the correct reply (by stating its number at any time) or reject them all (by stating “All wrong”). We did not provide any feedback about whether the replies they picked were correct or incorrect in order to avoid priming users to pay more or less attention in subsequent messages. Users did not have to finish listening to the entire list before making their selection.

Stimuli: Because we were interested in examining which was worse, verifying whether SMS response templates matched the meaning of an intended reply, or deciphering the sometimes nonsensical misrecognitions of dictation, we decided to experimentally control both the SMS reply uttered by the user as well as the 4-best list generated by the system. However, all SMS replies and 4-best lists were derived from the logs of an actual SMS Reply interface which implemented the dictation and the voice search approaches (see Ju & Paek, 2009). For each course, 5 of the SMS replies were short (with 3 or fewer words) and 5 were long (with 4 to 7 words). The mean length of the replies was 3.5 words (17.3 chars). The order of the short and long replies was randomized.

We selected 4-best lists where the correct answer was in each of four possible positions (1-4) or All Wrong; that is, there were as many 4-best lists with the first choice correct as there were with the second choice correct, and so forth. We then randomly ordered the presentation of different 4-best lists. Although one might argue that the four positions are not equally likely and that the top item of a 4-best list is most often the correct answer, we decided to experimentally control the position for two reasons: first, our previous research (Ju & Paek, 2009) had already demonstrated the superiority of the voice search approach with respect to the top position (i.e., 1-best), and second, our experimental design sought to identify whether the voice search approach was more usable than the dictation approach even when the ASR accuracy of the two approaches was the same.

In the *dictation* condition, the correct answer was not always an exact copy of the reply in 0-2 of the 10 SMS messages. For instance, a correct dictation answer for Example (1) above was “no I’m never with my GPS”. On the other hand, the *voice search* condition had more cases (2-4 messages) in which the correct answer was not an exact copy (e.g., “no I have GPS”) due to the nature of the template approach. To some degree, this could be seen as handicapping the *voice search* condition, though the results did not reflect the disadvantage, as we discuss later.

Measures: Performance for both the driving task and the SMS reply tasks were recorded. For the driving task, we measured the numbers of collisions, speeding (exceeding 10 mph above the limit), traffic light and stop sign violations, and missed or incorrect turns. For the SMS reply task, we measured duration (i.e., time elapsed between the beginning of the 4-best list and when users ultimately provided their answer) and the number of times users correctly identified which of the 4 reply candidates contained the correct answer.

Originally, we had an independent rater verify the position of the correct answer in all 4-best lists, however, we considered that some participants might be choosing replies that are semantically sufficient, even if they are not exactly correct. For example, a 4-best list generated by the dictation approach for Example (1) had: “One: no I’m never want my GPS. Two: no I’m never with my GPS. Three: no I’m never when my GPS. Or Four: no no I’m in my GPS.” Although the rater identified the second reply as being “correct”, a participant might view the first or third replies as sufficient. In order to avoid ambiguity about correctness, after the study, we showed the same 16 participants the SMS messages and replies as well as the 4-best lists they received during the study and asked them to select, for each SMS reply, any 4-best list items they felt sufficiently conveyed the same meaning, even if the items were ungrammatical. Participants were explicitly told that they could select multiple items from the 4-best list. We did not indicate which item they selected during the experiment and because this selection task occurred months after the experiment, it was unlikely that they would remember anyway. Participants were compensated with a cafeteria voucher.

In computing the number of “correct” answers, for each SMS reply, we counted an an-

swer to be correct if it was included among the participants' set of semantically sufficient 4-best list items. Hence, we calculated the number of correct items in a personalized fashion for every participant.

2.2 Results

We conducted a series of repeated measures ANOVAs on all driving task and SMS reply task measures. For the driving task, we did not find any statistically significant differences between the *voice search* and *dictation* conditions. In other words, we could not reject the null hypothesis that the two approaches were the same in terms of their influence on driving performance. However, for the SMS reply task, we did find a main effect for *SMS Reply Approach* ($F_{1,47} = 81.28, p < .001, \mu_{\text{Dictation}} = 2.13 (.19), \mu_{\text{VoiceSearch}} = .38 (.10)$). As shown in Figure 2, the average number of errors per driving course for *dictation* is roughly 6 times that for *voice search*. We also found a main effect for total duration ($F_{1,47} = 11.94, p < .01, \mu_{\text{Dictation}} = 113.75 \text{ sec } (3.54) \text{ or } 11.4 \text{ sec/reply}, \mu_{\text{VoiceSearch}} = 125.32 \text{ sec } (3.37) \text{ or } 12.5 \text{ sec/reply}$). We discuss our explanation for the shorter duration below. For both errors and duration, we did not find any interaction effects with *Driving Conditions*.

3 Discussion

We conducted a simulator study in order to examine which was worse while driving: verifying whether SMS response templates matched the meaning of an intended reply, or deciphering the sometimes nonsensical misrecognitions of dictation. Our results suggest that deciphering dictation results under the duress of driving leads to more errors. In conducting a post-hoc error analysis, we noticed that participants tended to err when the 4-best lists generated by the dictation approach contained phonetically similar candidate replies. Because it is not atypical for the dictation approach to have n-best list candidates differing from each other in this way, we recommend not utilizing this approach in speech-only user interfaces, unless the n-best list candidates can be made as distinct from each other as possible, phonetically, syntactically and most importantly, semantically. The voice search approach circumvents this problem in two ways: 1) templates were real responses and manually selected and cleaned up during the development phase so there were no grammatical mistakes, and 2) semantically redundant templates can be

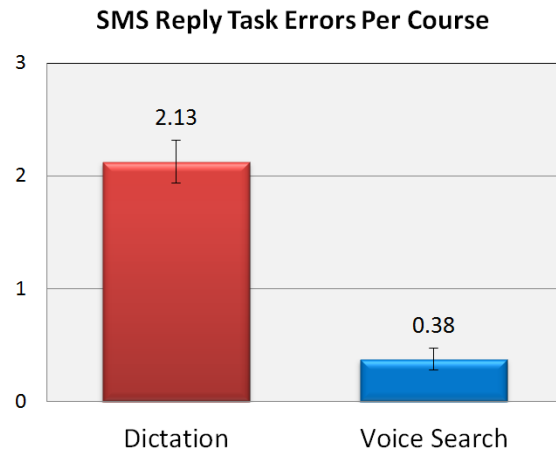


Figure 2. Mean number of errors for the dictation and voice search approaches. Error bars represent standard errors about the mean.

further discarded to only present the distinct concepts at the rendering time using the paraphrase detection algorithms reported in (Wu et al., 2010).

Given that users committed more errors in the *dictation* condition, we initially expected that *dictation* would exhibit higher duration than *voice search* since users might be spending more time figuring out the differences between the similar 4-best list candidates generated by the dictation approach. However, in our error analysis we observed that most likely users did not discover the misrecognitions, and prematurely selected a reply candidate, resulting in shorter durations. The slightly higher duration for the voice search approach does not constitute a problem if users are listening to all of their choices and correctly selecting their intended SMS reply. Note that the duration did not bring about any significant driving performance differences.

Although we did not find any significant driving performance differences, users experienced more difficulties confirming whether the dictation approach correctly interpreted their utterances than they did with the voice search approach. As such, if a user deems it absolutely necessary to respond to SMS messages while driving, our simulator study suggests that the most reliable (i.e., least error-prone) way to respond may just well be the voice search approach.

References

- Distracted Driving Summit. 2009. Department of Transportation. Retrieved Dec. 1: http://www.rita.dot.gov/distracted_driving_summit

- Y.C. Ju & T. Paek. 2009. A Voice Search Approach to Replying to SMS Messages in Automobiles. In *Proc. of Interspeech*.
- A. Kun, T. Paek & Z. Medenica. 2007. The Effect of Speech Interface Accuracy on Driving Performance, In *Proc. of Interspeech*.
- P. Natarajan, R. Prasad, R. Schwartz, & J. Makhoul. 2002. A Scalable Architecture for Directory Assistance Automation. In *Proc. of ICASSP*, pp. 21-24.
- L. Nunes & M. Recarte. 2002. Cognitive Demands of Hands-Free-Phone Conversation While Driving. *Transportation Research Part F*, 5: 133-144.
- N. Reed & R. Robbins. 2008. The Effect of Text Messaging on Driver Behaviour: A Simulator Study. Transport Research Lab Report, PPR 367.
- D. Strayer & W. Johnston. 2001. Driven to Distraction: Dual-task Studies of Simulated Driving and Conversing on a Cellular Phone. *Psychological Science*, 12: 462-466.
- R. Stross. 2008. "What carriers aren't eager to tell you about texting", New York Times, Dec. 26, 2008: http://www.nytimes.com/2008/12/28/business/28digi.html?_r=3
- D. Yu, Y.C. Ju, Y.-Y. Wang, G. Zweig, & A. Acero. 2007. Automated Directory Assistance System: From Theory to Practice. In *Proc. of Interspeech*.
- Wei Wu, Yun-Cheng Ju, Xiao Li, and Ye-Yi Wang, Paraphrase Detection on SMS Messages in Automobiles, in ICASSP, IEEE, March 2010

Classification of Feedback Expressions in Multimodal Data

Costanza Navarretta

University of Copenhagen
Centre for Language Technology (CST)
Njalsgade 140, 2300-DK Copenhagen
costanza@hum.ku.dk

Patrizia Paggio

University of Copenhagen
Centre for Language Technology (CST)
Njalsgade 140, 2300-DK Copenhagen
paggio@hum.ku.dk

Abstract

This paper addresses the issue of how linguistic feedback expressions, prosody and head gestures, i.e. head movements and face expressions, relate to one another in a collection of eight video-recorded Danish map-task dialogues. The study shows that in these data, prosodic features and head gestures significantly improve automatic classification of dialogue act labels for linguistic expressions of feedback.

1 Introduction

Several authors in communication studies have pointed out that head movements are relevant to feedback phenomena (see McClave (2000) for an overview). Others have looked at the application of machine learning algorithms to annotated multimodal corpora. For example, Jokinen and Ragni (2007) and Jokinen et al. (2008) find that machine learning algorithms can be trained to recognise some of the functions of head movements, while Reidsma et al. (2009) show that there is a dependence between focus of attention and assignment of dialogue act labels. Related are also the studies by Rieks op den Akker and Schulz (2008) and Murray and Renals (2008): both achieve promising results in the automatic segmentation of dialogue acts using the annotations in a large multimodal corpus.

Work has also been done on prosody and gestures in the specific domain of map-task dialogues, also targeted in this paper. Sridhar et al. (2009) obtain promising results in dialogue act tagging of the Switchboard-DAMSL corpus using lexical, syntactic and prosodic cues, while Gravano and Hirschberg (2009) examine the relation between particular acoustic and prosodic turn-yielding cues and turn taking in a large corpus of task-oriented dialogues. Louwerse et al. (2006) and Louwerse

et al. (2007) study the relation between eye gaze, facial expression, pauses and dialogue structure in annotated English map-task dialogues (Anderson et al., 1991) and find correlations between the various modalities both within and across speakers. Finally, feedback expressions (head nods and shakes) are successfully predicted from speech, prosody and eye gaze in interaction with Embodied Communication Agents as well as human communication (Fujie et al., 2004; Morency et al., 2005; Morency et al., 2007; Morency et al., 2009).

Our work is in line with these studies, all of which focus on the relation between linguistic expressions, prosody, dialogue content and gestures. In this paper, we investigate how feedback expressions can be classified into different dialogue act categories based on prosodic and gesture features. Our data are made up by a collection of eight video-recorded map-task dialogues in Danish, which were annotated with phonetic and prosodic information. We find that prosodic features improve the classification of dialogue acts and that head gestures, where they occur, contribute to the semantic interpretation of feedback expressions. The results, which partly confirm those obtained on a smaller dataset in Paggio and Navarretta (2010), must be seen in light of the fact that our gesture annotation scheme comprises more fine-grained categories than most of the studies mentioned earlier for both head movements and face expressions. The classification results improve, however, if similar categories such as head nods and jerks are collapsed into a more general category.

In Section 2 we describe the multimodal Danish corpus. In Section 3, we describe how the prosody of feedback expressions is annotated, how their content is coded in terms of dialogue act, turn and agreement labels, and we provide inter-coder agreement measures. In Section 4 we account for the annotation of head gestures, including inter-

coder agreements results. Section 5 contains a description of the resulting datasets and a discussion of the results obtained in the classification experiments. Section 6 is the conclusion.

2 The multimodal corpus

The Danish map-task dialogues from the DanPASS corpus (Grønnum, 2006) are a collection of dialogues in which 11 speaker pairs cooperate on a map task. The dialogue participants are seated in different rooms and cannot see each other. They talk through headsets, and one of them is recorded with a video camera. Each pair goes through four different sets of maps, and changes roles each time, with one subject giving instructions and the other following them. The material is transcribed orthographically with an indication of stress, articulatory hesitations and pauses. In addition to this, the acoustic signals are segmented into words, syllables and prosodic phrases, and annotated with POS-tags, phonological and phonetic transcriptions, pitch and intonation contours.

Phonetic and prosodic segmentation and annotation were performed independently and in parallel by two annotators and then an agreed upon version was produced with the supervision of an expert annotator, for more information see Grønnum (2006). The Praat tool was used (Boersma and Weenink, 2009).

The feedback expressions we analyse here are *Yes* and *No* expressions, i.e. in Danish words like *ja* (yes), *jo* (yes in a negative context), *jamen* (yes but, well), *nej* (no), *næh* (no). They can be single words or multi-word expressions.

Yes and *No* feedback expressions represent about 9% of the approximately 47,000 running words in the corpus. This is a rather high proportion compared to other corpora, both spoken and written, and a reason why we decided to use the DanPASS videos in spite of the fact that the gesture behaviour is relatively limited given the fact that the two dialogue participants cannot see each other. Furthermore, the restricted contexts in which feedback expressions occur in these dialogues allow for a very fine-grained analysis of the relation of these expressions with prosody and gestures. Feedback behaviour, both in speech and gestures, can be observed especially in the person who is receiving the instructions (the *follower*). Therefore, we decided to focus our analysis only on the follower’s part of the interaction. Because

of time restrictions, we limited the study to four different subject pairs and two interactions per pair, for a total of about an hour of video-recorded interaction.

3 Annotation of feedback expressions

As already mentioned, all words in DanPASS are phonetically and prosodically annotated. In the subset of the corpus considered here, 82% of the feedback expressions bear stress or tone information, and 12% are unstressed; 7% of them are marked with onset or offset hesitation, or both. For this study, we added semantic labels – including dialogue acts – and gesture annotation. Both kinds of annotation were carried out using ANVIL (Kipp, 2004). To distinguish among the various functions that feedback expressions have in the dialogues, we selected a subset of the categories defined in the emerging ISO 24617-2 standard for semantic annotation of language resources. This subset comprises the categories *Accept*, *Decline*, *RepeatRephrase* and *Answer*. Moreover, all feedback expressions were annotated with an agreement feature (*Agree*, *NonAgree*) where relevant. Finally, the two turn management categories *TurnTake* and *TurnElicit* were also coded.

It should be noted that the same expression may be annotated with a label for each of the three semantic dimensions. For example, a *yes* can be an *Answer* to a question, an *Agree* and a *TurnElicit* at the same time, thus making the semantic classification very fine-grained. Table 1 shows how the various types are distributed across the 466 feedback expressions in our data.

Dialogue Act		
Answer	70	15%
RepeatRephrase	57	12%
Accept	127	27%
None	212	46%
Agreement		
Agree	166	36%
NonAgree	14	3%
None	286	61%
Turn Management		
TurnTake	113	24%
TurnElicit	85	18%
None	268	58%

Table 1: Distribution of semantic categories

3.1 Inter-coder agreement on feedback expression annotation

In general, dialogue act, agreement and turn annotations were coded by an expert annotator and the annotations were subsequently checked by a second expert annotator. However, one dialogue was coded independently and in parallel by two expert annotators to measure inter-coder agreement. A measure was derived for each annotated feature using the agreement analysis facility provided in ANVIL. Agreement between two annotation sets is calculated here in terms of Cohen’s *kappa* (Cohen, 1960)¹ and corrected *kappa* (Brennan and Prediger, 1981)². Anvil divides the annotations in slices and compares each slice. We used slices of 0.04 seconds. The inter-coder agreement figures obtained for the three types of annotation are given in Table 2.

feature	Cohen’s <i>k</i>	corrected <i>k</i>
agreement	73.59	98.74
dial act	84.53	98.87
turn	73.52	99.16

Table 2: Inter-coder agreement on feedback expression annotation

Although researchers do not totally agree on how to measure agreement in various types of annotated data and on how to interpret the resulting figures, see Artstein and Poesio (2008), it is usually assumed that Cohen’s *kappa* figures over 60 are good while those over 75 are excellent (Fleiss, 1971). Looking at the cases of disagreement we could see that many of these are due to the fact that the annotators had forgotten to remove some of the features automatically proposed by ANVIL from the latest annotated element.

4 Gesture annotation

All communicative head gestures in the videos were found and annotated with ANVIL using a subset of the attributes defined in the MUMIN annotation scheme (Allwood et al., 2007). The MUMIN scheme is a general framework for the study of gestures in interpersonal communication. In this study, we do not deal with functional classification of the gestures in themselves, but rather

¹ $(P_a - P_e)/(1 - P_e)$.

² $(P_o - 1/c)/(1 - 1/c)$ where *c* is the number of categories.

with how gestures contribute to the semantic interpretations of linguistic expressions. Therefore, only a subset of the MUMIN attributes has been used, i.e. *Smile*, *Laughter*, *Scowl*, *FaceOther* for facial expressions, and *Nod*, *Jerk*, *Tilt*, *SideTurn*, *Shake*, *Waggle*, *Other* for head movements.

A link was also established in ANVIL between the gesture under consideration and the relevant speech sequence where appropriate. The link was then used to extract gesture information together with the relevant linguistic annotations on which to apply machine learning.

The total number of head gestures annotated is 264. Of these, 114 (43%) co-occur with feedback expressions, with *Nod* as by far the most frequent type (70 occurrences) followed by *FaceOther* as the second most frequent (16). The other tokens are distributed more or less evenly, with a few occurrences (2-8) per type. The remaining 150 gestures, linked to different linguistic expressions or to no expression at all, comprise many face expressions and a number of tilts. A rough preliminary analysis shows that their main functions are related to focusing or to different emotional attitudes. They will be ignored in what follows.

4.1 Measuring inter-coder agreement on gesture annotation

The head gestures in the DanPASS data have been coded by non expert annotators (one annotator per video) and subsequently controlled by a second annotator, with the exception of one video which was annotated independently and in parallel by two annotators. The annotations of this video were then used to measure inter-coder agreement in ANVIL as it was the case for the annotations on feedback expressions. In the case of gestures we also measured agreement on gesture segmentation. The figures obtained are given in Table 3.

feature	Cohen’s <i>k</i>	corrected <i>k</i>
face segment	69.89	91.37
face annotate	71.53	94.25
head mov segment	71.21	91.75
head mov annotate	71.65	95.14

Table 3: Inter-coder agreement on head gesture annotation

These results are slightly worse than those obtained in previous studies using the same annotation scheme (Jokinen et al., 2008), but are still sat-

isfactory given the high number of categories provided by the scheme.

A distinction that seemed particularly difficult was that between nods and jerks: although the direction of the two movement types is different (down-up and up-down, respectively), the movement quality is very similar, and makes it difficult to see the direction clearly. We return to this point below, in connection with our data analysis.

5 Analysis of the data

The multimodal data we obtained by combining the linguistic annotations from DanPASS with the gesture annotation created in ANVIL, resulted into two different groups of data, one containing all *Yes* and *No* expressions, and the other the subset of those that are accompanied by a face expression or a head movement, as shown in Table 4.

Expression	Count	%
<i>Yes</i>	420	90
<i>No</i>	46	10
Total	466	100
<i>Yes</i> with gestures	102	90
<i>No</i> with gestures	12	10
Total with gestures	114	100

Table 4: *Yes* and *No* datasets

These two sets of data were used for automatic dialogue act classification, which was run in the Weka system (Witten and Frank, 2005). We experimented with various Weka classifiers, comprising Hidden Naive Bayes, SMO, ID3, LADTree and Decision Table. The best results on most of our data were obtained using Hidden Naive Bayes (HNB) (Zhang et al., 2005). Therefore, here we show the results of this classifier. Ten-folds cross-validation was applied throughout.

In the first group of experiments we took into consideration all the *Yes* and *No* expressions (420 *Yes* and 46 *No*) without, however, considering gesture information. The purpose was to see how prosodic information contributes to the classification of dialogue acts. We started by totally leaving out prosody, i.e. only the orthographic transcription (*Yes* and *No* expressions) was considered; then we included information about stress (stressed or unstressed); in the third run we added tone attributes, and in the fourth information on hesitation. Agreement and turn attributes were used in all experiments, while Dialogue act anno-

tation was only used in the training phase. The baseline for the evaluation are the results provided by Weka’s ZeroR classifier, which always selects the most frequent nominal class.

In Table 5 we provide results in terms of precision (P), recall (R) and F-measure (F). These are calculated in Weka as weighted averages of the results obtained for each class.

dataset	Algor	P	R	F
YesNo	ZeroR	27.8	52.8	36.5
	HNB	47.2	53	46.4
+stress	HNB	47.5	54.1	47.1
+stress+tone	HNB	47.8	54.3	47.4
+stress+tone+hes	HNB	47.7	54.5	47.3

Table 5: Classification results with prosodic features

The results indicate that prosodic information improves the classification of dialogue acts with respect to the baseline in all four experiments with improvements of 10, 10.6, 10.9 and 10.8%, respectively. The best results are obtained using information on stress and tone, although the decrease in accuracy when hesitations are introduced is not significant. The confusion matrices show that the classifier is best at identifying *Accept*, while it is very bad at identifying *RepeatRephrase*. This result is not surprising since the former type is much more frequent in the data than the latter, and since prosodic information does not correlate with *RepeatRephrase* in any systematic way.

The second group of experiments was conducted on the dataset where feedback expressions are accompanied by gestures (102 *Yes* and 12 *No*). The purpose this time was to see whether gesture information improves dialogue act classification. We believe it makes sense to perform the test based on this restricted dataset, rather than the entire material, because the portion of data where gestures do accompany feedback expressions is rather small (about 20%). In a different domain, where subjects are less constrained by the technical setting, we expect gestures would make for a stronger and more widespread effect.

The Precision, Recall and F-measure of the ZeroR classifier on these data are 31.5, 56.1 and 40.4, respectively. For these experiments, however, we used as a baseline the results obtained based on stress, tone and hesitation information, the combination that gave the best results on the larger

dataset. Together with the prosodic information, Agreement and turn attributes were included just as earlier, while the dialogue act annotation was only used in the training phase. Face expression and head movement attributes were disregarded in the baseline. We then added face expression alone, head movement alone, and finally both gesture types together. The results are shown in Table 6.

dataset	Algor	P	R	F
YesNo	HNB	43.1	56.1	46.4
+face	HNB	43.7	56.1	46.9
+headm	HNB	44.7	55.3	48.2
+face+headm	HNB	49.9	57	50.3

Table 6: Classification results with head gesture features

These results indicate that adding head gesture information improves the classification of dialogue acts in this reduced dataset, although the improvement is not impressive. The best results are achieved when both face expressions and head movements are taken into consideration.

The confusion matrices show that although the recognition of both *Answer* and *None* improve, it is only the *None* class which is recognised quite reliably. We already explained that in our annotation a large number of feedback utterances have an agreement or turn label without necessarily having been assigned to one of our task-related dialogue act categories. This means that head gestures help distinguishing utterances with an agreement or turn function from other kinds. Looking closer at these utterances, we can see that nods and jerks often occur together with *TurnElicit*, while tilts, side turns and smiles tend to occur with *Agree*.

An issue that worries us is the granularity of the annotation categories. To investigate this, in a third group of experiments we collapsed *Nod* and *Jerk* into a more general category: the distinction had proven difficult for the annotators, and we don't have many jerks in the data. The results, displayed in Table 7, show as expected an improvement. The class which is recognised best is still *None*.

6 Conclusion

In this study we have experimented with the automatic classification of feedback expressions into different dialogue acts in a multimodal corpus of

dataset	Algor	P	R	F
YesNo	HNB	43.1	56.1	46.4
+face	HNB	43.7	56.1	46.9
+headm	HNB	47	57.9	51
+face+headm	HNB	51.6	57.9	53.9

Table 7: Classification results with fewer head movements

Danish. We have conducted three sets of experiments, first looking at how prosodic features contribute to the classification, then testing whether the use of head gesture information improved the accuracy of the classifier, finally running the classification on a dataset in which the head movement types were slightly more general. The results indicate that prosodic features improve the classification, and that in those cases where feedback expressions are accompanied by head gestures, gesture information is also useful. The results also show that using a more coarse-grained distinction of head movements improves classification in these data.

Slightly more than half of the head gestures in our data co-occur with other linguistic utterances than those targeted in this study. Extending our investigation to those, as we plan to do, will provide us with a larger dataset and therefore presumably with even more interesting and reliable results.

The occurrence of gestures in the data studied here is undoubtedly limited by the technical setup, since the two speakers do not see each other. Therefore, we want to investigate the role played by head gestures in other types of video and larger materials. Extending the analysis to larger datasets will also shed more light on whether our gesture annotation categories are too fine-grained for automatic classification.

Acknowledgements

This research has been done under the project VKK (Verbal and Bodily Communication) funded by the Danish Council for Independent Research in the Humanities, and the NOMCO project, a collaborative Nordic project with participating research groups at the universities of Gothenburg, Copenhagen and Helsinki which is funded by the NOS-HS NORDCORP programme. We would also like to thank Nina Grønnum for allowing us to use the DanPASS corpus, and our gesture annotators Josephine Bødker Arrild and Sara Andersen.

References

- Jens Allwood, Loredana Cerrato, Kristiina Jokinen, Costanza Navarretta, and Patrizia Paggio. 2007. The MUMIN Coding Scheme for the Annotation of Feedback, Turn Management and Sequencing. *Multimodal Corpora for Modelling Human Multimodal Behaviour. Special Issue of the International Journal of Language Resources and Evaluation*, 41(3–4):273–287.
- Anne H. Anderson, Miles Bader, Ellen Gurman Bard, Elizabeth Boyle, Gwyneth Doherty, Simon Garrod, Stephen Isard, Jacqueline Kowtko, Jan McAllister, Jim Miller, Catherine Sotillo, Henry S. Thompson, and Regina Weinert. 1991. The HCRC Map Task Corpus. *Language and Speech*, 34:351–366.
- Ron Artstein and Massimo Poesio. 2008. Inter-Coder Agreement for Computational Linguistics. *Computational Linguistics*, 34(4):555–596.
- Paul Boersma and David Weenink, 2009. *Praat: doing phonetics by computer*. Retrieved May 1, 2009, from <http://www.praat.org/>.
- Robert L. Brennan and Dale J. Prediger. 1981. Coefficient Kappa: Some uses, misuses, and alternatives. *Educational and Psychological Measurement*, 41:687–699.
- Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20(1):37–46.
- Joseph L. Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological Bulletin*, 76(5):378–382.
- Shinya Fujie, Y. Ejiri, K. Nakajima, Y Matsusaka, and Tetsunor Kobayashi. 2004. A conversation robot using head gesture recognition as para-linguistic information. In *Proceedings of the 13th IEEE International Workshop on Robot and Human Interactive Communication*, pages 159 – 164, september.
- Agustin Gravano and Julia Hirschberg. 2009. Turn-yielding cues in task-oriented dialogue. In *Proceedings of SIGDIAL 2009: the 10th Annual Meeting of the Special Interest Group in Discourse and Dialogue, September 2009*, pages 253–261, Queen Mary University of London.
- Nina Grønnum. 2006. DanPASS - a Danish phonetically annotated spontaneous speech corpus. In N. Calzolari, K. Choukri, A. Gangemi, B. Maegaard, J. Mariani, J. Odijk, and D. Tapias, editors, *Proceedings of the 5th LREC*, pages 1578–1583, Genoa, May.
- Kristiina Jokinen and Anton Ragni. 2007. Clustering experiments on the communicative properties of gaze and gestures. In *Proceeding of the 3rd. Baltic Conference on Human Language Technologies*, Kaunas, Lithuania, October.
- Kristiina Jokinen, Costanza Navarretta, and Patrizia Paggio. 2008. Distinguishing the communicative functions of gestures. In *Proceedings of the 5th MLMI*, LNCS 5237, pages 38–49, Utrecht, The Netherlands, September. Springer.
- Michael Kipp. 2004. *Gesture Generation by Imitation - From Human Behavior to Computer Character Animation*. Ph.D. thesis, Saarland University, Saarbruecken, Germany, Boca Raton, Florida, dissertation.com.
- Max M. Louwerse, Patrick Jeuniaux, Mohammed E. Hoque, Jie Wu, and Gwineth Lewis. 2006. Multimodal communication in computer-mediated map task scenarios. In R. Sun and N. Miyake, editors, *Proceedings of the 28th Annual Conference of the Cognitive Science Society*, pages 1717–1722, Mahwah, NJ: Erlbaum.
- Max M. Louwerse, Nick Benesh, Mohammed E. Hoque, Patrick Jeuniaux, Gwineth Lewis, Jie Wu, and Megan Zirnstein. 2007. Multimodal communication in face-to-face conversations. In R. Sun and N. Miyake, editors, *Proceedings of the 29th Annual Conference of the Cognitive Science Society*, pages 1235–1240, Mahwah, NJ: Erlbaum.
- Evelyn McClave. 2000. Linguistic functions of head movements in the context of speech. *Journal of Pragmatics*, 32:855–878.
- Louis-Philippe Morency, Candace Sidner, Christopher Lee, and Trevor Darrell. 2005. Contextual Recognition of Head Gestures. In *Proceedings of the International Conference on Multi-modal Interfaces*.
- Louis-Philippe Morency, Candace Sidner, Christopher Lee, and Trevor Darrell. 2007. Head gestures for perceptual interfaces: The role of context in improving recognition. *Artificial Intelligence*, 171(8–9):568–585.
- Louis-Philippe Morency, Iwan de Kok, and Jonathan Gratch. 2009. A probabilistic multimodal approach for predicting listener backchannels. *Autonomous Agents and Multi-Agent Systems*, 20:70–84, Springer.
- Gabriel Murray and Steve Renals. 2008. Detecting Action Meetings in Meetings. In *Proceedings of the 5th MLMI*, LNCS 5237, pages 208–213, Utrecht, The Netherlands, September. Springer.
- Harm Rieks op den Akker and Christian Schulz. 2008. Exploring features and classifiers for dialogue act segmentation. In *Proceedings of the 5th MLMI*, pages 196–207.
- Patrizia Paggio and Costanza Navarretta. 2010. Feedback in Head Gesture and Speech. To appear in *Proceedings of 7th Conference on Language Resources and Evaluation (LREC-2010)*, Malta, May.

- Dennis Reidsma, Dirk Heylen, and Harm Rieks op den Akker. 2009. On the Contextual Analysis of Agreement Scores. In Michael Kipp, Jean-Claude Martin, Patrizia Paggio, and Dirk Heylen, editors, *Multimodal Corpora From Models of Natural Interaction to Systems and Applications*, number 5509 in Lecture Notes in Artificial Intelligence, pages 122–137. Springer.
- Vivek Kumar Rangarajan Sridhar, Srinivas Bangaloreb, and Shrikanth Narayanan. 2009. Combining lexical, syntactic and prosodic cues for improved online dialog act tagging. *Computer Speech & Language*, 23(4):407–422.
- Ian H. Witten and Eibe Frank. 2005. *Data Mining: Practical machine learning tools and techniques*. Morgan Kaufmann, San Francisco, second edition.
- Harry Zhang, Liangxiao Jiang, and Jiang Su. 2005. Hidden Naive Bayes. In *Proceedings of the Twentieth National Conference on Artificial Intelligence*, pages 919–924.

Optimizing Informativeness and Readability for Sentiment Summarization

Hitoshi Nishikawa, Takaaki Hasegawa, Yoshihiro Matsuo and Genichiro Kikui

NTT Cyber Space Laboratories, NTT Corporation

1-1 Hikari-no-oka, Yokosuka, Kanagawa, 239-0847 Japan

{ nishikawa.hitoshi, hasegawa.takaaki }
{ matsuo.yoshihiro, kikui.genichiro } @lab.ntt.co.jp

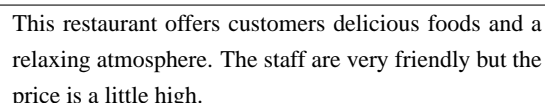
Abstract

We propose a novel algorithm for sentiment summarization that takes account of informativeness and readability, simultaneously. Our algorithm generates a summary by selecting and ordering sentences taken from multiple review texts according to two scores that represent the informativeness and readability of the sentence order. The informativeness score is defined by the number of sentiment expressions and the readability score is learned from the target corpus. We evaluate our method by summarizing reviews on restaurants. Our method outperforms an existing algorithm as indicated by its ROUGE score and human readability experiments.

1 Introduction

The Web holds a massive number of reviews describing the sentiments of customers about products and services. These reviews can help the user reach purchasing decisions and guide companies' business activities such as product improvements. It is, however, almost impossible to read all reviews given their sheer number.

These reviews are best utilized by the development of automatic text summarization, particularly sentiment summarization. It enables us to efficiently grasp the key bits of information. Sentiment summarizers are divided into two categories in terms of output style. One outputs lists of sentences (Hu and Liu, 2004; Blair-Goldensohn et al., 2008; Titov and McDonald, 2008), the other outputs texts consisting of ordered sentences (Carenini et al., 2006; Carenini and Cheung, 2008; Lerman et al., 2009; Lerman and McDonald, 2009). Our work lies in the latter category, and a typical summary is shown in Figure 1. Although visual representations such as bar or rader charts



This restaurant offers customers delicious foods and a relaxing atmosphere. The staff are very friendly but the price is a little high.

Figure 1: A typical summary.

are helpful, such representations necessitate some simplifications of information to presentation. In contrast, text can present complex information that can't readily be visualized, so in this paper we focus on producing textual summaries.

One crucial weakness of existing text-oriented summarizers is the poor readability of their results. Good readability is essential because readability strongly affects text comprehension (Barzilay et al., 2002).

To achieve readable summaries, the extracted sentences must be appropriately ordered (Barzilay et al., 2002; Lapata, 2003; Barzilay and Lee, 2004; Barzilay and Lapata, 2005). Barzilay et al. (2002) proposed an algorithm for ordering sentences according to the dates of the publications from which the sentences were extracted. Lapata (2003) proposed an algorithm that computes the probability of two sentences being adjacent for ordering sentences. Both methods delink sentence extraction from sentence ordering, so a sentence can be extracted that cannot be ordered naturally with the other extracted sentences.

To solve this problem, we propose an algorithm that chooses sentences and orders them simultaneously in such a way that the ordered sentences maximize the scores of informativeness and readability. Our algorithm efficiently searches for the best sequence of sentences by using dynamic programming and beam search. We verify that our method generates summaries that are significantly better than the baseline results in terms of ROUGE score (Lin, 2004) and subjective readability measures. As far as we know, this is the first work to

simultaneously achieve both informativeness and readability in the area of multi-document summarization.

This paper is organized as follows: Section 2 describes our summarization method. Section 3 reports our evaluation experiments. We conclude this paper in Section 4.

2 Optimizing Sentence Sequence

Formally, we define a summary $S^* = \langle s_0, s_1, \dots, s_n, s_{n+1} \rangle$ as a sequence consisting of n sentences where s_0 and s_{n+1} are symbols indicating the beginning and ending of the sequence, respectively. Summary S^* is also defined as follows:

$$S^* = \operatorname{argmax}_{S \in T} [\operatorname{Info}(S) + \lambda \operatorname{Read}(S)] \quad (1)$$

s.t. $\operatorname{length}(S) \leq K$

where $\operatorname{Info}(S)$ indicates the informativeness score of S , $\operatorname{Read}(S)$ indicates the readability score of S , T indicates possible sequences composed of sentences in the target documents, λ is a weight parameter balancing informativeness against readability, $\operatorname{length}(S)$ is the length of S , and K is the maximum size of the summary.

We introduce the informativeness score and the readability score, then describe how to optimize a sequence.

2.1 Informativeness Score

Since we attempt to summarize reviews, we assume that a good summary must involve as many sentiments as possible. Therefore, we define the informativeness score as follows:

$$\operatorname{Info}(S) = \sum_{e \in E(S)} f(e) \quad (2)$$

where e indicates sentiment $e = \langle a, p \rangle$ as the tuple of *aspect* a and *polarity* $p = \{-1, 0, 1\}$, $E(S)$ is the set of sentiments contained S , and $f(e)$ is the score of sentiment e . Aspect a represents a standpoint for evaluating products and services. With regard to restaurants, aspects include *food*, *atmosphere* and *staff*. Polarity represents whether the sentiment is positive or negative. In this paper, we define $p = -1$ as negative, $p = 0$ as neutral and $p = 1$ as positive sentiment.

Notice that Equation 2 defines the informativeness score of a summary as the sum of the score of the sentiments contained in S . To avoid duplicative sentences, each sentiment is counted only

once for scoring. In addition, the aspects are clustered and similar aspects (e.g. *air*, *ambience*) are treated as the same aspect (e.g. *atmosphere*). In this paper we define $f(e)$ as the frequency of e in the target documents.

Sentiments are extracted using a sentiment lexicon and pattern matched from dependency trees of sentences. The sentiment lexicon¹ consists of pairs of *sentiment expressions* and their polarities, for example, *delicious*, *friendly* and *good* are positive sentiment expressions, *bad* and *expensive* are negative sentiment expressions.

To extract sentiments from given sentences, first, we identify sentiment expressions among words consisting of parsed sentences. For example, in the case of the sentence ‘‘This restaurant offers customers delicious foods and a relaxing atmosphere.’’ in Figure 1, *delicious* and *relaxing* are identified as sentiment expressions. If the sentiment expressions are identified, the expressions and its aspects are extracted as aspect-sentiment expression pairs from dependency tree using some rules. In the case of the example sentence, *foods* and *delicious*, *atmosphere* and *relaxing* are extracted as aspect-sentiment expression pairs. Finally extracted sentiment expressions are converted to polarities, we acquire the set of sentiments from sentences, for example, $\langle \textit{foods}, 1 \rangle$ and $\langle \textit{atmosphere}, 1 \rangle$.

Note that since our method relies on only sentiment lexicon, extractable aspects are unlimited.

2.2 Readability Score

Readability consists of various elements such as conciseness, coherence, and grammar. Since it is difficult to model all of them, we approximate readability as the natural order of sentences.

To order sentences, Barzilay et al. (2002) used the publication dates of documents to catch temporally-ordered events, but this approach is not really suitable for our goal because reviews focus on entities rather than events. Lapata (2003) employed the probability of two sentences being adjacent as determined from a corpus. If the corpus consists of reviews, it is expected that this approach would be effective for sentiment summarization. Therefore, we adopt and improve Lapata’s approach to order sentences. We define the

¹Since we aim to summarize Japanese reviews, we utilize Japanese sentiment lexicon (Asano et al., 2008). However, our method is, except for sentiment extraction, language independent.

readability score as follows:

$$\text{Read}(S) = \sum_{i=0}^n \mathbf{w}^\top \phi(s_i, s_{i+1}) \quad (3)$$

where, given two adjacent sentences s_i and s_{i+1} , $\mathbf{w}^\top \phi(s_i, s_{i+1})$, which measures the connectivity of the two sentences, is the inner product of \mathbf{w} and $\phi(s_i, s_{i+1})$, \mathbf{w} is a parameter vector and $\phi(s_i, s_{i+1})$ is a feature vector of the two sentences. That is, the readability score of sentence sequence S is the sum of the connectivity of all adjacent sentences in the sequence.

As the features, Lapata (2003) proposed the Cartesian product of content words in adjacent sentences. To this, we add named entity tags (e.g. LOC, ORG) and connectives. We observe that the first sentence of a review of a restaurant frequently contains named entities indicating location. We aim to reproduce this characteristic in the ordering.

We also define feature vector $\Phi(S)$ of the entire sequence $S = \langle s_0, s_1, \dots, s_n, s_{n+1} \rangle$ as follows:

$$\Phi(S) = \sum_{i=0}^n \phi(s_i, s_{i+1}) \quad (4)$$

Therefore, the score of sequence S is $\mathbf{w}^\top \Phi(S)$. Given a training set, if a trained parameter \mathbf{w} assigns a score $\mathbf{w}^\top \Phi(S^+)$ to an correct order S^+ that is higher than a score $\mathbf{w}^\top \Phi(S^-)$ to an incorrect order S^- , it is expected that the trained parameter will give higher score to naturally ordered sentences than to unnaturally ordered sentences.

We use Averaged Perceptron (Collins, 2002) to find \mathbf{w} . Averaged Perceptron requires an argmax operation for parameter estimation. Since we attempt to order a set of sentences, the operation is regarded as solving the Traveling Salesman Problem; that is, we locate the path that offers maximum score through all n sentences as s_0 and s_{n+1} are starting and ending points, respectively. Thus the operation is NP-hard and it is difficult to find the global optimal solution. To alleviate this, we find an approximate solution by adopting the dynamic programming technique of the Held and Karp Algorithm (Held and Karp, 1962) and beam search.

We show the search procedure in Figure 2. \mathbf{S} indicates intended sentences and \mathbf{M} is a distance matrix of the readability scores of adjacent sentence pairs. $\mathbf{H}^i(\mathbf{C}, j)$ indicates the score of the hypothesis that has covered the set of i sentences \mathbf{C} and has the sentence j at the end of the path,

Sentences: $\mathbf{S} = \{s_1, \dots, s_n\}$
Distance matrix: $\mathbf{M} = [a_{i,j}]_{i=0 \dots n+1, j=0 \dots n+1}$
1: $\mathbf{H}^0(\{s_0\}, s_0) = 0$
2: for $i : 0 \dots n - 1$
3: for $j : 1 \dots n$
4: foreach $\mathbf{H}^i(\mathbf{C} \setminus \{j\}, k) \in \mathbf{b}$
5: $\mathbf{H}^{i+1}(\mathbf{C}, j) = \max_{\mathbf{H}^i(\mathbf{C} \setminus \{j\}, k) \in \mathbf{b}} \mathbf{H}^i(\mathbf{C} \setminus \{j\}, k)$
6: $\quad \quad \quad + \mathbf{M}_{k,j}$
7: $\mathbf{H}^* = \max_{\mathbf{H}_n(\mathbf{C}, k)} \mathbf{H}^n(\mathbf{C}, k) + \mathbf{M}_{k, n+1}$

Figure 2: Held and Karp Algorithm.

i.e. the last sentence of the summary being generated. For example, $\mathbf{H}^2(\{s_0, s_2, s_5\}, s_2)$ indicates a hypothesis that covers s_0, s_2, s_5 and the last sentence is s_2 . Initially, $\mathbf{H}^0(\{s_0\}, s_0)$ is assigned the score of 0, and new sentences are then added one by one. In the search procedure, our dynamic programming based algorithm retains just the hypothesis with maximum score among the hypotheses that have the same sentences and the same last sentence. Since this procedure is still computationally hard, only the top \mathbf{b} hypotheses are expanded.

Note that our method learns \mathbf{w} from texts automatically annotated by a POS tagger and a named entity tagger. Thus manual annotation isn't required.

2.3 Optimization

The argmax operation in Equation 1 also involves search, which is NP-hard as described in Section 2.2. Therefore, we adopt the Held and Karp Algorithm and beam search to find approximate solutions. The search algorithm is basically the same as parameter estimation, except for its calculation of the informativeness score and size limitation. Therefore, when a new sentence is added to a hypothesis, both the informativeness and the readability scores are calculated. The size of the hypothesis is also calculated and if the size exceeds the limit, the sentence can't be added. A hypothesis that can't accept any more sentences is removed from the search procedure and preserved in memory. After all hypotheses are removed, the best hypothesis is chosen from among the preserved hypotheses as the solution.

3 Experiments

This section evaluates our method in terms of ROUGE score and readability. We collected 2,940 reviews of 100 restaurants from a website. The

	R-2	R-SU4	R-SU9
Baseline	0.089	0.068	0.062
Method1	0.157	0.096	0.089
Method2	0.172	0.107	0.098
Method3	0.180	0.110	0.101
Human	0.258	0.143	0.131

Table 1: Automatic ROUGE evaluation.

average size of each document set (corresponds to one restaurant) was 5,343 bytes. We attempted to generate 300 byte summaries, so the summarization rate was about 6%. We used CRFs-based Japanese dependency parser (Imamura et al., 2007) and named entity recognizer (Suzuki et al., 2006) for sentiment extraction and constructing feature vectors for readability score, respectively.

3.1 ROUGE

We used ROUGE (Lin, 2004) for evaluating the content of summaries. We chose ROUGE-2, ROUGE-SU4 and ROUGE-SU9. We prepared four reference summaries for each document set.

To evaluate the effects of the informativeness score, the readability score and the optimization, we compared the following five methods.

Baseline: employs MMR (Carbonell and Goldstein, 1998). We designed the score of a sentence as term frequencies of the content words in a document set.

Method1: uses optimization without the informativeness score or readability score. It also used term frequencies to score sentences.

Method2: uses the informativeness score and optimization without the readability score.

Method3: the proposed method. Following Equation 1, the summarizer searches for a sequence with high informativeness and readability score. The parameter vector w was trained on the same 2,940 reviews in 5-fold cross validation fashion. λ was set to 6,000 using a development set.

Human is the reference summaries. To compare our summarizer to human summarization, we calculated ROUGE scores between each reference and the other references, and averaged them.

The results of these experiments are shown in Table 1. ROUGE scores increase in the order of Method1, Method2 and Method3 but no method could match the performance of Human. The methods significantly outperformed Baseline ac-

	Numbers
Baseline	1.76
Method1	4.32
Method2	10.41
Method3	10.18
Human	4.75

Table 2: Unique sentiment numbers.

ording to the Wilcoxon signed-rank test.

We discuss the contribution of readability to ROUGE scores. Comparing Method2 to Method3, ROUGE scores of the latter were higher for all criteria. It is interesting that the readability criterion also improved ROUGE scores.

We also evaluated our method in terms of sentiments. We extracted sentiments from the summaries using the above sentiment extractor, and averaged the unique sentiment numbers. Table 2 shows the results.

The references (Human) have fewer sentiments than the summaries generated by our method. In other words, the references included almost as many other sentences (e.g. reasons for the sentiments) as those expressing sentiments. Carenini et al. (2006) pointed out that readers wanted “detailed information” in summaries, and the reasons are one of such piece of information. Including them in summaries would greatly improve summarizer appeal.

3.2 Readability

Readability was evaluated by human judges. Three different summarizers generated summaries for each document set. Ten judges evaluated the thirty summaries for each. Before the evaluation the judges read evaluation criteria and gave points to summaries using a five-point scale. The judges weren’t informed of which method generated which summary.

We compared three methods; Ordering sentences according to publication dates and positions in which sentences appear after sentence extraction (**Method2**), Ordering sentences using the readability score after sentence extraction (**Method2+**) and searching a document set to discover the sequence with the highest score (**Method3**).

Table 3 shows the results of the experiment. Readability increased in the order of Method2, Method2+ and Method3. According to the

	Readability point
Method2	3.45
Method2+	3.54
Method3	3.74

Table 3: Readability evaluation.

Wilcoxon signed-rank test, there was no significance difference between Method2 and Method2+ but the difference between Method2 and Method3 was significant, $p < 0.10$.

One important factor behind the higher readability of Method3 is that it yields longer sentences on average (6.52). Method2 and Method2+ yielded averages of 7.23 sentences. The difference is significant as indicated by $p < 0.01$. That is, Method2 and Method2+ tended to select short sentences, which made their summaries less readable.

4 Conclusion

This paper proposed a novel algorithm for sentiment summarization that takes account of informativeness and readability, simultaneously. To summarize reviews, the informativeness score is based on sentiments and the readability score is learned from a corpus of reviews. The preferred sequence is determined by using dynamic programming and beam search. Experiments showed that our method generated better summaries than the baseline in terms of ROUGE score and readability.

One future work is to include important information other than sentiments in the summaries. We also plan to model the order of sentences globally. Although the ordering model in this paper is local since it looks at only adjacent sentences, a model that can evaluate global order is important for better summaries.

Acknowledgments

We would like to sincerely thank Tsutomu Hirao for his comments and discussions. We would also like to thank the reviewers for their comments.

References

Hisako Asano, Toru Hirano, Nozomi Kobayashi and Yoshihiro Matsuo. 2008. Subjective Information Indexing Technology Analyzing Word-of-mouth Content on the Web. *NTT Technical Review*, Vol.6, No.9.

Regina Barzilay, Noemie Elhadad and Kathleen McKeown. 2002. Inferring Strategies for Sentence Ordering in Multidocument Summarization. *Journal of Artificial Intelligence Research (JAIR)*, Vol.17, pp. 35–55.

Regina Barzilay and Lillian Lee. 2004. Catching the Drift: Probabilistic Content Models, with Applications to Generation and Summarization. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics (HLT-NAACL)*, pp. 113–120.

Regina Barzilay and Mirella Lapata. 2005. Modeling Local Coherence: An Entity-based Approach. In *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL)*, pp. 141–148.

Sasha Blair-Goldensohn, Kerry Hannan, Ryan McDonald, Tyler Neylon, George A. Reis and Jeff Reynar. 2008. Building a Sentiment Summarizer for Local Service Reviews. WWW Workshop NLP Challenges in the Information Explosion Era (NLPiX).

Jaime Carbonell and Jade Goldstein. 1998. The use of MMR, diversity-based reranking for reordering documents and producing summaries. In *Proceedings of the 21st annual international ACM SIGIR conference on Research and development in information retrieval (SIGIR)*, pp. 335–356.

Giuseppe Carenini, Raymond Ng and Adam Pauls. 2006. Multi-Document Summarization of Evaluative Text. In *Proceedings of the 11th European Chapter of the Association for Computational Linguistics (EACL)*, pp. 305–312.

Giuseppe Carenini and Jackie Chi Kit Cheung. 2008. Extractive vs. NLG-based Abstractive Summarization of Evaluative Text: The Effect of Corpus Controversiality. In *Proceedings of the 5th International Natural Language Generation Conference (INLG)*, pp. 33–41.

Michael Collins. 2002. Discriminative Training Methods for Hidden Markov Models: Theory and Experiments with Perceptron Algorithms. In *Proceedings of the 2002 Conference on Empirical Methods on Natural Language Processing (EMNLP)*, pp. 1–8.

Michael Held and Richard M. Karp. 1962. A dynamic programming approach to sequencing problems. *Journal of the Society for Industrial and Applied Mathematics (SIAM)*, Vol.10, No.1, pp. 196–210.

Minqing Hu and Bing Liu. 2004. Mining and Summarizing Customer Reviews. In *Proceedings of the 10th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, pp. 168–177.

- Kenji Imamura, Genichiro Kikui and Norihito Yasuda. 2007. Japanese Dependency Parsing Using Sequential Labeling for Semi-spoken Language. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics (ACL) Companion Volume Proceedings of the Demo and Poster Sessions*, pp. 225–228.
- Mirella Lapata. 2003. Probabilistic Text Structuring: Experiments with Sentence Ordering. In *Proceedings of the 41st Annual Meeting of the Association for Computational Linguistics (ACL)*, pp. 545–552.
- Kevin Lerman, Sasha Blair-Goldensohn and Ryan McDonald. 2009. Sentiment Summarization: Evaluating and Learning User Preferences. In *Proceedings of the 12th Conference of the European Chapter of the Association for Computational Linguistics (EACL)*, pp. 514–522.
- Kevin Lerman and Ryan McDonald. 2009. Contrastive Summarization: An Experiment with Consumer Reviews. In *Proceedings of Human Language Technologies: the 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL-HLT), Companion Volume: Short Papers*, pp. 113–116.
- Chin-Yew Lin. 2004. ROUGE: A Package for Automatic Evaluation of Summaries. In *Proceedings of the Workshop on Text Summarization Branches Out*, pp. 74–81.
- Jun Suzuki, Erik McDermott and Hideki Isozaki. 2006. Training Conditional Random Fields with Multivariate Evaluation Measures. In *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the ACL (COLING-ACL)*, pp. 217–224.
- Ivan Titov and Ryan McDonald. 2008. A Joint Model of Text and Aspect Ratings for Sentiment Summarization. In *Proceedings of the 46th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies (ACL-HLT)*, pp. 308–316.

Last but *Definitely* not Least:

On the Role of the Last Sentence in Automatic Polarity-Classification

Israela Becker and Vered Aharonson

AFEKA – Tel-Aviv Academic College of Engineering
218 Bney-Efraim Rd.
Tel-Aviv 69107, Israel
{IsraelaB,Vered}@afeka.ac.il

Abstract

Two psycholinguistic and psychophysical experiments show that in order to efficiently extract polarity of written texts such as customer-reviews on the Internet, one should concentrate computational efforts on messages in the final position of the text.

1 Introduction

The ever-growing field of polarity-classification of written texts may benefit greatly from linguistic insights and tools that will allow to efficiently (and thus economically) extract the polarity of written texts, in particular, online customer reviews.

Many researchers interpret “efficiently” as using better computational methods to resolve the polarity of written texts. We suggest that text units should be handled with tools of discourse linguistics too in order to reveal where, within texts, their polarity is best manifested. Specifically, we propose to focus on the last sentence of the given text in order to efficiently extract the polarity of the whole text. This will reduce computational costs, as well as improve the quality of polarity detection and classification when large databases of text units are involved.

This paper aims to provide psycholinguistic support to the hypothesis (which psycholinguistic literature lacks) that the last sentence of a customer review is a better predictor for the polarity of the whole review than other sentences in the review, in order to be later used for automatic polarity-classification. Therefore, we first briefly review the well-established structure of

text units while comparing notions of topic-extraction vs. *our* notion of polarity-classification. We then report the psycholinguistic experiments that we ran in order to support our prediction as to the role of the last sentence in polarity manifestation. Finally, we discuss the experimental results.

2 Topic-extraction

One of the basic features required to perform automatic topic-extraction is sentence position. The importance of sentence position for computational purposes was first indicated by Baxendale in the late 1950s (Baxendale, 1958): Baxendale hypothesized that the first and the last sentence of a given text are the potential topic-containing sentences. He tested this hypothesis on a corpus of 200 paragraphs extracted out of 6 technical articles. He found that in 85% of the documents, the first sentence was the topic sentence, whereas in only 7% of the documents, it was the last sentence. A large scale study supporting Baxendale’s hypothesis was conducted by Lin and Hovy (Lin and Hovy, 1997) who examined 13,000 documents of the Ziff-Davis newswire corpus of articles reviewing computer hardware and software. In this corpus, each document was accompanied by a set of topic keywords and a small abstract of six sentences. Lin and Hovy measured the yield of each sentence against the topic keywords and ranked the sentences by their average yield. They concluded that in $\sim 2/3$ of the documents, the topic keywords are indeed mentioned in the title and first five sentences of the document.

Baxendale’s theory gained further psycholinguistic support by the experimental results of Kieras (Kieras, 1978, Kieras, 1980) who showed that subjects re-constructed the content

of paragraphs they were asked to read by relying on sentences in initial positions. These findings subsequently gained extensive theoretical and experimental support by Giora (Giora, 1983, Giora, 1985) who correlated the position of a sentence within a text with its degree of informativeness.

Giora (Giora, 1985, Giora, 1988) defined a discourse topic (DT) as the least informative (*most uninformative*) yet *dominant* proposition of a text. The DT best represents the redundancy structure of the text. As such, this proposition functions as a reference point for processing the rest of the propositions. The text position which best benefits such processing is text initial; it facilitates processing of oncoming propositions (with respect to the DT) relative to when the DT is placed in text final position.

Furthermore, Giora and Lee showed (Giora and Lee, 1996) that when the DT appears *also* at the end of a text it is somewhat informationally redundant. However, functionally, it plays a role in wrapping the text up and marking its boundary. Authors often make reference to the DT at the end of a text in order to summarize and deliberately recapitulate what has been written up to that point while also signaling the end of discourse topic segment.

3 Polarity-classification vs. Topic-extraction

When dealing with polarity-classification (as with topic-extraction), one should again identify the *most uninformative* yet *dominant* proposition of the text. However, given the cognitive prominence of discourse final position in terms of memorability, known as “*recency effect*” (see below and see also (Giora, 1988)), we predict that when it comes to polarity-classification, the *last* proposition of a given text should be of greater importance than the *first* one (contrary to topic-extraction).

Based on preliminary investigations, we suggest that the DT of any customer review is the customer’s evaluation, whether negative or positive, of a product that s/he has purchased or a service s/he has used, rather than the details of the specific product or service. The message that customer reviews try to get across is, therefore, of evaluative nature. To best communicate this affect, the DT should appear at the end of the review (instead of the beginning of the review) as a means of recapitulating the point of the message, thereby guaranteeing that it is fully understood by the readership.

Indeed, the cognitive prominence of information in final position - the *recency-effect* - has been well established in numerous psychological experiments (see, for example, (Murdock, 1962)). Thus, the most *frequent* evaluation of the product (which is the *most uninformative* one) also should surface at the end of the text due to the ease of its retrieval, which is presumably what product review readers would refer to as “the bottom line”.

To the best of our knowledge, this psycholinguistic prediction has not been supported by psycholinguistic evidence to date. However, it has been somewhat supported by the computational results of Yang, Lin and Chen (Yang et al., 2007a, Yang et al., 2007b) who classified emotions of posts in blog corpora. Yang, Lin & Chen realized that bloggers tend to emphasize their feelings by using emoticons (such as: ☺, ☹ and 😊) and that these emoticons frequently appear in final sentences. Thus, they first focused on the last sentence of posts as representing the polarity of the entire posts. Then, they divided the positive category into 2 sub-categories - happy and joy, and the negative category - into angry and sad. They showed that extracting polarity and consequently sentiments from last sentences outperforms all other computational strategies.

4 Method

We aim to show that the last sentence of a customer review is a better predictor for the polarity of the whole review than any other sentence (assuming that the first sentence is devoted to presenting the product or service). To test our prediction, we ran two experiments and compared their results. In the first experiment we examined the readers’ rating of the polarity of reviews in their entirety, while in the second experiment we examined the readers’ rating of the same reviews based on reading single sentences extracted from these reviews: the last sentence *or* the second one. The second sentence could have been replaced by any other sentence, *but* the first one, as our preliminary investigations clearly show that the first sentence is in many cases devoted to presenting the product or service discussed and does not contain any polarity content. For example: “I read Isaac’s storm, by Erik Larson, around 1998. Recently I had occasion to thumb through it again which has prompted this review.....All in all a most interesting and rewarding book, one that I would recommend highly.” (Gerald T. Westbrook, “GTW”)

4.1 Materials

Sixteen customer-reviews were extracted from Blitzer, Dredze, and Pereira's sentiment database (Blitzer et al., 2007). This database contains product-reviews taken from Amazon¹ where each review is rated by its author on a 1-5 star scale. The database covers 4 product types (domains): Kitchen, Books, DVDs, and Electronics. Four reviews were selected from each domain. Of the 16 extracted reviews, 8 were positive (4-5 star rating) and the other 8 – negative (1-2 star rating).

Given that in this experiment we examine the polarity of the last sentence relative to that of the whole review or to a few other sentences, we focused on the first reviews (as listed in the aforementioned database) of at least 5 sentences or longer, rather than on too-short reviews. By “too-short” we refer to reviews in which such comparison would be meaningless; for example, ones that range between 1-3 sentences will not allow to compare the last sentence with any of the others.

4.2 Participants

Thirty-five subjects participated in the first experiment: 14 women and 21 men, ranging in age from 22 to 73. Thirty-six subjects participated in the second experiment: 23 women and 13 men ranging in age from 20 to 59. All participants were native speakers of English, had an academic education, and had normal or corrected-to-normal eye-vision.

4.3 Procedure

In the first experiment, subjects were asked to read 16 reviews; in the second experiment subjects were asked to read 32 single sentences extracted from the same 16 reviews: the last sentence and the second sentence of each review. The last and the second sentence of each review were *not* presented together but *individually*.

In both experiments subjects were asked to guess the ratings of the texts which were given by the authors on a 1-5 star scale, by clicking on a radio-button: “In each of the following screens you will be asked to read a customer review (or a sentence extracted out of a customer review). All the reviews were extracted from the www.amazon.com customer review section. Each review (or sentence) describes a different product. At the end of each review (or sentence)

you will be asked to decide whether the reviewer who wrote the review recommended or did not recommend the reviewed product on a 1-5 scale: Number 5 indicates that the reviewer highly recommended the product, while number 1 indicates that the reviewer was unsatisfied with the product and did not recommend it.”

In the second experiment, in addition to the psychological experiment, the latencies following reading of the texts up until the clicking of the mouse, as well as the biometric measurements of the mouse's trajectories, were recorded.

In both experiments each subject was run in an individual session and had an unlimited time to reflect and decide on the polarity of each text. Five seconds after a decision was made (as to whether the reviewer was in favor of the product or not), the subject was presented with the next text. The texts were presented in random order so as to prevent possible interactions between them.

In the initial design phase of the experiment we discussed the idea of adding an “irrelevant” option in addition to the 5-star scale of polarity. This option was meant to be used for sentences that carry no evaluation at all. Such an addition would have necessitated locating the extra-choice radio button at a separated remote place from the 5-star scale radio buttons, since conceptually it cannot be located on a nearby position. From the user interaction point of view, the mouse movement to that location would have been either considerably shorter or longer (depending on its distance from the initial location of the mouse cursor at the beginning of each trial), and the mouse trajectory and click time would have been, thus, very different and difficult to analyze.

Although the reviews were randomly selected, 32 sentences extracted out of 16 reviews might seem like a small sample. However, the upper time limit for *reliable* psycholinguistic experiments is 20-25 minute. Although tempted to extend the experiments in order to acquire more data, longer times result in subject impatience, which shows on lower scoring rates. Therefore, we chose to trade sample size for accuracy. Experimental times in both experiments ranged between 15-35 minutes.

5 Results

Results of the distribution of differences between the authors' and the readers' ratings of the texts are presented in Figure 1: The distribution of differences for whole reviews is (un-surprisingly) the narrowest (Figure 1a). The dis-

¹ <http://www.amazon.com>

tribution of differences for last sentences (Figure 1b) is somewhat wider than (but still quite similar to) the distribution of differences for whole reviews. The distribution of differences for second sentences is the widest of the three (Figure 1c).

Pearson correlation coefficient calculations (Table 1) show that both the correlation between authors' ratings and readers' rating for whole reviews and the correlation between authors' rating and readers' rating upon reading the last sentence are similar, while the correlation between authors' rating and readers' rating when presented with the second sentence of each review is significantly lower. Moreover, when correlating readers' rating of whole reviews with readers' rating of single sentences, the correlation coefficient for last sentences is significantly higher than for second sentences.

As for the biometric measurements performed in the second experiment, since all subjects were computer-skilled, hesitation revealed through mouse-movements was assumed to be attributed to difficulty of decision-making rather than to problems in operating the mouse. As previously stated, we recorded mouse latency times following the reading of the texts up until clicking the mouse. Mouse latency times were not normalized for each subject due to the limited number of results. However, the average latency time is shorter for last sentences ($19.61 \pm 12.23s$) than for second sentences ($22.06 \pm 14.39s$). Indeed, the difference between latency times is not significant, as a paired t-test could not reject the null hypothesis that those distributions have equal means, but might show some tendency.

We also used the WizWhy software (Meidan,

2005) to perform combined analyses of readers' rating and response times. The analyses showed that when the difference between authors' and readers' ratings was $\leq |1|$ and the response time *much* shorter than average (<14.1 sec), then 96% of the sentences were *last* sentences. Due to the small sample size, we cautiously infer that last sentences express polarity better than second sentences, bearing in mind that the second sentence in our experiment represents any other sentence in the text except for the first one.

We also predicted that hesitation in making a decision would effect not only latency times but also mouse trajectories. Namely, hesitation will be accompanied by moving the mouse here and there, while decisiveness will show a firm movement. However, no such difference between the responses to last sentences or to second sentences appeared in our analysis; most subjects laid their hand still while reading the texts and while reflecting upon their answers. They moved the mouse only to rate the texts.

6 Conclusions and Future Work

In 2 psycholinguistic and psychophysical experiments, we showed that rating whole customer-reviews as compared to rating final sentences of these reviews showed an (expected) insignificant difference. In contrast, rating whole customer-reviews as compared to rating second sentences of these reviews, showed a considerable difference. Thus, instead of focusing on whole texts, computational linguists should focus on the last sentences for efficient and accurate automatic polarity-classification. Indeed, last but definitely not least!

We are currently running experiments that

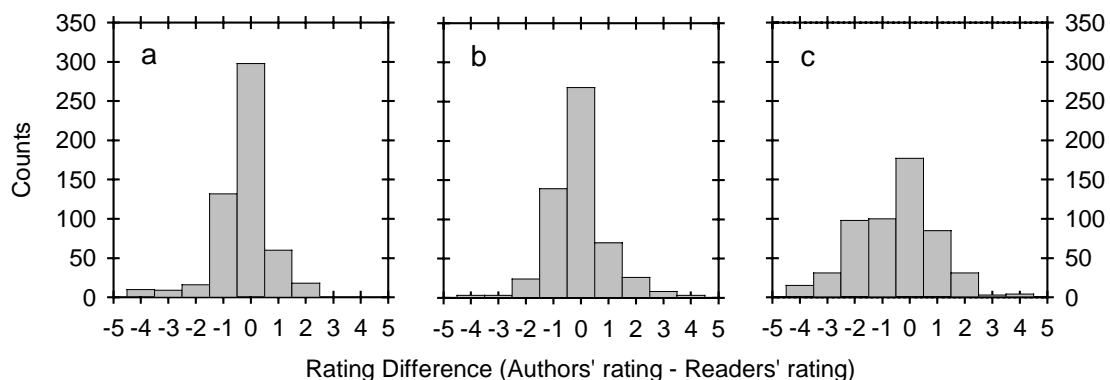


Figure 1. Histograms of the rating differences between the authors of reviews and their readers: for whole reviews (a), for last sentence only (b), and for second sentence only (c).

Readers' star rating of:	Correlated with:	Pearson Correlation Coefficient (P<0.0001)
Whole reviews	Authors' star rating	0.7891
Last sentences	of whole reviews	0.7616
Second sentences		0.4705
Last sentences	Readers' star rating	0.8463
Second sentences	of whole reviews	0.6563

Table 1. Pearson Correlation Coefficients

include hundreds of subjects in order to draw a profile of polarity evolution throughout customer reviews. Specifically, we present our subjects with sentences in various locations in customer reviews asking them to rate them. As the expanded experiment is *not* psychophysical, we added an additional *remote* radio button named “irrelevant” where subjects can judge a given text as lacking any evident polarity. Based on the rating results we will draw polarity profiles in order to see where, within customer reviews, polarity is best manifested and whether there are other “candidates” sentences that would serve as useful polarity indicators. The profiles will be used as a feature in our computational analysis.

Acknowledgments

We thank Prof. Rachel Giora and Prof. Ido Dagan for most valuable discussions, the 2 anonymous reviewers – for their excellent suggestions, and Thea Pagelson and Jason S. Henry - for their help with programming and running the psychophysical experiment.

References

- Baxendale, P. B. 1958. Machine-Made Index for Technical Literature - An Experiment. *IBM journal of research development* 2:263-311.
- Blitzer, John, Dredze, Mark, and Pereira, Fernando. 2007. Biographies, Bollywood, Boom-boxes and Blenders: Domain Adaptation for Sentiment Classification. Paper presented at *Association of Computational Linguistics (ACL)*.
- Giora, Rachel. 1983. Segmentation and Segment Cohesion: On the Thematic Organization of the Text. *Text* 3:155-182.
- Giora, Rachel. 1985. A Text-based Analysis of Non-narrative Texts. *Theoretical Linguistics* 12:115-135.
- Giora, Rachel. 1988. On the Informativeness Requirement. *Journal of Pragmatics* 12:547-565.
- Giora, Rachel, and Lee, Cher-Leng. 1996. Written Discourse Segmentation: The Function of Unstressed Pronouns in Mandarin Chinese. In *Refer-*

ence and Reference Accessibility ed. J. Gundel and T. Fretheim, 113-140. Amsterdam: Benjamins.

- Kieras, David E. 1978. Good and Bad Structure in Simple Paragraphs: Effects on Apparent Theme, Reading Time, and Recall. *Journal of Verbal Learning and Verbal Behavior* 17:13-28.
- Kieras, David E. 1980. Initial Mention as a Cue to the Main Idea and the Main Item of a Technical Passage. *Memory and Cognition* 8:345-353.
- Lin, Chen-Yew, and Hovy, Edward. 1997. Identifying Topic by Position. Paper presented at *Proceeding of the Fifth Conference on Applied Natural Language Processing*, San Francisco.
- Meidan, Abraham. 2005. Wizsoft's WizWhy. In *The Data Mining and Knowledge Discovery Handbook*, eds. Oded Maimon and Lior Rokach, 1365-1369: Springer.
- Murdock, B. B. Jr. 1962. The Serial Position Effect of Free Recall. *Journal of Experimental Psychology* 62:618-625.
- Yang, Changua, Lin, Kevin Hsin-Yih, and Chen, Hsin-Hsi. 2007a. Emotion Classification Using Web Blog Corpora. In *IEEE/WIC/ACM/ International Conference on Web Intelligence*. Silicon Valley, San Francisco.
- Yang, Changua, Lin, Kevin Hsin-Yih, and Chen, Hsin-Hsin. 2007b. Building Emotion Lexicon from Weblog Corpora. Paper presented at *Proceeding of the ACL 2007 Demo and Poster Session*, Prague.

Automatically generating annotator rationales to improve sentiment classification

Ainur Yessenalina Yejin Choi Claire Cardie

Department of Computer Science, Cornell University, Ithaca NY, 14853 USA

{ainur, ychoi, cardie}@cs.cornell.edu

Abstract

One of the central challenges in sentiment-based text categorization is that not every portion of a document is equally informative for inferring the overall sentiment of the document. Previous research has shown that enriching the sentiment labels with human annotators' "rationales" can produce substantial improvements in categorization performance (Zaidan et al., 2007). We explore methods to *automatically* generate annotator rationales for document-level sentiment classification. Rather unexpectedly, we find the automatically generated rationales just as helpful as human rationales.

1 Introduction

One of the central challenges in sentiment-based text categorization is that not every portion of a given document is equally informative for inferring its overall sentiment (e.g., Pang and Lee (2004)). Zaidan et al. (2007) address this problem by asking human annotators to mark (at least some of) the relevant text spans that *support each document-level sentiment decision*. The text spans of these "rationales" are then used to construct additional training examples that can guide the learning algorithm toward better categorization models.

But could we perhaps enjoy the performance gains of rationale-enhanced learning models without any additional human effort whatsoever (beyond the document-level sentiment label)? We hypothesize that in the area of sentiment analysis, where there has been a great deal of recent research attention given to various aspects of the task (Pang and Lee, 2008), this might be possible: using existing resources for sentiment analysis, we might be able to construct annotator rationales automatically.

In this paper, we explore a number of methods to automatically generate rationales for document-level sentiment classification. In particular, we investigate the use of off-the-shelf sentiment analysis components and lexicons for this purpose. Our approaches for generating annotator rationales can be viewed as *mostly unsupervised* in that we do not require manually annotated rationales for training.

Rather unexpectedly, our empirical results show that automatically generated rationales (91.78%) are just as good as human rationales (91.61%) for document-level sentiment classification of movie reviews. In addition, complementing the human annotator rationales with automatic rationales boosts the performance even further for this domain, achieving 92.5% accuracy. We further evaluate our rationale-generation approaches on product review data for which human rationales are not available: here we find that even randomly generated rationales can improve the classification accuracy although rationales generated from sentiment resources are not as effective as for movie reviews.

The rest of the paper is organized as follows. We first briefly summarize the SVM-based learning approach of Zaidan et al. (2007) that allows the incorporation of rationales (Section 2). We next introduce three methods for the automatic generation of rationales (Section 3). The experimental results are presented in Section 4, followed by related work (Section 5) and conclusions (Section 6).

2 Contrastive Learning with SVMs

Zaidan et al. (2007) first introduced the notion of *annotator rationales* — text spans highlighted by human annotators as support or evidence for each document-level sentiment decision. These rationales, of course, are only useful if the sentiment categorization algorithm can be extended to exploit the rationales effectively. With this in mind, Zaidan et al. (2007) propose the following *con-*

trastive learning extension to the standard SVM learning algorithm.

Let \vec{x}_i be movie review i , and let $\{\vec{r}_{ij}\}$ be the set of *annotator rationales* that support the positive or negative sentiment decision for \vec{x}_i . For each such rationale \vec{r}_{ij} in the set, construct a *contrastive training example* \vec{v}_{ij} , by removing the text span associated with the rationale \vec{r}_{ij} from the original review \vec{x}_i . Intuitively, the contrastive example \vec{v}_{ij} should not be as informative to the learning algorithm as the original review \vec{x}_i , since one of the supporting regions identified by the human annotator has been deleted. That is, the *correct* learned model should be *less confident* of its classification of a contrastive example vs. the corresponding original example, and the classification boundary of the model should be modified accordingly.

Zaidan et al. (2007) formulate exactly this intuition as SVM constraints as follows:

$$(\forall i, j) : y_i (\vec{w}\vec{x}_i - \vec{w}\vec{v}_{ij}) \geq \mu(1 - \xi_{ij})$$

where $y_i \in \{-1, +1\}$ is the negative/positive sentiment label of document i , \vec{w} is the weight vector, $\mu \geq 0$ controls the size of the margin between the original examples and the contrastive examples, and ξ_{ij} are the associated slack variables. After some re-writing of the equations, the resulting objective function and constraints for the SVM are as follows:

$$\frac{1}{2} \|\vec{w}\|^2 + C \sum_i \xi_i + C_{contrast} \sum_{ij} \xi_{ij} \quad (1)$$

subject to constraints:

$$(\forall i) : y_i \vec{w} \cdot \vec{x}_i \geq 1 - \xi_i, \quad \xi_i \geq 0$$

$$(\forall i, j) : y_i \vec{w} \cdot \vec{x}_{ij} \geq 1 - \xi_{ij} \quad \xi_{ij} \geq 0$$

where ξ_i and ξ_{ij} are the slack variables for \vec{x}_i (the original examples) and \vec{x}_{ij} (\vec{x}_{ij} are named as *pseudo examples* and defined as $\vec{x}_{ij} = \frac{\vec{x}_i - \vec{v}_{ij}}{\mu}$), respectively. Intuitively, the pseudo examples (\vec{x}_{ij}) represent the difference between the original examples (\vec{x}_i) and the contrastive examples (\vec{v}_{ij}), weighted by a parameter μ . C and $C_{contrast}$ are parameters to control the trade-offs between training errors and margins for the original examples \vec{x}_i and pseudo examples \vec{x}_{ij} respectively. As noted in Zaidan et al. (2007), $C_{contrast}$ values are generally smaller than C for noisy rationales.

In the work described below, we similarly employ Zaidan et al.’s (2007) contrastive learning method to incorporate rationales for document-level sentiment categorization.

3 Automatically Generating Rationales

Our goal in the current work, is to generate annotator rationales automatically. For this, we rely on the following two assumptions:

- (1) Regions marked as annotator rationales are more subjective than unmarked regions.
- (2) The sentiment of each annotator rationale coincides with the document-level sentiment.

Note that assumption 1 was not observed in the Zaidan et al. (2007) work: annotators were asked only to mark a few rationales, leaving other (also subjective) rationale sections unmarked.

And at first glance, assumption (2) might seem too obvious. But it is important to include as there can be subjective regions with seemingly conflicting sentiment in the same document (Pang et al., 2002). For instance, an author for a movie review might express a positive sentiment toward the movie, while also discussing a negative sentiment toward one of the fictional characters appearing in the movie. This implies that not all subjective regions will be relevant for the document-level sentiment classification — rather only those regions whose polarity matches that of the document should be considered.

In order to extract regions that satisfy the above assumptions, we first look for subjective regions in each document, then filter out those regions that exhibit a sentiment value (i.e., polarity) that conflicts with polarity of the document. Assumption 2 is important as there can be subjective regions with seemingly conflicting sentiment in the same document (Pang et al., 2002).

Because our ultimate goal is to reduce human annotation effort as much as possible, we do not employ supervised learning methods to directly learn to identify good rationales from human-annotated rationales. Instead, we opt for methods that make use of only the document-level sentiment and off-the-shelf utilities that were trained for slightly different sentiment classification tasks using a corpus from a different domain and of a different genre. Although such utilities might not be optimal for our task, we hoped that these basic resources from the research community would constitute an adequate source of sentiment information for our purposes.

We next describe three methods for the automatic acquisition of rationales.

3.1 Contextual Polarity Classification

The first approach employs OpinionFinder (Wilson et al., 2005a), an off-the-shelf opinion analysis utility.¹ In particular, OpinionFinder identifies phrases expressing positive or negative opinions. Because OpinionFinder models the task as a word-based classification problem rather than a sequence tagging task, most of the identified opinion phrases consist of a single word. In general, such short text spans cannot fully incorporate the contextual information relevant to the detection of subjective language (Wilson et al., 2005a). Therefore, we conjecture that good rationales should extend beyond short phrases.² For simplicity, we choose to extend OpinionFinder phrases to sentence boundaries.

In addition, to be consistent with our second operating assumption, we keep only those sentences whose polarity coincides with the document-level polarity. In sentences where OpinionFinder marks multiple opinion words with opposite polarities we perform a simple voting — if words with positive (or negative) polarity dominate, then we consider the entire sentence as positive (or negative). We ignore sentences with a tie. Each selected sentence is considered as a separate rationale.

3.2 Polarity Lexicons

Unfortunately, domain shift as well as task mismatch could be a problem with any opinion utility based on supervised learning.³ Therefore, we next consider an approach that does not rely on supervised learning techniques but instead explores the use of a manually constructed polarity lexicon. In particular, we use the lexicon constructed for Wilson et al. (2005b), which contains about 8000 words. Each entry is assigned one of three polarity values: positive, negative, neutral. We construct rationales from the polarity lexicon for every instance of positive and negative words in the lexicon that appear in the training corpus.

As in the OpinionFinder rationales, we extend the words found by the PolarityLexicon approach to sentence boundaries to incorporate potentially

¹Available at www.cs.pitt.edu/mpqa/opinionfinderrelease/.

²This conjecture is indirectly confirmed by the fact that human-annotated rationales are rarely a single word.

³It is worthwhile to note that OpinionFinder is trained on a newswire corpus whose prevailing sentiment is known to be negative (Wiebe et al., 2005). Furthermore, OpinionFinder is trained for a task (word-level sentiment classification) that is different from marking annotator rationales (sequence tagging or text segmentation).

relevant contextual information. We retain as rationales only those sentences whose polarity coincides with the document-level polarity as determined via the voting scheme of Section 3.1.

3.3 Random Selection

Finally, we generate annotator rationales randomly, selecting 25% of the sentences from each document⁴ and treating each as a separate rationale.

3.4 Comparison of Automatic vs. Human-annotated Rationales

Before evaluating the performance of the automatically generated rationales, we summarize in Table 1 the differences between automatic vs. human-generated rationales. All computations were performed on the same movie review dataset of Pang and Lee (2004) used in Zaidan et al. (2007). Note, that the Zaidan et al. (2007) annotation guidelines did not insist that annotators mark **all** rationales, only that some were marked for each document. Nevertheless, we report precision, recall, and F-score based on overlap with the human-annotated rationales of Zaidan et al. (2007), so as to demonstrate the degree to which the proposed approaches align with human intuition. Overlap measures were also employed by Zaidan et al. (2007).

As shown in Table 1, the annotator rationales found by OpinionFinder (F-score 49.5%) and the PolarityLexicon approach (F-score 52.6%) match the human rationales much better than those found by random selection (F-score 27.3%).

As expected, OpinionFinder’s positive rationales match the human rationales at a significantly lower level (F-score 31.9%) than negative rationales (59.5%). This is due to the fact that OpinionFinder is trained on a dataset biased toward negative sentiment (see Section 3.1 - 3.2). In contrast, all other approaches show a balanced performance for positive and negative rationales vs. human rationales.

4 Experiments

For our contrastive learning experiments we use *SVM^{light}* (Joachims, 1999). We evaluate the usefulness of automatically generated rationales on

⁴We chose the value of 25% to match the percentage of sentences per document, on average, that contain human-annotated rationales in our dataset (24.7%).

Method	% of sentences selected	Precision			Recall			F-Score		
		ALL	POS	NEG	ALL	POS	NEG	ALL	POS	NEG
OPINIONFINDER	22.8%	54.9	56.1	54.6	45.1	22.3	65.3	49.5	31.9	59.5
POLARITYLEXICON	38.7%	45.2	42.7	48.5	63.0	71.8	55.0	52.6	53.5	51.6
RANDOM	25.0%	28.9	26.0	31.8	25.9	24.9	26.7	27.3	25.5	29.0

Table 1: Comparison of Automatic vs. Human-annotated Rationales.

five different datasets. The first is the movie review data of Pang and Lee (2004), which was manually annotated with rationales by Zaidan et al. (2007)⁵; the remaining are four product review datasets from Blitzer et al. (2007).⁶ Only the movie review dataset contains human annotator rationales. We replicate the same feature set and experimental set-up as in Zaidan et al. (2007) to facilitate comparison with their work.⁷

The contrastive learning method introduced in Zaidan et al. (2007) requires three parameters: (C , μ , $C_{contrast}$). To set the parameters, we use a grid search with step 0.1 for the range of values of each parameter around the point (1,1,1). In total, we try around 3000 different parameter triplets for each type of rationales.

4.1 Experiments with the Movie Review Data

We follow Zaidan et al. (2007) for the training/test data splits. The top half of Table 2 shows the performance of a system trained with **no annotator rationales** vs. two variations of human annotator rationales. HUMANR treats each rationale in the same way as Zaidan et al. (2007). HUMANR@SENTENCE extends the human annotator rationales to sentence boundaries, and then treats each such sentence as a separate rationale. As shown in Table 2, we get almost the same performance from these two variations (91.33% and 91.61%).⁸ This result demonstrates that locking rationales to sentence boundaries was a reasonable

⁵Available at <http://www.cs.jhu.edu/~ozaidan/rationales/>.

⁶<http://www.cs.jhu.edu/~mdredze/datasets/sentiment/>.

⁷We use binary unigram features corresponding to the unstemmed words or punctuation marks with count greater or equal to 4 in the full 2000 documents, then we normalize the examples to the unit length. When computing the pseudo examples $\vec{x}_{ij} = \frac{\vec{x}_i - \vec{v}_{ij}}{\mu}$ we first compute $(\vec{x}_i - \vec{v}_{ij})$ using the binary representation. As a result, features (unigrams) that appeared in both vectors will be zeroed out in the resulting vector. We then normalize the resulting vector to a unit vector.

⁸The performance of HUMANR reported by Zaidan et al. (2007) is 92.2% which lies between the performance we get (91.61%) and the oracle accuracy we get if we knew the best parameters for the test set (92.67%).

Method	Accuracy
NORATIONALES	88.56
HUMANR	91.61 [•]
HUMANR@SENTENCE	91.33 ^{•†}
OPINIONFINDER	91.78 ^{•†}
POLARITYLEXICON	91.39 ^{•†}
RANDOM	90.00 [*]
OPINIONFINDER+HUMANR@SENTENCE	92.50 ^{•△}

Table 2: Experimental results for the movie review data.

- The numbers marked with [•] (or ^{*}) are statistically significantly better than NORATIONALES according to a paired t-test with $p < 0.001$ (or $p < 0.01$).
- The numbers marked with [△] are statistically significantly better than HUMANR according to a paired t-test with $p < 0.01$.
- The numbers marked with [†] are *not* statistically significantly worse than HUMANR according to a paired t-test with $p > 0.1$.

choice.

Among the approaches that make use of only automatic rationales (bottom half of Table 2), the best is OPINIONFINDER, reaching 91.78% accuracy. This result is slightly better than results exploiting human rationales (91.33-91.61%), although the difference is not statistically significant. This result demonstrates that automatically generated rationales are just as good as human rationales in improving document-level sentiment classification. Similarly strong results are obtained from the POLARITYLEXICON as well.

Rather unexpectedly, RANDOM also achieves statistically significant improvement over NORATIONALES (90.0% vs. 88.56%). However, notice that the performance of RANDOM is statistically significantly lower than those based on human rationales (91.33-91.61%).

In our experiments so far, we observed that some of the automatic rationales are just as good as human rationales in improving the document-level sentiment classification. Could we perhaps achieve an even better result if we combine the automatic rationales with human

rationales? The answer is yes! The accuracy of OPINIONFINDER+HUMANR@SENTENCE reaches 92.50%, which is statistically significantly better than HUMANR (91.61%). In other words, not only can our automatically generated rationales replace human rationales, but they can also improve upon human rationales when they are available.

4.2 Experiments with the Product Reviews

We next evaluate our approaches on datasets for which human annotator rationales do not exist. For this, we use some of the product review data from Blitzer et al. (2007): reviews for Books, DVDs, Videos and Kitchen appliances. Each dataset contains 1000 positive and 1000 negative reviews. The reviews, however, are substantially shorter than those in the movie review dataset: the average number of sentences in each review is 9.20/9.13/8.12/6.37 respectively vs. 30.86 for the movie reviews. We perform 10-fold cross-validation, where 8 folds are used for training, 1 fold for tuning parameters, and 1 fold for testing.

Table 3 shows the results. Rationale-based methods perform statistically significantly better than NORATIONALES for all but the Kitchen dataset. An interesting trend in product review datasets is that RANDOM rationales are just as good as other more sophisticated rationales. We suspect that this is because product reviews are generally shorter and more focused than the movie reviews, thereby any randomly selected sentence is likely to be a good rationale. Quantitatively, subjective sentences in the product reviews amount to 78% (McDonald et al., 2007), while subjective sentences in the movie review dataset are only about 25% (Mao and Lebanon, 2006).

4.3 Examples of Annotator Rationales

In this section, we examine an example to compare the automatically generated rationales (using OPINIONFINDER) with human annotator rationales for the movie review data. In the following positive document snippet, automatic rationales are underlined, while **human-annotated rationales** are in bold face.

...But a little niceness goes a long way these days, and **there’s no denying the entertainment value** of that thing you do! **It’s just about impossible to hate.** It’s an inoffensive, enjoyable piece of nostalgia that is sure to leave audiences smiling and humming, if not singing, “that thing you do!” –quite possibly for days...

Method	Books	DVDs	Videos	Kitchen
NORATIONALES	80.20	80.95	82.40	87.40
OPINIONFINDER	81.65*	82.35*	84.00*	88.40
POLARITYLEXICON	82.75•	82.85•	84.55•	87.90
RANDOM	82.05•	82.10•	84.15•	88.00

Table 3: Experimental results for subset of Product Review data

– The numbers marked with • (or *) are statistically significantly better than NORATIONALES according to a paired t-test with $p < 0.05$ (or $p < 0.08$).

Notice that, although OPINIONFINDER misses some human rationales, it avoids the inclusion of “impossible to hate”, which contains only negative terms and is likely to be confusing for the contrastive learner.

5 Related Work

In broad terms, constructing annotator rationales automatically and using them to formulate contrastive examples can be viewed as learning with prior knowledge (e.g., Schapire et al. (2002), Wu and Srihari (2004)). In our task, the prior knowledge corresponds to our operating assumptions given in Section 3. Those assumptions can be loosely connected to recognizing and exploiting discourse structure (e.g., Pang and Lee (2004), Taboada et al. (2009)). Our automatically generated rationales can be potentially combined with other learning frameworks that can exploit annotator rationales, such as Zaidan and Eisner (2008).

6 Conclusions

In this paper, we explore methods to automatically generate annotator rationales for document-level sentiment classification. Our study is motivated by the desire to retain the performance gains of rationale-enhanced learning models while eliminating the need for additional human annotation effort. By employing existing resources for sentiment analysis, we can create automatic annotator rationales that are as good as human annotator rationales in improving document-level sentiment classification.

Acknowledgments

We thank anonymous reviewers for their comments. This work was supported in part by National Science Foundation Grants BCS-0904822, BCS-0624277, IIS-0535099 and by the Department of Homeland Security under ONR Grant N0014-07-1-0152.

References

- John Blitzer, Mark Dredze, and Fernando Pereira. 2007. Biographies, bollywood, boom-boxes and blenders: Domain adaptation for sentiment classification. In *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, pages 440–447, Prague, Czech Republic, June. Association for Computational Linguistics.
- Thorsten Joachims. 1999. Making large-scale support vector machine learning practical. pages 169–184.
- Yi Mao and Guy Lebanon. 2006. Sequential models for sentiment prediction. In *Proceedings of the ICML Workshop: Learning in Structured Output Spaces Open Problems in Statistical Relational Learning Statistical Network Analysis: Models, Issues and New Directions*.
- Ryan McDonald, Kerry Hannan, Tyler Neylon, Mike Wells, and Jeff Reynar. 2007. Structured models for fine-to-coarse sentiment analysis. In *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, pages 432–439, Prague, Czech Republic, June. Association for Computational Linguistics.
- Bo Pang and Lillian Lee. 2004. A sentimental education: sentiment analysis using subjectivity summarization based on minimum cuts. In *ACL '04: Proceedings of the 42nd Annual Meeting on Association for Computational Linguistics*, page 271, Morristown, NJ, USA. Association for Computational Linguistics.
- Bo Pang and Lillian Lee. 2008. Opinion mining and sentiment analysis. *Found. Trends Inf. Retr.*, 2(1-2):1–135.
- Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. 2002. Thumbs up?: sentiment classification using machine learning techniques. In *EMNLP '02: Proceedings of the ACL-02 conference on Empirical methods in natural language processing*, pages 79–86, Morristown, NJ, USA. Association for Computational Linguistics.
- Robert E. Schapire, Marie Rochery, Mazin G. Rahim, and Narendra Gupta. 2002. Incorporating prior knowledge into boosting. In *ICML '02: Proceedings of the Nineteenth International Conference on Machine Learning*, pages 538–545, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.
- Maite Taboada, Julian Brooke, and Manfred Stede. 2009. Genre-based paragraph classification for sentiment analysis. In *Proceedings of the SIGDIAL 2009 Conference*, pages 62–70, London, UK, September. Association for Computational Linguistics.
- Janyce Wiebe, Theresa Wilson, and Claire Cardie. 2005. Annotating expressions of opinions and emotions in language. *Language Resources and Evaluation*, 1(2):0.
- Theresa Wilson, Paul Hoffmann, Swapna Somasundaran, Jason Kessler, Janyce Wiebe, Yejin Choi, Claire Cardie, Ellen Riloff, and Siddharth Patwardhan. 2005a. Opinionfinder: a system for subjectivity analysis. In *Proceedings of HLT/EMNLP on Interactive Demonstrations*, pages 34–35, Morristown, NJ, USA. Association for Computational Linguistics.
- Theresa Wilson, Janyce Wiebe, and Paul Hoffmann. 2005b. Recognizing contextual polarity in phrase-level sentiment analysis. In *HLT-EMNLP '05: Proceedings of the conference on Human Language Technology and Empirical Methods in Natural Language Processing*, pages 347–354, Morristown, NJ, USA. Association for Computational Linguistics.
- Xiaoyun Wu and Rohini Srihari. 2004. Incorporating prior knowledge with weighted margin support vector machines. In *KDD '04: Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 326–333, New York, NY, USA. ACM.
- Omar F. Zaidan and Jason Eisner. 2008. Modeling annotators: a generative approach to learning from annotator rationales. In *EMNLP '08: Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 31–40, Morristown, NJ, USA. Association for Computational Linguistics.
- Omar F. Zaidan, Jason Eisner, and Christine Piatko. 2007. Using “annotator rationales” to improve machine learning for text categorization. In *NAACL HLT 2007: Proceedings of the Main Conference*, pages 260–267, April.

Simultaneous Tokenization and Part-of-Speech Tagging for Arabic without a Morphological Analyzer

Seth Kulick

Linguistic Data Consortium
University of Pennsylvania
skulick@seas.upenn.edu

Abstract

We describe an approach to simultaneous tokenization and part-of-speech tagging that is based on separating the closed and open-class items, and focusing on the likelihood of the possible stems of the open-class words. By encoding some basic linguistic information, the machine learning task is simplified, while achieving state-of-the-art tokenization results and competitive POS results, although with a reduced tag set and some evaluation difficulties.

1 Introduction

Research on the problem of morphological disambiguation of Arabic has noted that techniques developed for lexical disambiguation in English do not easily transfer over, since the affixation present in Arabic creates a very different tag set than for English, in terms of the number and complexity of tags. In addition to inflectional morphology, the POS tags encode more complex tokenization sequences, such as *preposition + noun* or *noun + possessive pronoun*.

One approach taken to this problem is to use a morphological analyzer such as BAMA-v2.0 (Buckwalter, 2004) or SAMA-v3.1 (Maamouri et al., 2009c)¹, which generates a list of all possible morphological analyses for a given token. Machine learning approaches can model separate aspects of a solution (e.g., “has a pronominal clitic”) and then combine them to select the most appropriate solution from among this list. A benefit of this approach is that by picking a single solution from the morphological analyzer, the part-of-speech and tokenization comes as a unit (Habash and Rambow, 2005; Roth et al., 2008).

¹SAMA-v3.1 is an updated version of BAMA, with many significant differences in analysis.

In contrast, other approaches have used a pipelined approach, with separate models to first do tokenization and then part-of-speech tagging (Diab et al., 2007; Diab, 2009). While these approaches have somewhat lower performance than the joint approach, they have the advantage that they do not rely on the presence of a full-blown morphological analyzer, which may not always be available or appropriate as the data shifts to different genres or Arabic dialects.

In this work we present a novel approach to this problem that allows us to do simultaneous tokenization and core part-of-speech tagging with a simple classifier, without using a full-blown morphological analyzer. We distinguish between closed-class and open-class categories of words, and encode regular expressions that express the morphological patterns for the former, and simple regular expressions for the latter that provide only the generic templates for affixation. We find that a simple baseline for the closed-class words already works very well, and for the open-class words we classify only the possible stems for all such expressions. This is however sufficient for tokenization and core POS tagging, since the stem identifies the appropriate regular expression, which then in turn makes explicit, simultaneously, the tokenization and part-of-speech information.

2 Background

The Arabic Treebank (ATB) contains a full morphological analysis of each “source token”, a whitespace/punctuation-delimited string from the source text. The SAMA analysis includes four fields, as shown in the first part of Table 1.² TEXT is the actual source token text, to be analyzed. VOC is the vocalized form, including diacritics. Each VOC segment has associated with it a POS tag and

²This is the analysis for one particular instance of *ktbh*. The same source token may receive another analysis elsewhere in the treebank.

ATB analysis for one source token:

TEXT:	ktbh		
VOC:	kutub	u	hu
POS:	NOUN	CASE_NOM	POSS_PRON_3MS
GLOSS:	books	[def.nom.]	its/his
Results in two ATB tree tokens:			
TEXT:	ktb	h	
VOC:	kutub+u	hu	
POS:	NOUN+CASE_NOM	POSS_PRON_3MS	
Current work recovers:			
TEXT:	ktb	h	
POS:	NOA	POSS_PRON	

Table 1: Example analysis of one source token

NOUN, ADJ, NOUN.VN, ADJ.VN, NOUN_NUM, ADJ_NUM, NOUN_QUANT, ADJ_COMP, ABBREV	NOA (Noun or Adjective)
IV, IV_PASS	IV
PV, PV_PASS	PV
IVSUFF_DO, PVSUFF_DO	OBJ_PRON

Table 2: Collapsing of ATB core tags into reduced core tags

GLOSS. While “tokenization” can be done in different ways on top of this analysis, the ATB splits the VOC/POS/GLOSS segments up based on the POS tags to form the “tree tokens” necessary for treebanking. As shown in the second part of Table 1, the first two segments remain together as one tree token, and the pronoun is separated as a separate tree token. In addition, the input TEXT is separated among the two tree tokens.³

Each tree token’s POS tag therefore consists of what can be considered an “ATB core tag”, together with inflectional material (case, gender, number). For example, in Table 1, the “core tag” of the first tree token is NOUN. In this work, we aim to recover the separation of a source token TEXT into the corresponding separate tree token TEXTs, together with a “reduced core tag” for each tree token. By “reduced core tag”, we mean an ATB core tag that has been reduced in two ways:

(1) All inflectional material [infl] is stripped off six ATB core tags: PRON[infl], POSS_PRON[infl], DEM[infl], [IV|PV|CV]SUFF_DO[infl]

(2) Collapsing of some ATB core tags, as listed in Table 2.

These two steps result in a total of 40 reduced core tags, and each tree token has exactly one such reduced core tag. We work with the ATB3-v3.2 release of the ATB (Maamouri et al., 2009b), which

³See (Kulick et al., 2010) for a detailed discussion of how this splitting is done and how the tree token TEXT field (called INPUT STRING in the ATB releases) is created.

NOA	173938	PART	288
PREP	49894	RESTRIC_PART	237
PUNC	41398	DET	215
NOUN_PROP	29423	RC_PART	192
CONJ	28257	FOCUS_PART	191
PV	16669	TYPO	188
IV	15361	INTERROG_PART	187
POSS_PRON	9830	INTERROG_ADV	169
SUB_CONJ	8200	INTERROG_PRON	112
PRON	6995	CV	106
REL_PRON	5647	VOC_PART	74
DEM	3673	VERB	62
OBJ_PRON	2812	JUS_PART	56
NEG_PART	2649	FOREIGN	46
PSEUDO_VERB	1505	DIALECT	41
FUT_PART	1099	INTERJ	37
ADV	1058	EMPHATIC_PART	19
VERB_PART	824	CVSUFF_DO	15
REL_ADV	414	GRAMMAR_PROB	4
CONNEC_PART	405	LATIN	1

Table 3: The 40 “reduced core tags”, and their frequencies in ATB3-v3.2. The total count is 402291, which is the number of tree tokens in ATB3-v3.2.

has 339710 source tokens and 402291 tree tokens, where the latter are derived from the former as discussed above. Table 3 lists the 40 reduced tags we use, and their frequency among the ATB3-v3.2 tree tokens.

3 Description of Approach

Given a source token, we wish to recover (1) the tree tokens (which amounts to recovering the ATB tokenization), and (2) the reduced core POS tag for each tree token. For example, in Table 1, given the input source token TEXT *ktbh*, we wish to recover the tree tokens *ktb/NOA* and *h/POSS_PRON*.

As mentioned in the introduction, we use regular expressions that encode all the tokenization and POS tag possibilities. Each “group” (substring unit) in a regular expression (regex) is assigned an internal name, and a list is maintained of the possible reduced core POS tags that can occur with that regex group. It is possible, and indeed usually the case for groups representing affixes, that more than one such POS tag is possible. However, it is crucial for our approach that while some given source token TEXT may match many regular expressions (regexes), when the POS tag is also taken into account, there can be only one match among all the (open or closed-class) regexes. We say a source token “pos-matches” a regex if the TEXT matches and POS tags match, and “text-matches” if the TEXT matches the regex regardless of the POS. During training, the pos-matching

(REGEX #1) [w|f]lm
 w: [PART, CONJ, SUB_CONJ, PREP]
 f: [CONJ, SUB_CONJ, CONNEC_PART, RC_PART]
 lm: [NEG_PART]
 (REGEX #2) [w|f]lm
 w: and f: same as above
 lm: [REL_ADV, INTERROG_ADV]

Figure 1: Two sample closed-class regexes

regex for a source token TEXT is stored as the gold solution for closed-class patterns, or used to create the gold label for the open-class classifier.

We consider the open-class tags in Table 3 to be: NOUN_PROP, NOA, IV, CV, PV, VERB. A source token is considered to have an open-class solution if any of the tree tokens in that solution have an open-class tag. For example, *ktbh* in Table 1 has an open-class solution because one of the tree tokens has an open-class tag (NOA), even though the other is closed-class (POSS_PRON).

We encode the possible solutions for closed-class source tokens using the lists in the ATB morphological guidelines (Maamouri et al., 2009a). For example, Figure 1 shows two of the closed-class regexes. The text *wlm* can text-match either REGEX #1 or #2, but when the POS tag for *lm* is taken into account, only one can pos-match. We return to the closed-class regexes in Section 4.

We also encode regular expression for the open-class source tokens, but these are simply generic templates expressing the usual affix possibilities, such as:

```
[wf] [blk] stem.NOA poss.pronoun
```

where there is no list of possible strings for `stem_NOA`, but which instead can match anything. While all parts except for the stem are optional, we do not make such parts optional in a single expression. Instead, we multiple out the possibilities into different expressions with different parts (e.g., [wf]) being obligatory). The reason for this is that we give different names to the stem in each case, and this is the basis of the features for the classifier. As with the closed-class regexes, we associate a list of possible POS tags for each named group within a regular expression. Here the `stem_NOA` group can only have the tag NOA.

We create features for a classifier for the open-class words as follows. Each word is run through all of the open-class regular expressions. For each expression that text-matches, we make a feature

which is the name of the stem part of the regular expression, along with the characters that match the stem. The stem name encodes whether there is a prefix or suffix, but does not include a POS tag. However, the source token pos-matches exactly one of the regular expressions, and the pos tag for the stem is appended to the named stem for that expression to form the gold label for training and the target for testing.

For example, Table 4 lists the matching regular expression for three words. The first, *yjry*, text-matches the generic regular expressions for *any string/NOA*, *any string/IV*, etc. These are summarized in one listing, *yjry/all*. The name of the stem for all these expressions is the same, just *stem*, and so they all give rise to the same feature, `stem=yjry`. It also matches the expression for a NOA with a possessive pronoun⁴, and in this case the stem name in the regular expression is `stem_spp` (which stands for “stem with a possessive pronoun suffix”), and this gives rise to the feature `stem_spp=yjry`. Similarly, for *wAfAdt* the stem of the second expression has the name `p_stem`, for a prefix. The third example shows the different stem names that occur when there are both prefix and suffix possibilities. For each example, there is exactly one regex that not only text-matches, but also pos-matches. The combination of the stem name in these cases together with the gold tag forms the gold label, as indicated in column 3.

Therefore, for each source token TEXT, the features include the ones arising from the named stems of all the regexes that text-match that TEXT, as shown in column 4, and the gold label is the appropriate stem name together with the POS tag, as shown in column 3. We also include some typical features for each stem, such as first and last two letters of each stem, etc. For example, *wAfAdt* would also have the features `stem_fl=w`, `p_stem_fl=A`, indicating that the first letter of `stem` is *w* and the first letter of `p_stem` is *A*. We also extract a list of proper nouns from SAMA-v3.1 as a temporary proxy for a named entity list, and include a feature for a stem if that stem is in the list (`stem_in_list`, `p_stem_in_list`, etc.)

We do not model separate classifiers for prefix possibilities. There is a dependency between the

⁴The regex listed is slightly simplified. It actually contains a reference to the list of all possessive pronouns, not just *y*.

source TEXT	text-matching regular expressions	gold label	feature
yjry	yjry/all (<i>happens</i>)	stem:IV	stem=yjry
	yjr/NOA+y/POSS_PRON		stem_spp=yjr
wAfAdt	wAfAdt/all		stem=wAfAdt
	w + AfAdt/all (<i>and+reported</i>)	p_stem:PV	p_stem=AfAdt
lAstyDAHhm	lAstyDAHhm/all		stem=lAstyDAHhm
	l/PREP + AstyDAHhm/NOA		p_stem=AstyDAHhm
	l/PREP + AstyDAH/NOA + hm/POSS_PRON <i>for + request for clarification + their</i>	p_stem_spp:NOA	p_stem_spp=AstyDAH
	lAstyDAH/NOA + hm/POSS_PRON		stem_spp=lAstyDAH
	lAstyDAH/IV,PV,CV + hm/OBJ_PRON		stem_svop=lAstyDAH
	l/PREP,JUS_PART + AstyDAH/IV,PV,CV + hm/OBJ_PRON		p_stem_svop=AstyDAH

Table 4: Example features and gold labels for three words. Each text-matching regex gives rise to one feature shown in column 4, based on the stem of that regular expression. A *p_* before a stem means that it has a prefix, *_spp* after means that it has a possessive pronouns suffix, and *_svop* means that it has a (verbal) object pronoun suffix. “all” in the matching regular expression is shorthand for text-matching all the corresponding regular expressions with NOA, IV, etc. For each word, exactly one regex also pos-matches, which results in the gold label, shown in column 3.

possibility of a prefix and the likelihood of the remaining stem, and so we focus on the likelihood of the possible stems, where the open-class regexes enumerate the possible stems. A gold label together with the source token TEXT maps back to a single regex, and so for a given label, the TEXT is parsed by that regular expression, resulting in a tokenization along with list of possible POS tags for each affix group in the regex.⁵

During training and testing, we run each word through all the open and closed regexes. Text-matches for an open-class regex give rise to features as just described. Also, if the word matches any closed-class regex, it receives the feature `MATCHES_CLOSED`. During training, if the correct match for the word is one of the closed-class expressions, then the gold label is `CLOSED`. The classifier is used only to get solutions for the open-class words, although we wish to give the classifier all the words for the sentence. The cross-product of the stem name and (open-class) reduced core POS tags, plus the `CLOSED` tag, yields 24 labels for a CRF classifier in Mallet (McCallum, 2002).

4 Experiments and Evaluation

We worked with ATB3-v3.2, following the training/devtest split in (Roth et al., 2008) on a previous release of the same data. We keep a listing (List #1) of all (source token TEXT, solution) pairs seen during training. For an open-class solution, “solution” is the gold label as described in

⁵In Section 4 we discuss how these are narrowed down to one POS tag.

Section 3. For a closed-class solution, “solution” is the name of the single pos-matching regex. In addition, for every regex seen during training that pos-matches some source token TEXT, we keep a listing (List #2) of all ((regex-group-name, text), POS-tag) tuples. We use the information in List #1 to choose a solution for all words seen in training in the Baseline and Run 2 below, and in Run 3, for words text-matching a closed-class expression. We use List #2 to disambiguate all remaining cases of POS ambiguity, wherever a solution comes from.

For example, if *wlm* is seen during testing, List #1 will be consulted to find the most common solution (REGEX #1 or #2), and in either case, List #2 will be consulted to determine the most frequent tag for *w* as a prefix. While there is certainly room for improvement here, this works quite well since the tags for the affixes do not vary much.

We score the solution for a source token instance as correct for tokenization if it exactly matches the TEXT split for the tree tokens derived from that source token instance in the ATB. It is correct for POS if correct for tokenization and if each tree token has the same POS tag as the reduced core tag for that tree token in the ATB.

For a simple baseline, if a source token TEXT is in List #1 then we simply use the most frequent stored solution. Otherwise we run the TEXT through all the regexes. If it text-matches any closed-class expression, we pick a random choice from among those regexes and otherwise from the open-class regexes that it text-matches. Any POS ambiguities for a regex group are disambiguated

Solution Origin	Baseline			Run 2			Run 3		
	# tokens	Tok	POS	# tokens	Tok	POS	# tokens	Tok	POS
All	51664	96.0%	87.4%	51664	99.4%	95.1%	51664	99.3%	95.1%
Stored	46072	99.8%	96.6%	46072	99.8%	96.6%	16145	99.6%	96.4%
Open	5565	64.6%	11.6%	10	10.0%	0.0%	11	54.5%	0.0%
Closed	27	81.5%	59.3%	27	81.5%	63.0%	27	81.5%	63.0%
Mallet	0			5555	96.0%	83.8%	35481	99.1%	94.5%

Table 5: Results for Baseline and two runs. Origin “stored” means that the appropriate regex came from the list stored during training. Origins “open” and “closed” are random choices from the open or closed regexes for the source token. “Mallet” means that it comes from the label output by the CRF classifier.

using List #2, as discussed above. The results are shown in Table 5. The score is very high for the words seen during training, but much lower for open-class words that were not. As expected, almost all (except 27) instances of closed-class words were seen during training.

For run 2, we continue to use the stored solution if the token was seen in training. If not, then if the TEXT matches one or more closed-class regexes, we randomly choose one. Otherwise, if the CRF classifier has produced an open-class match for that token, we use that (and otherwise, in only 10 cases, use a random open-class match). There is a significant improvement in the score for the open-class items, and therefore in the overall results.

For run 3, we put more of a burden on the classifier. If a word matches any closed-class expression, we either use the most frequent occurrence during training (if it was seen), or use a random matching closed-class expression (if not). If the word doesn’t match a closed-class expression, we use the mallet result. The mallet score goes up, almost certainly because the score is now including results on words that were seen during training. The overall POS result for run 3 is slightly less than run 2. (95.099% compared to 95.147%).

It is not a simple matter to compare results with previous work, due to differing evaluation techniques, data sets, and POS tag sets. With different data sets and training sizes, Habash and Rambow (2005) report 99.3% word accuracy on tokenization, and Diab et al. (2007) reports a score of 99.1%. Habash and Rambow (2005) reported 97.6% on the LDC-supplied reduced tag set, and Diab et al. (2007) reported 96.6%. The LDC-supplied tag set used is smaller than the one in this paper (24 tags), but does distinguish between NOUN and ADJ. However, both (Habash and Rambow, 2005; Diab et al., 2007) assume gold

tokenization for evaluation of POS results, which we do not. The “MorphPOS” task in (Roth et al., 2008), 96.4%, is somewhat similar to ours in that it scores on a “core tag”, but unlike for us there is only one such tag for a source token (easier) but it distinguishes between NOUN and ADJ (harder).

We would like to do a direct comparison by simply running the above systems on the exact same data and evaluating them the same way. However, this unfortunately has to wait until new versions are released that work with the current version of the SAMA morphological analyzer and ATB.

5 Future Work

Obvious future work starts with the need to include determiner information in the POS tags and the important NOUN/ADJ distinction. There are various possibilities for recovering this information, such as (1) using a different module combining NOUN/ADJ disambiguation together with NP chunking, or (2) simply including NOUN/ADJ in the current classifier instead of NOA. We will be implementing and comparing these alternatives. We also will be using this system as a preprocessing step for a parser, as part of a complete Arabic NLP pipeline.

Acknowledgements

We thank Ann Bies, David Graff, Nizar Habash, Mohamed Maamouri, and Mitch Marcus for helpful discussions and comments. This work was supported by the Defense Advanced Research Projects Agency, GALE Program Grant No. HR0011-06-1-0003 and by the GALE program, DARPA/CMO Contract No. HR0011-06-C-0022. The content of this paper does not necessarily reflect the position or the policy of the Government, and no official endorsement should be inferred.

References

- Tim Buckwalter. 2004. Buckwalter Arabic morphological analyzer version 2.0. Linguistic Data Consortium LDC2004L02.
- Mona Diab, Kadri Hacioglu, and Daniel Jurafsky. 2007. Automatic processing of Modern Standard Arabic text. In Abdelhadi Soudi, Antal van den Bosch, and Gunter Neumann, editors, *Arabic Computational Morphology*, pages 159–179. Springer.
- Mona Diab. 2009. Second generation tools (AMIRA 2.0): Fast and robust tokenization, pos tagging, and base phrase chunking. In *Proceedings of 2nd International Conference on Arabic Language Resources and Tools (MEDAR)*, Cairo, Egypt, April.
- Nizar Habash and Owen Rambow. 2005. Arabic tokenization, part-of-speech tagging and morphological disambiguation in one fell swoop. In *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL'05)*, pages 573–580, Ann Arbor, Michigan, June. Association for Computational Linguistics.
- Seth Kulick, Ann Bies, and Mohamed Maamouri. 2010. Consistent and flexible integration of morphological annotation in the Arabic Treebank. In *Language Resources and Evaluation (LREC)*.
- Mohamed Maamouri, Ann Bies, Sondos Krouna, Fatma Gaddeche, Basma Bouziri, Seth Kulick, Widad Mekki, and Tim Buckwalter. 2009a. Arabic Treebank Morphological and Syntactic guidelines, July. <http://projects.ldc.upenn.edu/ArabicTreebank>.
- Mohamed Maamouri, Ann Bies, Seth Kulick, Sondos Krouna, Fatma Gaddeche, and Wajdi Zaghouni. 2009b. Arabic treebank part 3 - v3.2. Linguistic Data Consortium LDC2010T08, April.
- Mohammed Maamouri, Basma Bouziri, Sondos Krouna, David Graff, Seth Kulick, and Tim Buckwalter. 2009c. Standard Arabic morphological analyzer (SAMA) version 3.1. Linguistic Data Consortium LDC2009E73.
- Andrew McCallum. 2002. Mallet: A machine learning for language toolkit. <http://mallet.cs.umass.edu>.
- Ryan Roth, Owen Rambow, Nizar Habash, Mona Diab, and Cynthia Rudin. 2008. Arabic morphological tagging, diacritization, and lemmatization using lexeme models and feature ranking. In *Proceedings of ACL-08: HLT, Short Papers*, pages 117–120, Columbus, Ohio, June. Association for Computational Linguistics.

Hierarchical A* Parsing with Bridge Outside Scores

Adam Pauls and Dan Klein
Computer Science Division
University of California at Berkeley
{adpauls, klein}@cs.berkeley.edu

Abstract

Hierarchical A* (HA*) uses a hierarchy of coarse grammars to speed up parsing without sacrificing optimality. HA* prioritizes search in refined grammars using Viterbi outside costs computed in coarser grammars. We present Bridge Hierarchical A* (BHA*), a modified Hierarchical A* algorithm which computes a novel outside cost called a *bridge* outside cost. These bridge costs mix finer outside scores with coarser inside scores, and thus constitute tighter heuristics than entirely coarse scores. We show that BHA* substantially outperforms HA* when the hierarchy contains only very coarse grammars, while achieving comparable performance on more refined hierarchies.

1 Introduction

The Hierarchical A* (HA*) algorithm of Felzenszwalb and McAllester (2007) allows the use of a hierarchy of coarse grammars to speed up parsing without sacrificing optimality. Pauls and Klein (2009) showed that a hierarchy of coarse grammars outperforms standard A* parsing for a range of grammars. HA* operates by computing Viterbi inside and outside scores in an agenda-based way, using outside scores computed under coarse grammars as heuristics which guide the search in finer grammars. The outside scores computed by HA* are auxiliary quantities, useful only because they form admissible heuristics for search in finer grammars.

We show that a modification of the HA* algorithm can compute modified *bridge* outside scores which are tighter bounds on the true outside costs in finer grammars. These bridge outside scores mix inside and outside costs from finer grammars with inside costs from coarser grammars. Because the bridge costs represent tighter estimates of the

true outside costs, we expect them to reduce the work of computing inside costs in finer grammars. At the same time, because bridge costs mix computation from coarser and finer levels of the hierarchy, they are more expensive to compute than purely coarse outside costs. Whether the work saved by using tighter estimates outweighs the extra computation needed to compute them is an empirical question.

In this paper, we show that the use of bridge outside costs substantially outperforms the HA* algorithm when the coarsest levels of the hierarchy are very loose approximations of the target grammar. For hierarchies with tighter estimates, we show that BHA* obtains comparable performance to HA*. In other words, BHA* is more robust to poorly constructed hierarchies.

2 Previous Work

In this section, we introduce notation and review HA*. Our presentation closely follows Pauls and Klein (2009), and we refer the reader to that work for a more detailed presentation.

2.1 Notation

Assume we have input sentence $s_0 \dots s_{n-1}$ of length n , and a hierarchy of m weighted context-free grammars $\mathcal{G}_1 \dots \mathcal{G}_m$. We call the most refined grammar \mathcal{G}_m the *target* grammar, and all other (coarser) grammars *auxiliary* grammars. Each grammar \mathcal{G}_t has a set of symbols denoted with capital letters and a subscript indicating the level in the hierarchy, including a distinguished goal (root) symbol G_t . Without loss of generality, we assume Chomsky normal form, so each non-terminal rule r in \mathcal{G}_t has the form $r = A_t \rightarrow B_t C_t$ with weight w_r .

Edges are labeled spans $e = (A_t, i, j)$. The weight of a derivation is the sum of rule weights in the derivation. The weight of the best (minimum) inside derivation for an edge e is called the Viterbi inside score $\beta(e)$, and the weight of the

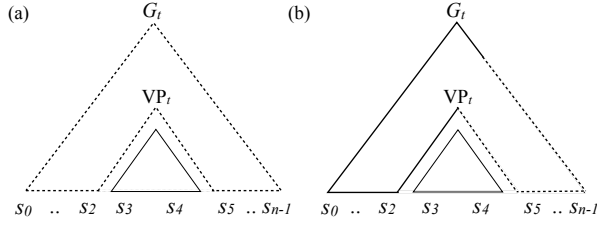


Figure 1: Representations of the different types of items used in parsing and how they depend on each other. (a) In HA*, the inside item $I(VP_t, 3, 5)$ relies on the coarse outside item $O(\pi_t(VP_t), 3, 5)$ for outside estimates. (b) In BHA*, the same inside item relies on the bridge outside item $\tilde{O}(VP_t, 3, 5)$, which mixes coarse and refined outside costs. The coarseness of an item is indicated with dotted lines.

best derivation of $G \rightarrow s_0 \dots s_{i-1} A_t s_j \dots s_{n-1}$ is called the Viterbi outside score $\alpha(e)$. The goal of a 1-best parsing algorithm is to compute the Viterbi inside score of the edge $(G_m, 0, n)$; the actual best parse can be reconstructed from backpointers in the standard way.

We assume that each auxiliary grammar \mathcal{G}_{t-1} forms a *relaxed projection* of \mathcal{G}_t . A grammar \mathcal{G}_{t-1} is a *projection* of \mathcal{G}_t if there exists some many-to-one onto function π_t which maps each symbol in \mathcal{G}_t to a symbol in \mathcal{G}_{t-1} ; hereafter, we will use A'_t to represent $\pi_t(A_t)$. A projection is *relaxed* if, for every rule $r = A_t \rightarrow B_t C_t$ with weight w_r the projection $r' = A'_t \rightarrow B'_t C'_t$ has weight $w_{r'} \leq w_r$ in \mathcal{G}_{t-1} . In other words, the weight of r' is a lower bound on the weight of all rules r in \mathcal{G}_t which project to r' .

2.2 Deduction Rules

HA* and our modification BHA* can be formulated in terms of prioritized weighted deduction rules (Shieber et al., 1995; Felzenszwalb and McAllester, 2007). A *prioritized weighted deduction rule* has the form

$$\phi_1 : w_1, \dots, \phi_n : w_n \xrightarrow{p(w_1, \dots, w_n)} \phi_0 : g(w_1, \dots, w_n)$$

where ϕ_1, \dots, ϕ_n are the *antecedent items* of the deduction rule and ϕ_0 is the *conclusion item*. A deduction rule states that, given the antecedents ϕ_1, \dots, ϕ_n with weights w_1, \dots, w_n , the conclusion ϕ_0 can be formed with weight $g(w_1, \dots, w_n)$ and priority $p(w_1, \dots, w_n)$.

These deduction rules are “executed” within a generic agenda-driven algorithm, which constructs items in a prioritized fashion. The algorithm maintains an *agenda* (a priority queue of

items), as well as a *chart* of items already processed. The fundamental operation of the algorithm is to pop the highest priority item ϕ from the agenda, put it into the chart with its current weight, and form using deduction rules any items which can be built by combining ϕ with items already in the chart. If new or improved, resulting items are put on the agenda with priority given by $p(\cdot)$. Because all antecedents must be constructed before a deduction rule is executed, we sometimes refer to particular conclusion item as “waiting” on an other item(s) before it can be built.

2.3 HA*

HA* can be formulated in terms of two types of items. *Inside* items $I(A_t, i, j)$ represent possible derivations of the edge (A_t, i, j) , while *outside* items $O(A_t, i, j)$ represent derivations of $G \rightarrow s_1 \dots s_{i-1} A_t s_j \dots s_n$ rooted at $(G_t, 0, n)$. See Figure 1(a) for a graphical depiction of these edges. Inside items are used to compute Viterbi inside scores under grammar \mathcal{G}_t , while outside items are used to compute Viterbi outside scores.

The deduction rules which construct inside and outside items are given in Table 1. The IN deduction rule combines two inside items over smaller spans with a grammar rule to form an inside item over larger spans. The weight of the resulting item is the sum of the weights of the smaller inside items and the grammar rule. However, the IN rule also requires that an outside score in the coarse grammar¹ be computed before an inside item is built. Once constructed, this coarse outside score is added to the weight of the conclusion item to form the priority of the resulting item. In other words, the coarse outside score computed by the algorithm plays the same role as a heuristic in standard A* parsing (Klein and Manning, 2003).

Outside scores are computed by the OUT-L and OUT-R deduction rules. These rules combine an outside item over a large span and inside items over smaller spans to form outside items over smaller spans. Unlike the IN deduction, the OUT deductions only involve items from the same level of the hierarchy. That is, whereas inside scores wait on coarse outside scores to be constructed, outside scores wait on inside scores at the same level in the hierarchy.

Conceptually, these deduction rules operate by

¹For the coarsest grammar \mathcal{G}_1 , the IN rule builds rules using 0 as an outside score.

HA*

IN:	$I(B_t, i, l) : w_1$	$I(C_t, l, j) : w_2$	<u>$O(A'_t, i, j) : w_3$</u>	$\xrightarrow{w_1+w_2+w_r+w_3}$	$I(A_t, i, j) : w_1 + w_2 + w_r$
OUT-L:	$O(A_t, i, j) : w_1$	$I(B_t, i, l) : w_2$	$I(C_t, l, j) : w_3$	$\xrightarrow{w_1+w_3+w_r+w_2}$	$O(B_t, i, l) : w_1 + w_3 + w_r$
OUT-R:	$O(A_t, i, j) : w_1$	$I(B_t, i, l) : w_2$	$I(C_t, l, j) : w_3$	$\xrightarrow{w_1+w_2+w_r+w_3}$	$O(C_t, l, j) : w_1 + w_2 + w_r$

Table 1: HA* deduction rules. Red underline indicates items constructed under the previous grammar in the hierarchy.

BHA*

B-IN:	$I(B_t, i, l) : w_1$	$I(C_t, l, j) : w_2$	$\tilde{O}(A_t, i, j) : w_3$	$\xrightarrow{w_1+w_2+w_r+w_3}$	$I(A_t, i, j) : w_1 + w_2 + w_r$
B-OUT-L:	$\tilde{O}(A_t, i, j) : w_1$	<u>$I(B'_t, i, l) : w_2$</u>	<u>$I(C'_t, l, j) : w_3$</u>	$\xrightarrow{w_1+w_r+\underline{w_2+w_3}}$	$\tilde{O}(B_t, i, l) : w_1 + w_r + \underline{w_3}$
B-OUT-R:	$\tilde{O}(A_t, i, j) : w_1$	$I(B_t, i, l) : w_2$	<u>$I(C'_t, l, j) : w_3$</u>	$\xrightarrow{w_1+w_2+w_r+w_3}$	$\tilde{O}(C_t, l, j) : w_1 + w_2 + w_r$

Table 2: BHA* deduction rules. Red underline indicates items constructed under the previous grammar in the hierarchy.

first computing inside scores bottom-up in the coarsest grammar, then outside scores top-down in the same grammar, then inside scores in the next finest grammar, and so on. However, the crucial aspect of HA* is that items from all levels of the hierarchy compete on the same queue, interleaving the computation of inside and outside scores at all levels. The HA* deduction rules come with three important guarantees. The first is a *monotonicity* guarantee: each item is popped off the agenda in order of its *intrinsic priority* $\hat{p}(\cdot)$. For inside items $I(e)$ over edge e , this priority $\hat{p}(I(e)) = \beta(e) + \alpha(e')$ where e' is the projection of e . For outside items $O(\cdot)$ over edge e , this priority is $\hat{p}(O(e)) = \beta(e) + \alpha(e)$.

The second is a *correctness* guarantee: when an inside/outside item is popped of the agenda, its weight is its true Viterbi inside/outside cost. Taken together, these two imply an *efficiency* guarantee, which states that only items x whose intrinsic priority $\hat{p}(x)$ is less than or equal to the Viterbi inside score of the goal are removed from the agenda.

2.4 HA* with Bridge Costs

The outside scores computed by HA* are useful for prioritizing computation in more refined grammars. The key property of these scores is that they form consistent and admissible heuristic costs for more refined grammars, but coarse outside costs are not the only quantity which satisfy this requirement. As an alternative, we propose a novel “bridge” outside cost $\tilde{\alpha}(e)$. Intuitively, this cost represents the cost of the best derivation where rules “above” and “left” of an edge e come from \mathcal{G}_t , and rules “below” and “right” of the e come from \mathcal{G}_{t-1} ; see Figure 2 for a graphical depiction. More formally, let the *spine* of an edge $e = (A_t, i, j)$ for some derivation d be

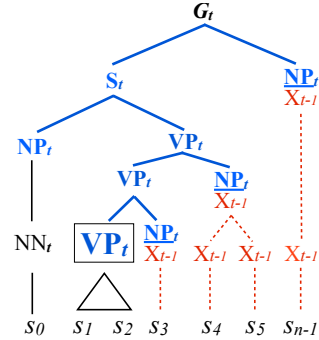


Figure 2: A concrete example of a possible bridge outside derivation for the bridge item $\tilde{O}(VP_t, 1, 4)$. This edge is boxed for emphasis. The spine of the derivation is shown in bold and colored in blue. Rules from a coarser grammar are shown with dotted lines, and colored in red. Here we have the simple projection $\pi_t(A) = X, \forall A$.

the sequence of rules between e and the root edge $(G_t, 0, n)$. A *bridge outside derivation* of e is a derivation d of $G \rightarrow s_1 \dots s_i A_t s_{j+1} \dots s_n$ such that every rule on or left of the spine comes from \mathcal{G}_t , and all other rules come from \mathcal{G}_{t-1} . The score of the best such derivation for e is the bridge outside cost $\tilde{\alpha}(e)$.

Like ordinary outside costs, bridge outside costs form consistent and admissible estimates of the true Viterbi outside score $\alpha(e)$ of an edge e . Because bridge costs mix rules from the finer and coarser grammar, bridge costs are at least as good an estimate of the true outside score as entirely coarse outside costs, and will in general be much tighter. That is, we have

$$\alpha(e') \leq \tilde{\alpha}(e) \leq \alpha(e)$$

In particular, note that the bridge costs become better approximations farther right in the sentence, and the bridge cost of the last word in the sentence is equal to the Viterbi outside cost of that word.

To compute bridge outside costs, we introduce

bridge outside items $\tilde{O}(A_t, i, j)$, shown graphically in Figure 1(b). The deduction rules which build both inside items and bridge outside items are shown in Table 2. The rules are very similar to those which define HA^* , but there are two important differences. First, inside items wait for bridge outside items *at the same level*, while outside items wait for inside items from the previous level. Second, the left and right outside deductions are no longer symmetric – bridge outside items can be extended to the left given two coarse inside items, but can only be extended to the right given an exact inside item on the left and coarse inside item on the right.

2.5 Guarantees

These deduction rules come with guarantees analogous to those of HA^* . The monotonicity guarantee ensures that inside and (bridge) outside items are processed in order of:

$$\begin{aligned}\hat{p}(I(e)) &= \beta(e) + \tilde{\alpha}(e) \\ \hat{p}(\tilde{O}(e)) &= \tilde{\alpha}(e) + \beta(e')\end{aligned}$$

The correctness guarantee ensures that when an item is removed from the agenda, its weight will be equal to $\beta(e)$ for inside items and $\tilde{\alpha}(e)$ for bridge items. The efficiency guarantee remains the same, though because the intrinsic priorities are different, the set of items processed will be different from those processed by HA^* .

A proof of these guarantees is not possible due to space restrictions. The proof for BHA^* follows the proof for HA^* in Felzenszwalb and McAllester (2007) with minor modifications. The key property of HA^* needed for these proofs is that coarse outside costs form consistent and admissible heuristics for inside items, and exact inside costs form consistent and admissible heuristics for outside items. BHA^* also has this property, with bridge outside costs forming admissible and consistent heuristics for inside items, and coarse inside costs forming admissible and consistent heuristics for outside items.

3 Experiments

The performance of BHA^* is determined by the efficiency guarantee given in the previous section. However, we cannot determine in advance whether BHA^* will be faster than HA^* . In fact, BHA^* has the potential to be slower – BHA^*

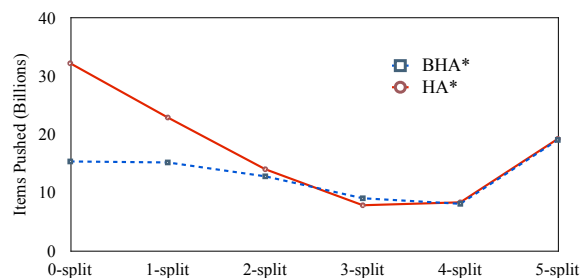


Figure 3: Performance of HA^* and BHA^* as a function of increasing refinement of the coarse grammar. Lower is faster.

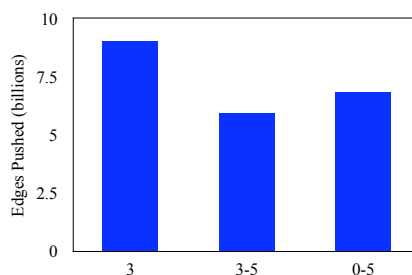


Figure 4: Performance of BHA^* on hierarchies of varying size. Lower is faster. Along the x-axis, we show which coarse grammars were used in the hierarchy. For example, 3-5 indicates the 3-,4-, and 5-split grammars were used as coarse grammars.

builds both inside and bridge outside items under the target grammar, where HA^* only builds inside items. It is an empirical, grammar- and hierarchy-dependent question whether the increased tightness of the outside estimates outweighs the additional cost needed to compute them. We demonstrate empirically in this section that for hierarchies with very loosely approximating coarse grammars, BHA^* can outperform HA^* , while for hierarchies with good approximations, performance of the two algorithms is comparable.

We performed experiments with the grammars of Petrov et al. (2006). The training procedure for these grammars produces a hierarchy of increasingly refined grammars through state-splitting, so a natural projection function π_t is given. We used the Berkeley Parser² to learn such grammars from Sections 2-21 of the Penn Treebank (Marcus et al., 1993). We trained with 6 split-merge cycles, producing 7 grammars. We tested these grammars on 300 sentences of length ≤ 25 of Section 23 of the Treebank. Our “target grammar” was in all cases the most split grammar.

²<http://berkeleyparser.googlecode.com>

In our first experiment, we construct 2-level hierarchies consisting of one coarse grammar and the target grammar. By varying the coarse grammar from the 0-split (X-bar) through 5-split grammars, we can investigate the performance of each algorithm as a function of the coarseness of the coarse grammar. We follow Pauls and Klein (2009) in using the number of items pushed as a machine- and implementation-independent measure of speed. In Figure 3, we show the performance of HA* and BHA* as a function of the total number of items pushed onto the agenda. We see that for very coarse approximating grammars, BHA* substantially outperforms HA*, but for more refined approximating grammars the performance is comparable, with HA* slightly outperforming BHA* on the 3-split grammar.

Finally, we verify that BHA* can benefit from multi-level hierarchies as HA* can. We constructed two multi-level hierarchies: a 4-level hierarchy consisting of the 3-,4-,5-, and 6- split grammars, and 7-level hierarchy consisting of all grammars. In Figure 4, we show the performance of BHA* on these multi-level hierarchies, as well as the best 2-level hierarchy from the previous experiment. Our results echo the results of Pauls and Klein (2009): although the addition of the reasonably refined 4- and 5-split grammars produces modest performance gains, the addition of coarser grammars can actually hurt overall performance.

Acknowledgements

This project is funded in part by the NSF under grant 0643742 and an NSERC Postgraduate Fellowship.

References

- P. Felzenszwalb and D. McAllester. 2007. The generalized A* architecture. *Journal of Artificial Intelligence Research*.
- Dan Klein and Christopher D. Manning. 2003. A* parsing: Fast exact Viterbi parse selection. In *Proceedings of the Human Language Technology Conference and the North American Association for Computational Linguistics (HLT-NAACL)*, pages 119–126.
- M. Marcus, B. Santorini, and M. Marcinkiewicz. 1993. Building a large annotated corpus of English: The Penn Treebank. In *Computational Linguistics*.
- Adam Pauls and Dan Klein. 2009. Hierarchical search for parsing. In *Proceedings of The Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*.
- Slav Petrov, Leon Barrett, Romain Thibaux, and Dan Klein. 2006. Learning accurate, compact, and interpretable tree annotation. In *Proceedings of the Association for Computational Linguistics (ACL)*.
- Stuart M. Shieber, Yves Schabes, and Fernando C. N. Pereira. 1995. Principles and implementation of deductive parsing. *Journal of Logic Programming*, 24:3–36.

Using Parse Features for Preposition Selection and Error Detection

Joel Tetreault
Educational Testing Service
Princeton
NJ, USA
JTetreault@ets.org

Jennifer Foster
NCLT
Dublin City University
Ireland
jfoster@computing.dcu.ie

Martin Chodorow
Hunter College of CUNY
New York, NY, USA
martin.chodorow
@hunter.cuny.edu

Abstract

We evaluate the effect of adding parse features to a leading model of preposition usage. Results show a significant improvement in the preposition selection task on native speaker text and a modest increment in precision and recall in an ESL error detection task. Analysis of the parser output indicates that it is robust enough in the face of noisy non-native writing to extract useful information.

1 Introduction

The task of preposition error detection has received a considerable amount of attention in recent years because selecting an appropriate preposition poses a particularly difficult challenge to learners of English as a second language (ESL). It is not only ESL learners that struggle with English preposition usage — automatically detecting preposition errors made by ESL speakers is a challenging task for NLP systems. Recent state-of-the-art systems have precision ranging from 50% to 80% and recall as low as 10% to 20%.

To date, the conventional wisdom in the error detection community has been to avoid the use of statistical parsers under the belief that a WSJ-trained parser's performance would degrade too much on noisy learner texts and that the traditionally hard problem of prepositional phrase attachment would be even harder when parsing ESL writing. However, there has been little substantial research to support or challenge this view. In this paper, we investigate the following research question: *Are parser output features helpful in modeling preposition usage in well-formed text and learner text?*

We recreate a state-of-the-art preposition usage system (Tetreault and Chodorow (2008), henceforth T&C08) originally trained with lexical features and augment it with parser output features. We employ the Stanford parser in our experiments because it consists of a competitive phrase structure parser *and* a constituent-to-dependency conversion tool (Klein and Manning, 2003a; Klein and Manning, 2003b; de Marneffe et al., 2006; de Marneffe and Manning, 2008). We compare the original model with the parser-augmented model on the tasks of preposition selection in well-formed text (fluent writers) and preposition error detection in learner texts (ESL writers).

This paper makes the following contributions:

- We demonstrate that parse features have a significant impact on preposition selection in well-formed text. We also show which features have the greatest effect on performance.
- We show that, despite the noisiness of learner text, parse features can actually make small, albeit non-significant, improvements to the performance of a state-of-the-art preposition error detection system.
- We evaluate the accuracy of parsing and especially preposition attachment in learner texts.

2 Related Work

T&C08, De Felice and Pulman (2008) and Gamon *et al.* (2008) describe very similar preposition error detection systems in which a model of correct prepositional usage is trained from well-formed text and a writer's preposition is compared with the predictions of this model. It is difficult to directly compare these systems since they are trained and tested on different data sets

but they achieve accuracy in a similar range. Of these systems, only the DAPPER system (De Felice and Pulman, 2008; De Felice and Pulman, 2009; De Felice, 2009) uses a parser, the C&C parser (Clark and Curran, 2007)), to determine the head and complement of the preposition. De Felice and Pulman (2009) remark that the parser tends to be misled more by spelling errors than by grammatical errors. The parser is fundamental to their system and they do not carry out a comparison of the use of a parser to determine the preposition’s attachments versus the use of shallower techniques. T&C08, on the other hand, reject the use of a parser because of the difficulties they foresee in applying one to learner data. Hermet *et al.* (2008) make only limited use of the Xerox Incremental Parser in their preposition error detection system. They split the input sentence into the chunks before and after the preposition, and parse both chunks separately. Only very shallow analyses are extracted from the parser output because they do not trust the full analyses.

Lee and Knutsson (2008) show that knowledge of the PP attachment site helps in the task of preposition selection by comparing a classifier trained on lexical features (the verb before the preposition, the noun between the verb and the preposition, if any, and the noun after the preposition) to a classifier trained on attachment features which explicitly state whether the preposition is attached to the preceding noun or verb. They also argue that a parser which is capable of distinguishing between arguments and adjuncts is useful for generating the correct preposition.

3 Augmenting a Preposition Model with Parse Features

To test the effects of adding parse features to a model of preposition usage, we replicated the lexical and combination feature model used in T&C08, training on 2M events extracted from a corpus of news and high school level reading materials. Next, we added the parse features to this model to create a new model “+Parse”. In 3.1 we describe the T&C08 system and features, and in 3.2 we describe the parser output features used to augment the model. We illustrate our features using the example phrase *many local groups around the country*. Fig. 1 shows the phrase structure tree and dependency triples returned by the Stanford parser for this phrase.

3.1 Baseline System

The work of Chodorow *et al.* (2007) and T&C08 treat the tasks of preposition selection and error detection as a classification problem. That is, given the context around a preposition and a model of correct usage, a classifier determines which of the 34 prepositions covered by the model is most appropriate for the context. A model of correct preposition usage is constructed by training a Maximum Entropy classifier (Ratnaparkhi, 1998) on millions of preposition contexts from well-formed text.

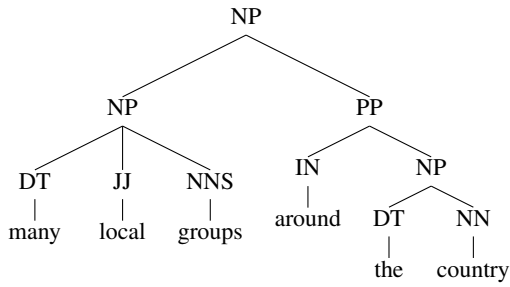
A context is represented by 25 lexical features and 4 combination features:

Lexical Token and POS n-grams in a 2 word window around the preposition, plus the head verb in the preceding verb phrase (PV), the head noun in the preceding noun phrase (PN) and the head noun in the following noun phrase (FN) when available (Chodorow *et al.*, 2007). Note that these are determined not through full syntactic parsing but rather through the use of a heuristic chunker. So, for the phrase *many local groups around the country*, examples of lexical features for the preposition *around* include: *FN = country*, *PN = groups*, *left-2-word-sequence = local-groups*, and *left-2-POS-sequence = JJ-NNS*.

Combination T&C08 expand on the lexical feature set by combining the PV, PN and FN features, resulting in features such as *PN-FN* and *PV-PN-FN*. POS and token versions of these features are employed. The intuition behind creating combination features is that the Maximum Entropy classifier does not automatically model the interactions between individual features. An example of the *PN-FN* feature is *groups-country*.

3.2 Parse Features

To augment the above model we experimented with 14 features divided among five main classes. Table 1 shows the features and their values for our *around* example. The **Preposition Head and Complement** feature represents the two basic attachment relations of the preposition, i.e. its head (what it is attached to) and its complement (what is attached to it). **Relation** specifies the relation between the head and complement. The **Preposition Head and Complement Combined** features are similar to the T&C08 Combination features except that they are extracted from parser output.



amod(groups-3, many-1)
 amod(groups-3, local-2)
 prep(groups-3, around-4)
 det(country-6, the-5)
 pobj(around-4, country-6)

Figure 1: Phrase structure tree and dependency triples produced by the Stanford parser for the phrase *many local groups around the country*

Prep. Head & Complement 1. head of the preposition: <i>groups</i> 2. POS of the head: NNS 3. complement of the preposition: <i>country</i> 4. POS of the complement: NN
Prep. Head & Complement Relation 5. Prep-Head relation name: <code>prep</code> 6. Prep-Comp relation name: <code>pobj</code>
Prep. Head & Complement Combined 7. Head-Complement tokens: <i>groups-country</i> 8. Head-Complement tags: NNS-NN
Prep. Head & Complement Mixed 9. Head Tag and Comp Token: NNS- <i>country</i> 10. Head Token and Comp Tag: <i>groups</i> -NN
Phrase Structure 11. Preposition Parent: PP 12. Preposition Grandparent: NP 13. Left context of preposition parent: NP 14. Right context of preposition parent: -

Table 1: Parse Features

Model	Accuracy
combination only	35.2
parse only	60.6
combination+parse	61.9
lexical only	64.4
combination+lexical (T&C08)	65.2
lexical+parse	68.1
all features (+Parse)	68.5

Table 2: Accuracy on preposition selection task for various feature combinations

The **Preposition Head and Complement Mixed** features are created by taking the first feature in the previous set and backing-off either the head or the complement to its POS tag. This mix of tags and tokens in a word-word dependency has proven to be an effective feature in sentiment analysis (Joshi and Penstein-Rosé, 2009). All the features described so far are extracted from the set of dependency triples output by the Stanford parser. The final set of features (**Phrase Structure**), however, is extracted directly from the phrase structure trees themselves.

4 Evaluation

In Section 4.1, we compare the T&C08 and +Parse models on the task of preposition selection on well-formed texts written by native speakers. For every preposition in the test set, we compare the system’s top preposition for that context to the writer’s preposition, and report accuracy rates. In Section 4.2, we evaluate the two models on ESL data. The task here is slightly different - if the most likely preposition according to the model differs from the likelihood of the writer’s preposition by a certain threshold amount, a preposition error is flagged.

4.1 Native Speaker Test Data

Our test set consists of 259K preposition events from the same source as the original training data. The T&C08 model performs at **65.2%** and when the parse features are added, the +Parse model improves performance by more than 3% to **68.5%**.¹ The improvement is statistically significant.

¹Prior research has shown preposition selection performance accuracy ranging from 65% to nearly 80%. The differences are largely due to different test sets and also training sizes. Given the time required to train large models, we report here experiments with a relatively small model.

Model	Accuracy
T&C08	65.2
+Phrase Structure Only	67.1
+Dependency Only	68.2
+Parse	68.5
+head-tag+comp-tag	66.9
+left	66.8
+grandparent	66.6
+head-token+comp-tag	66.6
+head-tag	66.5
+head-token	66.4
+head-tag+comp-token	66.1

Table 3: Which parse features are important? Feature Addition Experiment

Table 2 shows the effect of various feature class combinations on prediction accuracy. The results are clear: a significant performance improvement is obtained on the preposition selection task when features from parser output are added. The two best models in Table 2 contain parse features. The table also shows that the non-parser-based feature classes are not entirely subsumed by the parse features but rather provide, to varying degrees, complementary information.

Having established the effectiveness of parse features, we investigate which parse feature classes contribute the most. To test each contribution, we perform a feature addition experiment, separately adding features to the T&C08 model (see Table 3). We make three observations. First, while there is overlapping information between the dependency features and the phrase structure features, the phrase structure features *are* making a contribution. This is interesting because it suggests that a pure dependency parser might be less useful than a parser which explicitly produces both constituent and dependency information. Second, using a parser to identify the preposition head seems to be more useful than using it to identify the preposition complement.² Finally, as was the case for the T&C08 features, the combination parse features are also important (particularly the tag-tag or tag/token pairs).

4.2 ESL Test Data

Our test data consists of 5,183 preposition events extracted from a set of essays written by non-

²De Felice (2009) observes the same for the DAPPER system.

Method	Precision	Recall
T&C08	0.461	0.215
+Parse	0.486	0.225

Table 4: ESL Error Detection Results

native speakers for the Test of English as a Foreign Language (TOEFL®). The prepositions were judged by two trained annotators and checked by the authors using the preposition annotation scheme described in Tetreault and Chodorow (2008b). 4,881 of the prepositions were judged to be correct and the remaining 302 were judged to be incorrect.

The writer’s preposition is flagged as an error by the system if its likelihood according to the model satisfied a set of criteria (e.g., the difference between the probability of the system’s choice and the writer’s preposition is 0.8 or higher). Unlike the selection task where we use accuracy as the metric, we use precision and recall with respect to error detection. To date, performance figures that have been reported in the literature have been quite low, reflecting the difficulty of the task. Table 4 shows the performance figures for the T&C08 and +Parse models. Both precision and recall are higher for the +Parse model, however, given the low number of errors in our annotated test set, the difference is not statistically significant.

5 Parser Accuracy on ESL Data

To evaluate parser performance on ESL data, we manually inspected the phrase structure trees and dependency graphs produced by the Stanford parser for 210 ESL sentences, split into 3 groups: the sentences in the first group are fluent and contain no obvious grammatical errors, those in the second contain at least one preposition error and the sentences in the third are clearly ungrammatical with a variety of error types. For each preposition we note whether the parser was successful in determining its head and complement. The results for the three groups are shown in Table 5. The figures in the first row are for correct prepositions and those in the second are for incorrect ones.

The parser tends to do a better job of determining the preposition’s complement than its head which is not surprising given the well-known problem of PP attachment ambiguity. Given the preposition, the preceding noun, the preceding

	OK	
	Head	Comp
Prep Correct	86.7% (104/120)	95.0% (114/120)
Prep Incorrect	-	-
	Preposition Error	
	Head	Comp
Prep Correct	89.0% (65/73)	97.3% (71/73)
Prep Incorrect	87.1% (54/62)	96.8% (60/62)
	Ungrammatical	
	Head	Comp
Prep Correct	87.8% (115/131)	89.3% (117/131)
Prep Incorrect	70.8% (17/24)	87.5% (21/24)

Table 5: Parser Accuracy on Prepositions in a Sample of ESL Sentences

verb and the following noun, Collins (1999) reports an accuracy rate of 84.5% for a PP attachment classifier. When confronted with the same information, the accuracy of three trained annotators is 88.2%. Assuming 88.2% as an approximate PP-attachment upper bound, the Stanford parser appears to be doing a good job. Comparing the results over the three sentence groups, its ability to identify the preposition’s head is quite robust to grammatical noise.

Preposition errors in isolation do not tend to mislead the parser: in the second group which contains sentences which are largely fluent apart from preposition errors, there is little difference between the parser’s accuracy on the correctly used prepositions and the incorrectly used ones. Examples are

```
(S (NP I)
  (VP had
    (NP (NP a trip)
      (PP for (NP Italy)))
    )
  )
)
```

in which the erroneous preposition *for* is correctly attached to the noun *trip*, and

```
(S (NP A scientist)
  (VP devotes
    (NP (NP his prime part)
      (PP of (NP his life)))
    )
  (PP in (NP research))
  )
)
```

in which the erroneous preposition *in* is correctly attached to the verb *devotes*.

6 Conclusion

We have shown that the use of a parser can boost the accuracy of a preposition selection model tested on well-formed text. In the error detection task, the improvement is less marked. Nevertheless, examination of parser output shows the parse features can be extracted reliably from ESL data.

For our immediate future work, we plan to carry out the ESL evaluation on a larger test set to better gauge the usefulness of a parser in this context, to carry out a detailed error analysis to understand why certain parse features are effective and to explore a larger set of features.

In the longer term, we hope to compare different types of parsers in both the preposition selection and error detection tasks, i.e. a task-based parser evaluation in the spirit of that carried out by Miyao *et al.* (2008) on the task of protein pair interaction extraction. We would like to further investigate the role of parsing in error detection by looking at other error types and other text types, e.g. machine translation output.

Acknowledgments

We would like to thank Rachele De Felice and the reviewers for their very helpful comments.

References

- Martin Chodorow, Joel Tetreault, and Na-Rae Han. 2007. Detection of grammatical errors involving prepositions. In *Proceedings of the 4th ACL-SIGSEM Workshop on Prepositions*, Prague, Czech Republic, June.
- Stephen Clark and James R. Curran. 2007. Wide-coverage efficient statistical parsing with CCG and log-linear models. *Computational Linguistics*, 33(4):493–552.
- Michael Collins. 1999. *Head-driven Statistical Models for Natural Language Parsing*. Ph.D. thesis, University of Pennsylvania.
- Rachele De Felice and Stephen G. Pulman. 2008. A classifier-based approach to preposition and determiner error correction in L2 english. In *Proceedings of the 22nd COLING*, Manchester, United Kingdom.
- Rachele De Felice and Stephen Pulman. 2009. Automatic detection of preposition errors in learning writing. *CALICO Journal*, 26(3):512–528.
- Rachele De Felice. 2009. *Automatic Error Detection in Non-native English*. Ph.D. thesis, Oxford University.

- Marie-Catherine de Marneffe and Christopher D. Manning. 2008. The stanford typed dependencies representation. In *Proceedings of the COLING08 Workshop on Cross-framework and Cross-domain Parser Evaluation*, Manchester, United Kingdom.
- Marie-Catherine de Marneffe, Bill MacCartney, and Christopher D. Manning. 2006. Generating typed dependency parses from phrase structure parses. In *Proceedings of LREC*, Genoa, Italy.
- Michael Gamon, Jianfeng Gao, Chris Brockett, Alexandre Klementiev, William B. Dolan, Dmitriy Belenko, and Lucy Vanderwende. 2008. Using contextual speller techniques and language modelling for ESL error correction. In *Proceedings of the International Joint Conference on Natural Language Processing*, Hyderabad, India.
- Matthieu Hermet, Alain Désilets, and Stan Szpakowicz. 2008. Using the web as a linguistic resource to automatically correct lexico-syntactic errors. In *Proceedings of LREC*, Marrekech, Morocco.
- Mahesh Joshi and Carolyn Penstein-Rosé. 2009. Generalizing dependency features for opinion mining. In *Proceedings of the ACL-IJCNLP 2009 Conference Short Papers*, pages 313–316, Singapore.
- Dan Klein and Christopher D. Manning. 2003a. Accurate unlexicalized parsing. In *Proceedings of the 41st Annual Meeting of the ACL*, pages 423–430, Sapporo, Japan.
- Dan Klein and Christopher D. Manning. 2003b. Fast exact inference with a factored model for exact parsing. In *Advances in Neural Information Processing Systems*, pages 3–10. MIT Press, Cambridge, MA.
- John Lee and Ola Knutsson. 2008. The role of PP attachment in preposition generation. In *Proceedings of CICLing*. Springer-Verlag Berlin Heidelberg.
- Yusuke Miyao, Rune Saetre, Kenji Sagae, Takuya Matsuzaki, and Jun'ichi Tsujii. 2008. Task-oriented evaluation of syntactic parsers and their representations. In *Proceedings of the 46th Annual Meeting of the ACL*, pages 46–54, Columbus, Ohio.
- Adwait Ratnaparkhi. 1998. *Maximum Entropy Models for natural language ambiguity resolution*. Ph.D. thesis, University of Pennsylvania.
- Joel Tetreault and Martin Chodorow. 2008. The ups and downs of preposition error detection in ESL writing. In *Proceedings of the 22nd COLING*, Manchester, United Kingdom.
- Joel Tetreault and Martin Chodorow. 2008b. Native Judgments of non-native usage: Experiments in preposition error detection. In *COLING Workshop on Human Judgments in Computational Linguistics*, Manchester, United Kingdom.

Distributional Similarity vs. PU Learning for Entity Set Expansion

Xiao-Li Li

Institute for Infocomm Research,
1 Fusionopolis Way #21-01 Connexis
Singapore 138632
xlli@i2r.a-star.edu.sg

Lei Zhang

University of Illinois at Chicago,
851 South Morgan Street, Chicago,
Chicago, IL 60607-7053, USA
zhang3@cs.uic.edu

Bing Liu

University of Illinois at Chicago,
851 South Morgan Street, Chicago,
Chicago, IL 60607-7053, USA
liub@cs.uic.edu

See-Kiong Ng

Institute for Infocomm Research,
1 Fusionopolis Way #21-01 Connexis
Singapore 138632
skng@i2r.a-star.edu.sg

Abstract

Distributional similarity is a classic technique for entity set expansion, where the system is given a set of seed entities of a particular class, and is asked to expand the set using a corpus to obtain more entities of the same class as represented by the seeds. This paper shows that a machine learning model called *positive and unlabeled learning* (PU learning) can model the set expansion problem better. Based on the test results of 10 corpora, we show that a PU learning technique outperformed distributional similarity significantly.

In machine learning, there is a class of semi-supervised learning algorithms that learns from *positive* and *unlabeled* examples (PU learning for short). The key characteristic of PU learning is that there is no negative training example available for learning. This class of algorithms is less known to the natural language processing (NLP) community compared to some other semi-supervised learning models and algorithms.

PU learning is a two-class classification model. It is stated as follows (Liu *et al.* 2002): Given a set P of positive examples of a particular class and a set U of unlabeled examples (containing hidden positive and negative cases), a classifier is built using P and U for classifying the data in U or future test cases. The results can be either binary decisions (whether each test case belongs to the positive class or not), or a ranking based on how likely each test case belongs to the positive class represented by P . Clearly, the set expansion problem can be mapped into PU learning exactly, with S and D as P and U respectively.

This paper shows that a PU learning method called S-EM (Liu *et al.* 2002) outperforms distributional similarity considerably based on the results from 10 corpora. The experiments involved extracting named entities (e.g., product and organization names) of the same type or class as the given seeds. Additionally, we also compared S-EM with a recent method, called *Bayesian Sets* (Ghahramani and Heller, 2005), which was designed specifically for set expansion. It also does not perform as well as PU learning. We will explain why PU learning performs better than both methods in Section 5. We believe that this finding is of interest to the NLP community.

1 Introduction

The entity set expansion problem is defined as follows: Given a set S of seed entities of a particular class, and a set D of candidate entities (e.g., extracted from a text corpus), we wish to determine which of the entities in D belong to S . In other words, we “expand” the set S based on the given seeds. This is clearly a classification problem which requires arriving at a binary decision for each entity in D (belonging to S or not). However, in practice, the problem is often solved as a ranking problem, i.e., ranking the entities in D based on their likelihoods of belonging to S .

The classic method for solving this problem is based on *distributional similarity* (Pantel *et al.* 2009; Lee, 1998). The approach works by comparing the similarity of the surrounding word distributions of each candidate entity with the seed entities, and then ranking the candidate entities using their similarity scores.

There is another approach used in the Web environment for entity set expansion. It exploits Web page structures to identify lists of items using wrapper induction or other techniques. The idea is that items in the same list are often of the same type. This approach is used by Google Sets (Google, 2008) and Boo!Wa! (Wang and Cohen, 2008). However, as it relies on Web page structures, it is not applicable to general free texts.

2 Three Different Techniques

2.1 Distributional Similarity

Distributional similarity is a classic technique for the entity set expansion problem. It is based on the hypothesis that words with similar meanings tend to appear in similar contexts (Harris, 1985). As such, a method based on distributional similarity typically fetches the surrounding contexts for each term (i.e. both seeds and candidates) and represents them as vectors by using TF-IDF or PMI (*Pointwise Mutual Information*) values (Lin, 1998; Gorman and Curran, 2006; Paşca *et al.* 2006; Agirre *et al.* 2009; Pantel *et al.* 2009). Similarity measures such as *Cosine*, *Jaccard*, *Dice*, etc, can then be employed to compute the similarities between each candidate vector and the seeds centroid vector (one centroid vector for all seeds). Lee (1998) surveyed and discussed various distribution similarity measures.

2.2 PU Learning and S-EM

PU learning is a semi-supervised or partially supervised learning model. It learns from positive and unlabeled examples as opposed to the model of learning from a small set of labeled examples of every class and a large set of unlabeled examples, which we call LU learning (L and U stand for *labeled* and *unlabeled* respectively) (Blum and Mitchell, 1998; Nigam *et al.* 2000)

There are several PU learning algorithms (Liu *et al.* 2002; Yu *et al.* 2002; Lee and Liu, 2003; Li *et al.* 2003; Elkan and Noto, 2008). In this work, we used the S-EM algorithm given in (Liu *et al.* 2002). S-EM is efficient as it is based on naïve Bayesian (NB) classification and also performs well. The main idea of S-EM is to use a *spy* technique to identify some *reliable negatives* (RN) from the unlabeled set U , and then use an EM algorithm to learn from P , RN and $U-RN$.

The *spy* technique in S-EM works as follows (Figure 1): First, a small set of positive examples (denoted by SP) from P is randomly sampled (line 2). The default sampling ratio in S-EM is $s = 15\%$, which we also used in our experiments.

The positive examples in SP are called “spies”. Then, a NB classifier is built using the set $P-SP$ as positive and the set $U \cup SP$ as negative (line 3, 4, and 5). The NB classifier is applied to classify each $u \in U \cup SP$, i.e., to assign a probabilistic class label $p(+|u)$ (+ means positive). The probabilistic labels of the spies are then used to decide reliable negatives (RN). In particular, a probability threshold t is determined using the probabilistic labels of spies in SP and the input parameter l (noise level). Due to space constraints, we are unable to explain l . Details can be found in (Liu *et al.* 2002). t is then used to find RN from U (lines 8-10). The idea of the spy technique is clear. Since spy examples are from P and are put into U in building the NB classifier, they should behave similarly to the hidden positive cases in U . Thus, they can help us find the set RN .

Algorithm Spy(P, U, s, l)

1. $RN \leftarrow \emptyset$; // Reliable negative set
2. $SP \leftarrow \text{Sample}(P, s\%)$;
3. Assign each example in $P - SP$ the class label +1;
4. Assign each example in $U \cup SP$ the class label -1;
5. $C \leftarrow \text{NB}(P - S, U \cup SP)$; // Produce a NB classifier
6. Classify each $u \in U \cup SP$ using C ;
7. Decide a probability threshold t using SP and l ;
8. **for** each $u \in U$ **do**
9. **if** its probability $p(+|u) < t$ **then**
10. $RN \leftarrow RN \cup \{u\}$;

Figure 1. Spy technique for extracting reliable negatives (RN) from U .

Given the positive set P , the reliable negative set RN and the remaining unlabeled set $U-RN$, an Expectation-Maximization (EM) algorithm is run. In S-EM, EM uses the naïve Bayesian classification as its base method. The detailed algorithm is given in (Liu *et al.* 2002).

2.3 Bayesian Sets

Bayesian Sets, as its name suggests, is based on Bayesian inference, and was designed specifically for the set expansion problem (Ghahramani and Heller, 2005). The algorithm learns from a seeds set (i.e., a positive set P) and an unlabeled candidate set U . Although it was not designed as a PU learning method, it has similar characteristics and produces similar results as PU learning. However, there is a major difference. PU learning is a classification model, while Bayesian Sets is a ranking method. This difference has a major implication on the results that they produce as we will discuss in Section 5.3.

In essence, Bayesian Sets learns a score func-

tion using P and U to generate a score for each unlabeled case $u \in U$. The function is as follows:

$$score(u) = \frac{p(u|P)}{p(u)} \quad (1)$$

where $p(u|P)$ represents how probable u belongs to the positive class represented by P . $p(u)$ is the prior probability of u . Using the Bayes' rule, equation (1) can be re-written as:

$$score(u) = \frac{p(u, P)}{p(u)p(P)} \quad (2)$$

Following the idea, Ghahramani and Heller (2005) proposed a computable score function. The scores can be used to rank the unlabeled candidates in U to reflect how likely each $u \in U$ belongs to P . The mathematics for computing the score is involved. Due to the limited space, we cannot discuss it here. See (Ghahramani and Heller, 2005) for details. In (Heller and Ghahramani, 2006), Bayesian Sets was also applied to an image retrieval application.

3 Data Generation for Distributional Similarity, Bayesian Sets and S-EM

Preparing the data for distributional similarity is fairly straightforward. Given the seeds set S , a seeds centroid vector is produced using the surrounding word contexts (see below) of all occurrences of all the seeds in the corpus (Pantel et al, 2009). In a similar way, a centroid is also produced for each candidate (or unlabeled) entity.

Candidate entities: Since we are interested in named entities, we select single words or phrases as candidate entities based on their corresponding part-of-speech (POS) tags. In particular, we choose the following POS tags as entity indicators — NNP (proper noun), NNPS (plural proper noun), and CD (cardinal number). We regard a phrase (could be one word) with a sequence of NNP, NNPS and CD POS tags as one candidate entity (CD cannot be the first word unless it starts with a letter), e.g., “Windows/NNP 7/CD” and “Nokia/NNP N97/CD” are regarded as two candidates “Windows 7” and “Nokia N97”.

Context: For each seed or candidate occurrence, the *context* is its set of surrounding words within a window of size w , i.e. we use w words right before the seed or the candidate and w words right after it. Stop words are removed.

For S-EM and Bayesian Sets, both the positive set P (based on the seeds set S) and the unlabeled candidate set U are generated differently. They are not represented as centroids.

Positive and unlabeled sets: For each seed $s_i \in S$, each occurrence in the corpus forms a vector as a positive example in P . The vector is formed based on the surrounding words context (see above) of the seed mention. Similarly, for each candidate $d \in D$ (see above; D denotes the set of all candidates), each occurrence also forms a vector as an unlabeled example in U . Thus, each unique seed or candidate entity may produce multiple feature vectors, depending on the number of times that it appears in the corpus.

The components in the feature vectors are term frequencies for S-EM as S-EM uses naïve Bayesian classification as its base classifier. For Bayesian Sets, they are 1's and 0's as Bayesian Sets only takes binary vectors based on whether a term occurs in the context or not.

4 Candidate Ranking

For distributional similarity, ranking is done using the similarity value of each candidate's centroid and the seeds' centroid (one centroid vector for all seeds). Rankings for S-EM and Bayesian Sets are more involved. We discuss them below.

After it ends, S-EM produces a Bayesian classifier C , which is used to classify each vector $u \in U$ and to assign a probability $p(+|u)$ to indicate the likelihood that u belongs to the positive class. Similarly, Bayesian Sets produces a score $score(u)$ for each u (not a probability).

Recall that for both S-EM and Bayesian Sets, each unique candidate entity may generate multiple feature vectors, depending on the number of times that the candidate entity occurs in the corpus. As such, the rankings produced by S-EM and Bayesian Sets are not the rankings of the entities, but rather the rankings of the entities' occurrences. Since different vectors representing the same candidate entity can have very different probabilities (for S-EM) or scores (for Bayesian Sets), we need to combine them and compute a single score for each unique candidate entity for ranking.

To this end, we also take the entity frequency into consideration. Typically, it is highly desirable to rank those correct and frequent entities at the top because they are more important than the infrequent ones in applications. With this in mind, we define a ranking method.

Let the probabilities (or scores) of a candidate entity $d \in D$ be $V_d = \{v_1, v_2, \dots, v_n\}$ for the n feature vectors of the candidate. Let M_d be the median of V_d . The final score (fs) for d is defined as:

$$fs(d) = M_d \times \log(1 + n) \quad (3)$$

The use of the median of V_d can be justified based on the statistical *skewness* (Neter *et al.* 1993). If the values in V_d are skewed towards the high side (negative skew), it means that the candidate entity is very likely to be a true entity, and we should take the median as it is also high (higher than the mean). However, if the skew is towards the low side (positive skew), it means that the candidate entity is unlikely to be a true entity and we should again use the median as it is low (lower than the mean) under this condition.

Note that here n is the frequency count of candidate entity d in the corpus. The constant 1 is added to smooth the value. The idea is to push the frequent candidate entities up by multiplying the logarithm of frequency. \log is taken in order to reduce the effect of big frequency counts.

The final score $fs(d)$ indicates candidate d 's overall likelihood to be a relevant entity. A high $fs(d)$ implies a high likelihood that d is in the expanded entity set. We can then rank all the candidates based on their $fs(d)$ values.

5 Experimental Evaluation

We empirically evaluate the three techniques in this section. We implemented *distribution similarity* and *Bayesian Sets*. S-EM was downloaded from <http://www.cs.uic.edu/~liub/S-EM/S-EM-download.html>. For both Bayesian Sets and S-EM, we used their default parameters. EM in S-EM ran only two iterations. For distributional similarity, we tested TF-IDF and PMI as feature values of vectors, and Cosine and Jaccard as similarity measures. Due to space limitations, we only show the results of the PMI and Cosine combination as it performed the best. This combination was also used in (Pantel *et al.*, 2009).

5.1 Corpora and Evaluation Metrics

We used 10 diverse corpora to evaluate the techniques. They were obtained from a commercial company. The data were crawled and extracted from multiple online message boards and blogs discussing different products and services. We split each message into sentences, and the sentences were POS-tagged using Brill's tagger (Brill, 1995). The tagged sentences were used to extract candidate entities and their contexts. Table 1 shows the domains and the number of sentences in each corpus, as well as the three seed entities used in our experiments for each corpus. The three seeds for each corpus were randomly selected from a set of common entities in the application domain.

Table 1. Descriptions of the 10 corpora

Domains	# Sentences	Seed Entities
Bank	17394	Citi, Chase, Wesabe
Blu-ray	7093	S300, Sony, Samsung
Car	2095	Honda, A3, Toyota
Drug	1504	Enbrel, Hurmia, Methotrexate
Insurance	12419	Cobra, Cigna, Kaiser
LCD	1733	PZ77U, Samsung, Sony
Mattress	13191	Simmons, Serta, Heavenly
Phone	14884	Motorola, Nokia, N95
Stove	25060	Kenmore, Frigidaire, GE
Vacuum	13491	Dc17, Hoover, Roomba

The regular evaluation metrics for named entity recognition such as precision and recall are not suitable for our purpose as we do not have the complete sets of gold standard entities to compare with. We adopt rank precision, which is commonly used for evaluation of entity set expansion techniques (Pantel *et al.*, 2009):

Precision @ N: The percentage of correct entities among the top N entities in the ranked list.

5.2 Experimental Results

The detailed experimental results for window size 3 ($w=3$) are shown in Table 2 for the 10 corpora. We present the precisions at the top 15-, 30- and 45-ranked positions (i.e., precisions @15, 30 and 45) for each corpus, with the average given in the last column. For distributional similarity, to save space Table 2 only shows the results of **Distr-Sim-freq**, which is the distributional similarity method with term frequency considered in the same way as for Bayesian Sets and S-EM, instead of the original distributional similarity, which is denoted by **Distr-Sim**. This is because on average, **Distr-Sim-freq** performs better than **Distr-Sim**. However, the summary results of both **Distr-Sim-freq** and **Distr-Sim** are given in Table 3.

From Table 2, we observe that on average **S-EM** outperforms **Distr-Sim-freq** by about 12 – 20% in terms of *Precision @ N*. **Bayesian-Sets** is also more accurate than **Distr-Sim-freq**, but **S-EM** outperforms **Bayesian-Sets** by 9 – 10%.

To test the sensitivity of window size w , we also experimented with $w = 6$ and $w = 9$. Due to space constraints, we present only their average results in Table 3. Again, we can see the same performance pattern as in Table 2 ($w = 3$): **S-EM** performs the best, **Bayesian-Sets** the second, and the two distributional similarity methods the third and the fourth, with **Distr-Sim-freq** slightly better than **Distr-Sim**.

Table 2. Precision @ top N (with 3 seeds, and window size $w = 3$)

	Bank	Blu-ray	Car	Drug	Insurance	LCD	Mattress	Phone	Stove	Vacuum	Avg.
Top 15											
Distr-Sim-freq	0.466	0.333	0.800	0.666	0.666	0.400	0.666	0.533	0.666	0.733	0.592
Bayesian-Sets	0.533	0.266	0.600	0.666	0.600	0.733	0.666	0.533	0.800	0.800	0.617
S-EM	0.600	0.733	0.733	0.733	0.533	0.666	0.933	0.533	0.800	0.933	0.720
Top 30											
Distr-Sim-freq	0.466	0.266	0.700	0.600	0.500	0.333	0.500	0.466	0.600	0.566	0.499
Bayesian-Sets	0.433	0.300	0.633	0.666	0.400	0.566	0.700	0.333	0.833	0.700	0.556
S-EM	0.500	0.700	0.666	0.666	0.566	0.566	0.733	0.600	0.600	0.833	0.643
Top 45											
Distr-Sim-freq	0.377	0.288	0.555	0.500	0.377	0.355	0.444	0.400	0.533	0.400	0.422
Bayesian-Sets	0.377	0.333	0.666	0.555	0.377	0.511	0.644	0.355	0.733	0.600	0.515
S-EM	0.466	0.688	0.644	0.733	0.533	0.600	0.644	0.555	0.644	0.688	0.620

Table 3. Average precisions over the 10 corpora of different window size (3 seeds)

Top Results	Window-size $w = 3$			Window-size $w = 6$			Window-size $w = 9$		
	Top 15	Top 30	Top 45	Top 15	Top 30	Top 45	Top 15	Top 30	Top 45
Distr-Sim	0.579	0.466	0.410	0.553	0.483	0.439	0.519	0.473	0.412
Distr-Sim-freq	0.592	0.499	0.422	0.553	0.492	0.441	0.559	0.476	0.410
Bayesian-Sets	0.617	0.556	0.515	0.593	0.539	0.524	0.539	0.522	0.497
S-EM	0.720	0.643	0.620	0.666	0.606	0.597	0.666	0.620	0.604

5.3 Why does S-EM Perform Better?

From the tables, we can see that both S-EM and Bayesian Sets performed better than distributional similarity. S-EM is better than Bayesian Sets. We believe that the reason is as follows: Distributional similarity does not use any information in the candidate set (or the unlabeled set U). It tries to rank the candidates solely through similarity comparisons with the given seeds (or positive cases). Bayesian Sets is better because it considers U . Its learning method produces a weight vector for features based on their occurrence differences in the positive set P and the unlabeled set U (Ghahramani and Heller 2005). This weight vector is then used to compute the final scores used in ranking. In this way, Bayesian Sets is able to exploit the useful information in U that was ignored by distributional similarity. S-EM also considers these differences in its NB classification; in addition, it uses the reliable negative set (RN) to help distinguish negative and positive cases, which both Bayesian Sets and distributional similarity do not do. We believe this balanced attempt by S-EM to distinguish the positive and negative cases is the reason for the better performance of S-EM. This raises an interesting question. Since Bayesian Sets is a ranking method and S-EM is a classification method, can we say even for ranking (our evaluation is based

on ranking) classification methods produce better results than ranking methods? Clearly, our single experiment cannot answer this question. But intuitively, classification, which separates positive and negative cases by pulling them towards two opposite directions, should perform better than ranking which only pulls the data in one direction. Further research on this issue is needed.

6 Conclusions and Future Work

Although distributional similarity is a classic technique for entity set expansion, this paper showed that PU learning performs considerably better on our diverse corpora. In addition, PU learning also outperforms Bayesian Sets (designed specifically for the task). In our future work, we plan to experiment with various other PU learning methods (Liu *et al.* 2003; Lee and Liu, 2003; Li *et al.* 2007; Elkan and Noto, 2008) on this entity set expansion task, as well as other tasks that were tackled using distributional similarity. In addition, we also plan to combine some syntactic patterns (Etzioni *et al.* 2005; Sarmiento *et al.* 2007) to further improve the results.

Acknowledgements: Bing Liu and Lei Zhang acknowledge the support of HP Labs Innovation Research Grant 2009-1062-1-A, and would like to thank Suk Hwan Lim and Eamonn O'Brien-Strain for many helpful discussions.

References

- Agirre, E., Alfonseca, E., Hall, K., Kravalova, J., Pasca, M., and Soroa, A. 2009. *A study on similarity and relatedness using distributional and WordNet-based approaches*. NAACL HLT.
- Blum, A. and Mitchell, T. 1998. *Combining labeled and unlabeled data with co-training*. In Proc. of Computational Learning Theory, pp. 92–100, 1998.
- Brill, E. 1995. *Transformation-Based error-Driven learning and natural language processing: a case study in part of speech tagging*. *Computational Linguistics*.
- Bunescu, R. and Mooney, R. 2004. *Collective information extraction with relational Markov Networks*. ACL.
- Cheng T., Yan X. and Chang C. K. 2007. *Entity-Rank: searching entities directly and holistically*. VLDB.
- Chieu, H.L. and Ng, H. Tou. 2002. *Name entity recognition: a maximum entropy approach using global information*. In The 6th Workshop on Very Large Corpora.
- Downey, D., Broadhead, M. and Etzioni, O. 2007. *Locating complex named entities in Web Text*. IJCAI.
- Elkan, C. and Noto, K. 2008. *Learning classifiers from only positive and unlabeled data*. KDD, 213-220.
- Etzioni, O., Cafarella, M., Downey, D., Popescu, A., Shaked, T., Soderland, S., Weld, D. Yates. 2005. *A. Unsupervised named-entity extraction from the Web: An Experimental Study*. *Artificial Intelligence*, 165(1):91-134.
- Ghahramani, Z and Heller, K.A. 2005. *Bayesian sets*. NIPS.
- Google Sets. 2008. *System and methods for automatically creating lists*. US Patent: US7350187, March 25.
- Gorman, J. and Curran, J. R. 2006. *Scaling distributional similarity to large corpora*. ACL.
- Harris, Z. Distributional Structure. 1985. In: Katz, J. J. (ed.), *The philosophy of linguistics*. Oxford University Press.
- Heller, K. and Ghahramani, Z. 2006. *A simple Bayesian framework for content-based image retrieval*. CVPR.
- Isozaki, H. and Kazawa, H. 2002. *Efficient support vector classifiers for named entity recognition*. COLING.
- Jiang, J. and Zhai, C. 2006. *Exploiting domain structure for named entity recognition*. HLT-NAACL.
- Lafferty J., McCallum A., and Pereira F. 2001. *Conditional random fields: probabilistic models for segmenting and labeling sequence data*. ICML.
- Lee, L. 1999. *Measures of distributional similarity*. ACL.
- Lee, W-S. and Liu, B. 2003. *Learning with Positive and Unlabeled Examples Using Weighted Logistic Regression*. ICML.
- Li, X., Liu, B. 2003. *Learning to classify texts using positive and unlabeled data*, IJCAI.
- Li, X., Liu, B., Ng, S. 2007. *Learning to identify unexpected instances in the test sSet*. IJCAI.
- Lin, D. 1998. *Automatic retrieval and clustering of similar words*. COLING/ACL.
- Liu, B, Lee, W-S, Yu, P. S, and Li, X. 2002. *Partially supervised text classification*. ICML, 387-394.
- Liu, B, Dai, Y., Li, X., Lee, W-S., and Yu. P. 2003. *Building text classifiers using positive and unlabeled examples*. ICDM, 179-188.
- Neter, J., Wasserman, W., and Whitmore, G. A. 1993. *Applied Statistics*. Allyn and Bacon.
- Nigam, K., McCallum, A., Thrun, S. and Mitchell, T. 2000. *Text classification from labeled and unlabeled documents using EM*. *Machine Learning*, 39(2/3), 103–134.
- Pantel, P., Eric Crestan, Arkady Borkovsky, Ana-Maria Popescu, Vishnu, Vyas. 2009. *Web-Scale Distributional similarity and entity set expansion*, EMNLP.
- Paşca, M. Lin, D. Bigham, J. Lifchits, A. Jain, A. 2006. *Names and similarities on the web: fast extraction in the fast lane*. ACL.
- Sarmiento, L., Jijkuon, V. de Rijke, M. and Oliveira, E. 2007. *“More like these”: growing entity classes from seeds*. CIKM.
- Wang, R. C. and Cohen, W. W. 2008. *Iterative set expansion of named entities using the web*. ICDM.
- Yu, H., Han, J., K. Chang. 2002. *PEBL: Positive example based learning for Web page classification using SVM*. KDD, 239-248.

Active Learning-Based Elicitation for Semi-Supervised Word Alignment

Vamshi Ambati, Stephan Vogel and Jaime Carbonell

{vamshi, vogel, jgc}@cs.cmu.edu

Language Technologies Institute, Carnegie Mellon University
5000 Forbes Avenue, Pittsburgh, PA 15213, USA

Abstract

Semi-supervised word alignment aims to improve the accuracy of automatic word alignment by incorporating full or partial manual alignments. Motivated by standard active learning query sampling frameworks like uncertainty-, margin- and query-by-committee sampling we propose multiple query strategies for the alignment link selection task. Our experiments show that by active selection of uncertain and informative links, we reduce the overall manual effort involved in elicitation of alignment link data for training a semi-supervised word aligner.

1 Introduction

Corpus-based approaches to machine translation have become predominant, with phrase-based statistical machine translation (PB-SMT) (Koehn et al., 2003) being the most actively progressing area. The success of statistical approaches to MT can be attributed to the IBM models (Brown et al., 1993) that characterize *word-level* alignments in parallel corpora. Parameters of these alignment models are learnt in an unsupervised manner using the EM algorithm over *sentence-level* aligned parallel corpora. While the ease of automatically aligning sentences at the word-level with tools like GIZA++ (Och and Ney, 2003) has enabled fast development of SMT systems for various language pairs, the quality of alignment is typically quite low for language pairs like Chinese-English, Arabic-English that diverge from the independence assumptions made by the generative models. Increased parallel data enables better estimation of the model parameters, but a large number of language pairs still lack such resources.

Two directions of research have been pursued for improving generative word alignment. The first is to relax or update the independence assumptions based on more information, usually syntactic, from the language pairs (Cherry and Lin, 2006; Fraser and Marcu, 2007a). The second is to use extra annotation, typically *word-level* human alignment for some sentence pairs, in conjunction with the parallel data to learn alignment in a semi-supervised manner. Our research is in the direction of the latter, and aims to reduce the effort involved in hand-generation of word alignments by using active learning strategies for careful selection of word pairs to seek alignment.

Active learning for MT has not yet been explored to its full potential. Much of the literature has explored one task – selecting sentences to translate and add to the training corpus (Haffari and Sarkar, 2009). In this paper we explore active learning for word alignment, where the input to the active learner is a sentence pair (S, T) and the annotation elicited from human is a set of links $\{a_{ij}, \forall s_i \in S, t_j \in T\}$. Unlike previous approaches, our work does not require elicitation of full alignment for the sentence pair, which could be effort-intensive. We propose active learning query strategies to selectively elicit partial alignment information. Experiments in Section 5 show that our selection strategies reduce alignment error rates significantly over baseline.

2 Related Work

Researchers have begun to explore models that use both labeled and unlabeled data to build word-alignment models for MT. Fraser and Marcu (2006) pose the problem of alignment as a search problem in log-linear space with features coming from the IBM alignment models. The log-

linear model is trained on available labeled data to improve performance. They propose a semi-supervised training algorithm which alternates between discriminative error training on the labeled data to learn the weighting parameters and maximum-likelihood EM training on unlabeled data to estimate the parameters. Callison-Burch et al. (2004) also improve alignment by interpolating human alignments with automatic alignments. They observe that while working with such data sets, alignments of higher quality should be given a much higher weight than the lower-quality alignments. Wu et al. (2006) learn separate models from labeled and unlabeled data using the standard EM algorithm. The two models are then interpolated to use as a learner in the semi-supervised algorithm to improve word alignment. To our knowledge, there is no prior work that has looked at reducing human effort by selective elicitation of partial word alignment using active learning techniques.

3 Active Learning for Word Alignment

Active learning attempts to optimize performance by selecting the most informative instances to label where ‘informativeness’ is defined as maximal expected improvement in accuracy. The objective is to select optimal instance for an external expert to label and then run the learning method on the newly-labeled and previously-labeled instances to minimize prediction or translation error, repeating until either the maximal number of external queries is reached or a desired accuracy level is achieved. Several studies (Tong and Koller, 2002; Nguyen and Smeulders, 2004; Donmez and Carbonell, 2008) show that active learning greatly helps to reduce the labeling effort in various classification tasks.

3.1 Active Learning Setup

We discuss our active learning setup for word alignment in Algorithm 1. We start with an unlabeled dataset $U = \{(S_k, T_k)\}$, indexed by k , and a seed pool of partial alignment links $A_0 = \{a_{ij}^k, \forall s_i \in S_k, t_j \in T_k\}$. This is usually an empty set at iteration $t = 0$. We iterate for T iterations. We take a pool-based active learning strategy, where we have access to all the automatically aligned links and we can score the links based on our active learning query strategy. The query strategy uses the automatically trained alignment

model M_t from current iteration t for scoring the links. Re-training and re-tuning an SMT system for each link at a time is computationally infeasible. We therefore perform batch learning by selecting a set of N links scored high by our query strategy. We seek manual corrections for the selected links and add the alignment data to the current labeled data set. The word-level aligned labeled data is provided to our semi-supervised word alignment algorithm for training an alignment model M_{t+1} over U .

Algorithm 1 AL FOR WORD ALIGNMENT

- 1: Unlabeled Data Set: $U = \{(S_k, T_k)\}$
 - 2: Manual Alignment Set : $A_0 = \{a_{ij}^k, \forall s_i \in S_k, t_j \in T_k\}$
 - 3: Train Semi-supervised Word Alignment using $(U, A_0) \rightarrow M_0$
 - 4: N : batch size
 - 5: **for** $t = 0$ to T **do**
 - 6: $L_t = \text{LinkSelection}(U, A_t, M_t, N)$
 - 7: Request Human Alignment for L_t
 - 8: $A_{t+1} = A_t + L_t$
 - 9: Re-train Semi-Supervised Word Alignment on $(U, A_{t+1}) \rightarrow M_{t+1}$
 - 10: **end for**
-

We can iteratively perform the algorithm for a defined number of iterations T or until a certain desired performance is reached, which is measured by alignment error rate (AER) (Fraser and Marcu, 2007b) in the case of word alignment. In a more typical scenario, since reducing human effort or cost of elicitation is the objective, we iterate until the available budget is exhausted.

3.2 Semi-Supervised Word Alignment

We use an extended version of MGIZA++ (Gao and Vogel, 2008) to perform the constrained semi-supervised word alignment. Manual alignments are incorporated in the EM training phase of these models as constraints that restrict the summation over all possible alignment paths. Typically in the EM procedure for IBM models, the training procedure requires for each source sentence position, the summation over all positions in the target sentence. The manual alignments allow for one-to-many alignments and many-to-many alignments in both directions. For each position i in the source sentence, there can be more than one manually aligned target word. The restricted training will allow only those paths, which are consistent with

the manual alignments. Therefore, the restriction of the alignment paths reduces to restricting the summation in EM.

4 Query Strategies for Link Selection

We propose multiple query selection strategies for our active learning setup. The scoring criteria is designed to select alignment links across sentence pairs that are highly uncertain under current automatic translation models. These links are difficult to align correctly by automatic alignment and will cause incorrect phrase pairs to be extracted in the translation model, in turn hurting the translation quality of the SMT system. Manual correction of such links produces the maximal benefit to the model. We would ideally like to elicit the least number of manual corrections possible in order to reduce the cost of data acquisition. In this section we discuss our link selection strategies based on the standard active learning paradigm of ‘uncertainty sampling’(Lewis and Catlett, 1994). We use the automatically trained translation model θ_t for scoring each link for uncertainty, which consists of bidirectional translation lexicon tables computed from the bidirectional alignments.

4.1 Uncertainty Sampling: Bidirectional Alignment Scores

The automatic Viterbi alignment produced by the alignment models is used to obtain translation lexicons. These lexicons capture the conditional distributions of source-given-target $P(s/t)$ and target-given-source $P(t/s)$ probabilities at the word level where $s_i \in S$ and $t_j \in T$. We define certainty of a link as the harmonic mean of the bidirectional probabilities. The selection strategy selects the least scoring links according to the formula below which corresponds to links with maximum uncertainty:

$$Score(a_{ij}/s_1^I, t_1^J) = \frac{2 * P(t_j/s_i) * P(s_i/t_j)}{P(t_j/s_i) + P(s_i/t_j)} \quad (1)$$

4.2 Confidence Sampling: Posterior Alignment probabilities

Confidence estimation for MT output is an interesting area with meaningful initial exploration (Blatz et al., 2004; Ueffing and Ney, 2007). Given a sentence pair (s_1^I, t_1^J) and its word alignment, we compute two confidence metrics at alignment link level – based on the posterior link probability as seen in Equation 5. We select the alignment

links that the initial word aligner is least confident according to our metric and seek manual correction of the links. We use $t2s$ to denote computation using higher order (IBM4) target-given-source models and $s2t$ to denote source-given-target models. Targeting some of the uncertain parts of word alignment has already been shown to improve translation quality in SMT (Huang, 2009). We use confidence metrics as an active learning sampling strategy to obtain most informative links. We also experimented with other confidence metrics as discussed in (Ueffing and Ney, 2007), especially the IBM 1 model score metric, but it did not show significant improvement in this task.

$$P_{t2s}(a_{ij}, t_1^J/s_1^I) = \frac{p_{t2s}(t_j/s_i, a_{ij} \in A)}{\sum_i^M p_{t2s}(t_j/s_i)} \quad (2)$$

$$P_{s2t}(a_{ij}, s_1^I/t_1^J) = \frac{p_{s2t}(s_i/t_j, a_{ij} \in A)}{\sum_i^N p_{s2t}(s_i/t_j)} \quad (3)$$

$$Conf1(a_{ij}/S, T) = \frac{2 * P_{t2s} * P_{s2t}}{P_{t2s} + P_{s2t}} \quad (4)$$

4.3 Query by Committee

The generative alignments produced differ based on the choice of direction of the language pair. We use A_{s2t} to denote alignment in the source to target direction and A_{t2s} to denote the target to source direction. We consider these alignments to be two experts that have two different views of the alignment process. We formulate our query strategy to select links where the agreement differs across these two alignments. In general query by committee is a standard sampling strategy in active learning(Freund et al., 1997), where the committee consists of any number of experts, in this case alignments, with varying opinions. We formulate a query by committee sampling strategy for word alignment as shown in Equation 6. In order to break ties, we extend this approach to select the link with higher average frequency of occurrence of words involved in the link.

$$Score(a_{ij}) = \alpha \quad (6)$$

$$where \quad \alpha = \begin{cases} 2 & a_{ij} \in A_{s2t} \cap A_{t2s} \\ 1 & a_{ij} \in A_{s2t} \cup A_{t2s} \\ 0 & otherwise \end{cases}$$

4.4 Margin Sampling

The strategy for confidence based sampling only considers information about the best scoring link

$conf(a_{ij}/S, T)$. However we could benefit from information about the second best scoring link as well. In typical multi-class classification problems, earlier work shows success using such a ‘margin based’ approach (Scheffer et al., 2001), where the difference between the probabilities assigned by the underlying model to the first best and second best labels is used as a sampling criteria. We adapt such a margin-based approach to link-selection using the *Conf1* scoring function discussed in the earlier sub-section. Our *margin* technique is formulated below, where \hat{a}_{1ij} and \hat{a}_{2ij} are potential first best and second best scoring alignment links for a word at position i in the source sentence S with translation T . The word with minimum margin value is chosen for human alignment. Intuitively such a word is a possible candidate for mis-alignment due to the inherent confusion in its target translation.

$$Margin(i) = Conf1(\hat{a}_{1ij}/S, T) - Conf1(\hat{a}_{2ij}/S, T)$$

5 Experiments

5.1 Data Setup

Our aim in this paper is to show that active learning can help select the most informative alignment links that have high uncertainty according to a given automatically trained model. We also show that fixing such alignments leads to the maximum reduction of error in word alignment, as measured by AER. We compare this with a baseline where links are selected at random for manual correction. To run our experiments iteratively, we automate the setup by using a parallel corpus for which the gold-standard human alignment is already available. We select the Chinese-English language pair, where we have access to 21,863 sentence pairs along with complete manual alignment.

5.2 Results

We first automatically align the Cn-En corpus using GIZA++ (Och and Ney, 2003). We then use the learned model in running our link selection algorithm over the entire corpus to determine the most uncertain links according to each active learning strategy. The links are then looked up in the gold-standard human alignment database and corrected. In case a link is not present in the gold-standard data, we introduce a NULL alignment, else we propose the alignment as given in

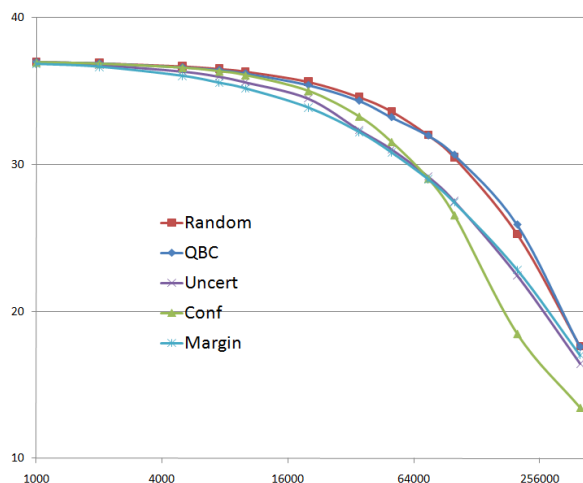


Figure 1: Performance of active sampling strategies for link selection

the gold standard. We select the partial alignment as a set of alignment links and provide it to our semi-supervised word aligner. We plot performance curves as number of links used in each iteration vs. the overall reduction of AER on the corpus.

Query by committee performs worse than random indicating that two alignments differing in direction are not sufficient in deciding for uncertainty. We will be exploring alternative formulations to this strategy. We observe that confidence based metrics perform significantly better than the baseline. From the scatter plots in Figure 1¹ we can say that using our best selection strategy one achieves similar performance to the baseline, but at a much lower cost of elicitation assuming cost per link is uniform.

We also perform end-to-end machine translation experiments to show that our improvement of alignment quality leads to an improvement of translation scores. For this experiment, we train a standard phrase-based SMT system (Koehn et al., 2007) over the entire parallel corpus. We tune on the MT-Eval 2004 dataset and test on a subset of MT-Eval 2004 dataset consisting of 631 sentences. We first obtain the baseline score where no manual alignment was used. We also train a configuration using gold standard manual alignment data for the parallel corpus. This is the maximum translation accuracy that we can achieve by any link selection algorithm. We now take the best link selection criteria, which is the confidence

¹X axis has number of links elicited on a log-scale

System	BLEU	METEOR
Baseline	18.82	42.70
Human Alignment	19.96	44.22
Active Selection 20%	19.34	43.25

Table 1: Alignment and Translation Quality

based method and train a system by only selecting 20% of all the links. We observe that at this point we have reduced the AER from 37.09 AER to 26.57 AER. The translation accuracy as measured by BLEU (Papineni et al., 2002) and METEOR (Lavie and Agarwal, 2007) also shows improvement over baseline and approaches gold standard quality. Therefore we achieve 45% of the possible improvement by only using 20% elicitation effort.

5.3 Batch Selection

Re-training the word alignment models after eliciting every individual alignment link is infeasible. In our data set of 21,863 sentences with 588,075 links, it would be computationally intensive to re-train after eliciting even 100 links in a batch. We therefore sample links as a discrete batch, and train alignment models to report performance at fixed points. Such a batch selection is only going to be sub-optimal as the underlying model changes with every alignment link and therefore becomes ‘stale’ for future selections. We observe that in some scenarios while fixing one alignment link could potentially fix all the mis-alignments in a sentence pair, our batch selection mechanism still samples from the rest of the links in the sentence pair. We experimented with an exponential decay function over the number of links previously selected, in order to discourage repeated sampling from the same sentence pair. We performed an experiment by selecting one of our best performing selection strategies (*conf*) and ran it in both configurations - one with the decay parameter (*batchdecay*) and one without it (*batch*). As seen in Figure 2, the decay function has an effect in the initial part of the curve where sampling is sparse but the effect gradually fades away as we observe more samples. In the reported results we do not use batch decay, but an optimal estimation of ‘staleness’ could lead to better gains in batch link selection using active learning.

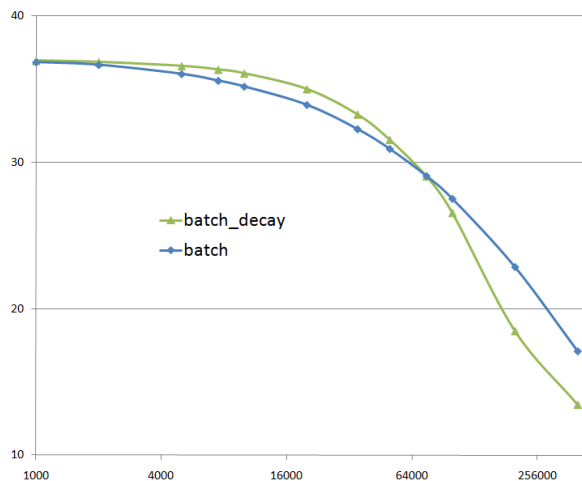


Figure 2: Batch decay effects on Conf-posterior sampling strategy

6 Conclusion and Future Work

Word-Alignment is a particularly challenging problem and has been addressed in a completely unsupervised manner thus far (Brown et al., 1993). While generative alignment models have been successful, lack of sufficient data, model assumptions and local optimum during training are well known problems. Semi-supervised techniques use partial manual alignment data to address some of these issues. We have shown that active learning strategies can reduce the effort involved in eliciting human alignment data. The reduction in effort is due to careful selection of maximally uncertain links that provide the most benefit to the alignment model when used in a semi-supervised training fashion. Experiments on Chinese-English have shown considerable improvements. In future we wish to work with word alignments for other language pairs like Arabic and English. We have tested out the feasibility of obtaining human word alignment data using Amazon Mechanical Turk and plan to obtain more data reduce the cost of annotation.

Acknowledgments

This research was partially supported by DARPA under grant NBCHC080097. Any opinions, findings, and conclusions expressed in this paper are those of the authors and do not necessarily reflect the views of the DARPA. The first author would like to thank Qin Gao for the semi-supervised word alignment software and help with running experiments.

References

- John Blatz, Erin Fitzgerald, George Foster, Simona Gandrabur, Cyril Goutte, Alex Kulesza, Alberto Sanchis, and Nicola Ueffing. 2004. Confidence estimation for machine translation. In *Proceedings of Coling 2004*, pages 315–321, Geneva, Switzerland, Aug 23–Aug 27. COLING.
- Peter F. Brown, Vincent J. Della Pietra, Stephen A. Della Pietra, and Robert L. Mercer. 1993. The mathematics of statistical machine translation: parameter estimation. *Computational Linguistics*, 19(2):263–311.
- Chris Callison-Burch, David Talbot, and Miles Osborne. 2004. Statistical machine translation with word- and sentence-aligned parallel corpora. In *ACL 2004*, page 175, Morristown, NJ, USA. Association for Computational Linguistics.
- Colin Cherry and Dekang Lin. 2006. Soft syntactic constraints for word alignment through discriminative training. In *Proceedings of the COLING/ACL on Main conference poster sessions*, pages 105–112, Morristown, NJ, USA.
- Pinar Donmez and Jaime G. Carbonell. 2008. Optimizing estimated loss reduction for active sampling in rank learning. In *ICML '08: Proceedings of the 25th international conference on Machine learning*, pages 248–255, New York, NY, USA. ACM.
- Alexander Fraser and Daniel Marcu. 2006. Semi-supervised training for statistical word alignment. In *ACL-44: Proceedings of the 21st International Conference on Computational Linguistics and the 44th annual meeting of the Association for Computational Linguistics*, pages 769–776, Morristown, NJ, USA. Association for Computational Linguistics.
- Alexander Fraser and Daniel Marcu. 2007a. Getting the structure right for word alignment: LEAF. In *Proceedings of the 2007 Joint Conference on EMNLP-CoNLL*, pages 51–60.
- Alexander Fraser and Daniel Marcu. 2007b. Measuring word alignment quality for statistical machine translation. *Comput. Linguist.*, 33(3):293–303.
- Yoav Freund, Sebastian H. Seung, Eli Shamir, and Naftali Tishby. 1997. Selective sampling using the query by committee algorithm. *Machine Learning.*, 28(2-3):133–168.
- Qin Gao and Stephan Vogel. 2008. Parallel implementations of word alignment tool. In *Software Engineering, Testing, and Quality Assurance for Natural Language Processing*, pages 49–57, Columbus, Ohio, June. Association for Computational Linguistics.
- Gholamreza Haffari and Anoop Sarkar. 2009. Active learning for multilingual statistical machine translation. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, pages 181–189, Suntec, Singapore, August. Association for Computational Linguistics.
- Fei Huang. 2009. Confidence measure for word alignment. In *Proceedings of the Joint ACL and IJCNLP*, pages 932–940, Suntec, Singapore, August. Association for Computational Linguistics.
- Philipp Koehn, Franz Josef Och, and Daniel Marcu. 2003. Statistical phrase-based translation. In *Proc. of the HLT/NAACL*, Edmonton, Canada.
- Philipp Koehn, Hieu Hoang, Alexandra Birch Mayne, Christopher Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondrej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In *ACL Demonstration Session*.
- Alon Lavie and Abhaya Agarwal. 2007. Meteor: an automatic metric for mt evaluation with high levels of correlation with human judgments. In *WMT 2007*, pages 228–231, Morristown, NJ, USA.
- David D. Lewis and Jason Catlett. 1994. Heterogeneous uncertainty sampling for supervised learning. In *In Proceedings of the Eleventh International Conference on Machine Learning*, pages 148–156. Morgan Kaufmann.
- Hieu T. Nguyen and Arnold Smeulders. 2004. Active learning using pre-clustering. In *ICML*.
- Franz Josef Och and Hermann Ney. 2003. A systematic comparison of various statistical alignment models. *Computational Linguistics*, pages 19–51.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *ACL 2002*, pages 311–318, Morristown, NJ, USA.
- Tobias Scheffer, Christian Decomain, and Stefan Wrobel. 2001. Active hidden markov models for information extraction. In *IDA '01: Proceedings of the 4th International Conference on Advances in Intelligent Data Analysis*, pages 309–318, London, UK. Springer-Verlag.
- Simon Tong and Daphne Koller. 2002. Support vector machine active learning with applications to text classification. *Journal of Machine Learning*, pages 45–66.
- Nicola Ueffing and Hermann Ney. 2007. Word-level confidence estimation for machine translation. *Comput. Linguist.*, 33(1):9–40.
- Hua Wu, Haifeng Wang, and Zhanyi Liu. 2006. Boosting statistical word alignment using labeled and unlabeled data. In *Proceedings of the COLING/ACL on Main conference poster sessions*, pages 913–920, Morristown, NJ, USA. Association for Computational Linguistics.

An Active Learning Approach to Finding Related Terms

David Vickrey

Stanford University

dvickrey@cs.stanford.edu

Oscar Kipersztok

Boeing Research & Technology

oscar.kipersztok@boeing.com

Daphne Koller

Stanford University

koller@cs.stanford.edu

Abstract

We present a novel system that helps non-experts find sets of similar words. The user begins by specifying one or more seed words. The system then iteratively suggests a series of candidate words, which the user can either accept or reject. Current techniques for this task typically bootstrap a classifier based on a fixed seed set. In contrast, our system involves the user throughout the labeling process, using active learning to intelligently explore the space of similar words. In particular, our system can take advantage of negative examples provided by the user. Our system combines multiple pre-existing sources of similarity data (a standard thesaurus, WordNet, contextual similarity), enabling it to capture many types of similarity groups (“synonyms of crash,” “types of car,” etc.). We evaluate on a hand-labeled evaluation set; our system improves over a strong baseline by 36%.

1 Introduction

Set expansion is a well-studied NLP problem where a machine-learning algorithm is given a fixed set of *seed words* and asked to find additional members of the implied set. For example, given the seed set {“elephant,” “horse,” “bat”}, the algorithm is expected to return other mammals. Past work, e.g. (Roark & Charniak, 1998; Ghahramani & Heller, 2005; Wang & Cohen, 2007; Pantel et al., 2009), generally focuses on semi-automatic acquisition of the remaining members of the set by mining large amounts of unlabeled data.

State-of-the-art set expansion systems work well for well-defined sets of nouns, e.g. “US Presidents,” particularly when given a large seed set. Set expansion is more difficult with fewer seed words and for other kinds of sets. The seed words may have multiple senses and the user may have in mind a variety of attributes that the answer must match. For example, suppose the seed word is

“jaguar”. First, there is sense ambiguity; we could be referring to either a “large cat” or a “car.” Beyond this, we might have in mind various more (or less) specific groups: “Mexican animals,” “predators,” “luxury cars,” “British cars,” etc.

We propose a system which addresses several shortcomings of many set expansion systems. First, these systems can be difficult to use. As explored by Vyas et al. (2009), non-expert users produce seed sets that lead to poor quality expansions, for a variety of reasons including ambiguity and lack of coverage. Even for expert users, constructing seed sets can be a laborious and time-consuming process. Second, most set expansion systems do not use negative examples, which can be very useful for weeding out other bad answers. Third, many set expansion systems concentrate on noun classes such as “US Presidents” and are not effective or do not apply to other kinds of sets.

Our system works as follows. The user initially thinks of at least one seed word belonging to the desired set. One at a time, the system presents candidate words to the user and asks whether the candidate fits the concept. The user’s answer is fed back into the system, which takes into account this new information and presents a new candidate to the user. This continues until the user is satisfied with the compiled list of “Yes” answers. Our system uses both positive and negative examples to guide the search, allowing it to recover from initially poor seed words. By using multiple sources of similarity data, our system captures a variety of kinds of similarity. Our system replaces the potentially difficult problem of thinking of many seed words with the easier task of answering yes/no questions. The downside is a possibly increased amount of user interaction (although standard set expansion requires a non-trivial amount of user interaction to build the seed set).

There are many practical uses for such a system. Building a better, more comprehensive thesaurus/gazetteer is one obvious application. Another application is in high-precision query expansion, where a human manually builds a list of ex-

pansion terms. Suppose we are looking for pages discussing “public safety.” Then synonyms (or near-synonyms) of “safety” would be useful (e.g. “security”) but also non-synonyms such as “precautions” or “prevention” are also likely to return good results. In this case, the concept we are interested in is “Words which imply that safety is being discussed.” Another interesting direction not pursued in this paper is using our system as part of a more-traditional set expansion system to build seed sets more quickly.

2 Set Expansion

As input, we are provided with a small set of seed words s . The desired output is a target set of words G , consisting of all words that fit the desired concept. A particular seed set s can belong to many possible goal sets G , so additional information may be required to do well.

Previous work tries to do as much as possible using only s . Typically s is assumed to contain at least 2 words and often many more. Pantel et al. (2009) discusses the issue of seed set size in detail, concluding that 5-20 seed words are often required for good performance.

There are several problems with the fixed seed set approach. It is not always easy to think of even a single additional seed word (e.g., the user is trying to find “German automakers” and can only think of “Volkswagen”). Even if the user can think of additional seed words, time and effort might be saved by using active learning to find good suggestions. Also, as Vyas et al. (2009) show, non-expert users often produce poor-quality seed sets.

3 Active Learning System

Any system for this task relies on information about similarity between words. Our system takes as input a rectangular matrix M . Each column corresponds to a particular word. Each row corresponds to a unique *dimension of similarity*; the j^{th} entry in row i m_{ij} is a number between 0 and 1 indicating the degree to which w_j belongs to the i^{th} similarity group. Possible similarity dimensions include “How similar is word w_j to the verb jump?” “Is w_j a type of cat?” and “Are the words which appear in the context of w_j similar to those that appear in the context of boat?” Each row r_i of M is labeled with a word l_i . This may follow intuitively from the similarity axis (e.g., “jump,” “cat,” and “boat”, respectively), or it can be generated automatically (e.g. the word w_j with the highest membership m_{ij}).

Let θ be a vector of weights, one per row, which

correspond to how well each row aligns with the goal set G . Thus, θ_i should be large and positive if row i has large entries for positive but not negative examples; and it should be large and negative if row i has large entries for negative but not positive examples. Suppose that we have already chosen an appropriate weight vector θ . We wish to rank all possible words (i.e., the columns of M) so that the most promising word gets the highest score. A natural way to generate a score z_j for column j is to take the dot product of θ with column j , $z_j = \sum_i \theta_i m_{ij}$. This rewards word w_j for having high membership in rows with positive θ , and low membership in rows with negative θ .

Our system uses a “batch” approach to active learning. At iteration i , it chooses a new θ based on all data labeled so far (for the 1^{st} iteration, this data consists of the seed set s). It then chooses the column (word) with the highest score (among words not yet labeled) as the candidate word w_i . The user answers “Yes” or “No,” indicating whether or not w_i belongs to G . w_i is added to the positive set \mathbf{p} or the negative set \mathbf{n} based on the user’s answer. Thus, we have a labeled data set that grows from iteration to iteration as the user labels each candidate word. Unlike set expansion, this procedure generates (and uses) both positive and negative examples.

We explore two options for choosing θ . Recall that each row i is associated with a label l_i . The first method is to set $\theta_i = 1$ if $l_i \in \mathbf{p}$ (that is, the set of positively labeled words includes label l_i), $\theta_i = -1$ if $l_i \in \mathbf{n}$, and $\theta_i = 0$ otherwise. We refer to this method as “Untrained”, although it is still *adaptive* — it takes into account the labeled examples the user has provided so far.

The second method uses a standard machine learning algorithm, logistic regression. As before, the final ranking over words is based on the score z_j . However, z_j is passed through the logistic function to produce a score between 0 and 1, $z'_j = \frac{1}{1+e^{-z_j}}$. We can interpret this score as the probability that w_j is a positive example, $P_\theta(Y|w_j)$. This leads to the objective function

$$L(\theta) = \log\left(\prod_{w_j \in \mathbf{p}} P_\theta(Y|w_j) \prod_{w_j \in \mathbf{n}} (1-P_\theta(Y|w_j))\right).$$

This objective is convex and can be optimized using standard methods such as L-BFGS (Liu & Nocedal, 1989). Following standard practice we add an L_2 regularization term $-\frac{\theta^T \theta}{2\sigma^2}$ to the objective. This method does not use the row labels l_i .

Data	Word	Similar words
Moby	arrive	accomplish, achieve, achieve success, advance, appear, approach, arrive at, arrive in, attain,...
WordNet	factory	(plant,-1.9);(arsenal,-2.8);(mill,-2.9);(sweatshop,-4.1);(refinery,-4.2);(winery,-4.5);...
DistSim	watch	(jewelry,.137),(wristwatch,.115),(shoe,0.09),(appliance,0.09),(household appliance,0.089),...

Table 1: Examples of unprocessed similarity entries from each data source.

4 Data Sources

We consider three similarity data sources: the Moby thesaurus¹, WordNet (Fellbaum, 1998), and distributional similarity based on a large corpus of text (Lin, 1998). Table 1 shows similarity lists from each. These sources capture different kinds of similarity information, which increases the representational power of our system. For all sources, the similarity of a word with itself is set to 1.0.

It is worth noting that our system is not strictly limited to choosing from pre-existing groups. For example, if we have a list of luxury items, and another list of cars, our system can learn weights so that it prefers items in the intersection, luxury cars.

Moby thesaurus consists of a list of word-based thesaurus entries. Each word w_i has a list of similar words sim_j^i . Moby has a total of about 2.5 million related word pairs. Unlike some other thesauri (such as WordNet and thesaurus.com), entries are not broken down by word sense.

In the raw format, the similarity relation is not symmetric; for example, there are many words that occur only in similarity lists but do not have their own entries. We augmented the thesaurus to make it symmetric: if “dog” is in the similarity entry for “cat,” we add “cat” to the similarity entry for “dog” (creating an entry for “dog” if it does not exist yet). We then have a row i for every similarity entry in the augmented thesaurus; m_{ij} is 1 if w_j appears in the similarity list of w_i , and 0 otherwise. The label l_i of row i is simply word w_i . Unlike some other thesauri (including WordNet and thesaurus.com), the entries are not broken down by word sense or part of speech. For polysemic words, there will be a mix of the words similar to each sense and part of speech.

WordNet is a well-known dictionary/thesaurus/ontology often used in NLP applications. It consists of a large number of synsets; a synset is a set of one or more similar word senses. The synsets are then connected with hypernym/hyponym links, which represent IS-A relationships. We focused on measuring similarity in WordNet using the hypernym hierarchy.² There are many methods for

converting this hierarchy into a similarity score; we chose to use the Jiang-Conrath distance (Jiang & Conrath, 1997) because it tends to be more robust to the exact structure of WordNet. The number of types of similarity in WordNet tends to be less than that captured by Moby, because synsets in WordNet are (usually) only allowed to have a single parent. For example, “murder” is classified as a type of killing, but not as a type of crime.

The Jiang-Conrath distance gives scores for pairs of word senses, not pairs of words. We handle this by adding one row for every word sense with the right part of speech (rather than for every word); each row measures the similarity of every word to a particular word sense. The label of each row is the (undisambiguated) word; multiple rows can have the same label. For the columns, we do need to collapse the word senses into words; for each word, we take a maximum across all of its senses. For example, to determine how similar (the only sense of) “factory” is to the word “plant,” we compute the similarity of “factory” to the “industrial plant” sense of “plant” and to the “living thing” sense of “plant” and take the higher of the two (in this case, the former).

The Jiang-Conrath distance is a number between $-\infty$ and 0. By examination, we determined that scores below -12.0 indicate virtually no similarity. We cut off scores below this point and linearly mapped each score x to the range 0 to 1, yielding a final similarity of $\frac{\min(0,x+12)}{12}$. This greatly sparsified the similarity matrix M .

Distributional similarity. We used Dekang Lin’s dependency-based thesaurus, available at www.cs.ualberta.ca/~lindek/downloads.htm. This resource groups words based on the words they co-occur with in normal text. The words most similar to “cat” are “dog,” “animal,” and “monkey,” presumably because they all “eat,” “walk,” etc. Like Moby, similarity entries are not divided by word sense; usually, only the dominant sense of each word is represented. This type of similarity is considerably different from the other two types, tending to focus less on minor details and more on broad patterns.

Each similarity entry corresponds to a single

¹Available at icon.shef.ac.uk/Moby/.

²A useful similarity metric we did not explore in this paper is similarity between WordNet dictionary definitions

word w_i and is a list of *scored* similar words sim_j^i . The scores vary between 0 and 1, but usually the highest-scored word in a similarity list gets a score of no more than 0.3. To calibrate these scores with the previous two types, we divided all scores by the score of the highest-scored word in that list. Since each row is normalized individually, the similarity matrix M is not symmetric. Also, there are separate similarity lists for each of nouns, verbs, and modifiers; we only used the lists matching the seed word’s part of speech.

5 Experimental Setup

Given a seed set s and a complete target set G , it is easy to evaluate our system; we say “Yes” to anything in G , “No” to everything else, and see how many of the candidate words are in G . However, building a complete gold-standard G is in practice prohibitively difficult; instead, we are only capable of saying whether or not a word belongs to G when presented with that word.

To evaluate a particular active learning algorithm, we can just run the algorithm manually, and see how many candidate words we say “Yes” to (note that this will not give us an accurate estimate of the recall of our algorithm). Evaluating several different algorithms for the same s and G is more difficult. We could run each algorithm separately, but there are several problems with this approach. First, we might unconsciously (or consciously) bias the results in favor of our preferred algorithms. Second, it would be fairly difficult to be consistent across multiple runs. Third, it would be inefficient, since we would label the same words multiple times for different algorithms.

We solved this problem by building a labeling system which runs all algorithms that we wish to test in parallel. At each step, we pick a random algorithm and either present its current candidate to the user or, if that candidate has already been labeled, we supply that algorithm with the given answer. We do NOT ever give an algorithm a labeled training example unless it actually asks for it – this guarantees that the combined system is equivalent to running each algorithm separately. This procedure has the property that the user cannot tell which algorithms presented which words.

To evaluate the relative contribution of active learning, we consider a version of our system where active learning is disabled. Instead of re-training the system every iteration, we train it once on the seed set s and keep the weight vector θ fixed from iteration to iteration.

We evaluated our algorithms along three axes. First, the method for choosing θ : Untrained and Logistic (U and L). Second, the data sources used: each source separately (M for Moby, W for WordNet, D for distributional similarity), and all three in combination (MWD). Third, whether active learning is used (+/-). Thus, logistic regression using Moby and no active learning is L(M,-). For logistic regression, we set the regularization penalty σ^2 to 1, based on qualitative analysis during development (before seeing the test data).

We also compared the performance of our algorithms to the popular online thesaurus <http://thesaurus.com>. The entries in this thesaurus are similar to Moby, except that each word may have multiple sense-disambiguated entries. For each seed word w , we downloaded the page for w and extracted a set of synonyms entries for that word. To compare fairly with our algorithms, we propose a word-by-word method for exploring the thesaurus, intended to model a user scanning the thesaurus. This method checks the first 3 words from each entry; if none of these are labeled “Yes,” it moves on to the next entry. We omit details for lack of space.

6 Experimental Results

We designed a test set containing different types of similarity. Table 2 shows each category, with examples of specific similarity queries. For each type, we tested on five different queries. For each query, the first author built the seed set by writing down the first three words that came to mind. For most queries this was easy. However, for the similarity type Hard Synonyms, coming up with more than one seed word was considerably more difficult. To build seed sets for these queries, we ran our evaluation system using a single seed word and took the first two positive candidates; this ensured that we were not biasing our seed set in favor of a particular algorithm or data set.

For each query, we ran our evaluation system until each algorithm had suggested 25 candidate words, for a total of 625 labeled words per algorithm. We measured performance using mean average precision (MAP), which corresponds to area under the precision-recall curve. It gives an overall assessment across different stopping points.

Table 3 shows results for an informative subset of the tested algorithms. There are many conclusions we can draw. Thesaurus.Com performs poorly overall; our best system, L(MWD,+), outscores it by 164%. The next group of al-

Category Name	Example Similarity Queries
Simple Groups (SG)	car brands, countries, mammals, crimes
Complex Groups (CG)	luxury car brands, sub-Saharan countries
Synonyms (Syn)	syn of {scandal, helicopter, arrogant, slay}
Hard Synonyms (HS)	syn of {(stock-market) crash, (legal) maneuver}
Meronym/Material (M)	parts of a car, things made of wood

Table 2: Categories and examples

	SG	CG	Syn	HS	M
Thesaurus.Com	.041	.060	.275	.173	.060
L(D,+)	.377	.344	.211	.329	.177
L(M,-)	.102	.118	.393	.279	.119
U(W,+)	.097	.136	.296	.277	.165
U(MWD,+)	.194	.153	.438	.357	.213
L(MWD,-)	.344	.207	.360	.345	.173
L(MWD,+)	.366	.335	.379	.372	.158

Table 4: Results by category

gorithms, U(*,-), add together the similarity entries of the seed words for a particular similarity source. The best of these uses distributional similarity; L(MWD,+) outcores it by 53%. Combining all similarity types, U(MWD,-) improves by 10% over U(D,-). L(MWD,+) improves over the best single-source, L(D,+), by a similar margin.

Using logistic regression instead of the untrained weights significantly improves performance. For example, L(MWD,+) outcores U(MWD,+) by 19%. Using active learning also significantly improves performance: L(MWD,+) outcores L(MWD,-) by 13%. This shows that active learning is useful even when a reasonable amount of initial information is available (three seed words for each test case). The gains from logistic regression and active learning are cumulative; L(MWD,+) outcores U(MWD,-) by 38%.

Finally, our best system, L(MWD,+) improves over L(D,-), the best system using a single data source and no active learning, by 36%. We consider L(D,-) to be a strong baseline; this comparison demonstrates the usefulness of the main contributions of this paper, the use of multiple data sources and active learning. L(D,-) is still fairly sophisticated, since it combines information from the similarity entries for different words.

Table 4 shows the breakdown of results by category. For this chart, we chose the best setting for each similarity type. Broadly speaking, the thesauri work reasonably well for synonyms, but poorly for groups. Meronyms were difficult

Algorithm	MAP
Thesaurus.Com	.122
U(M,-)	.176
U(W,-)	.182
U(D,-)	.211
L(D,-)	.236
L(D,+)	.288
U(MWD,-)	.233
U(MWD,+)	.271
L(MWD,-)	.286
L(MWD,+)	.322

Table 3: Comparison of algorithms

across the board. Neither logistic regression nor active learning always improved performance, but L(MWD,+) performs near the top for every category. The complex groups category is particularly interesting, because achieving high performance on this category required using both logistic regression and active learning. This makes sense since negative evidence is particularly important for this category.

7 Discussion and Related Work

The biggest difference between our system and previous work is the use of active learning, especially in allowing the use of negative examples. Most previous set expansion systems use bootstrapping from a small set of positive examples. Recently, the use of negative examples for set expansion was proposed by Vyas and Pantel (2009), although in a different way. First, set expansion is run as normal using a fixed seed set. Then, human annotators label a small number of negative examples from the returned results, which are used to weed out other bad answers. Our method incorporates negative examples at an earlier stage. Also, we use a logistic regression model to robustly incorporate negative information, rather than deterministically ruling out words and features.

Our system is limited by our data sources. Suppose we want actors who appeared in Star Wars. If we only know that Harrison Ford and Mark Hamill are actors, we have little to go on. There has been a large amount of work on other sources of word-similarity. Hughes and Ramage (2007) use random walks over WordNet, incorporating information such as meronymy and dictionary glosses. Snow et al. (2006) extract hypernyms from free text. Wang and Cohen (2007) exploit web-page structure, while Pasca and Durme (2008) examine query logs. We expect that adding these types of data would significantly improve our system.

References

- Fellbaum, C. (Ed.). (1998). *Wordnet: An electronic lexical database*. MIT Press.
- Ghahramani, Z., & Heller, K. (2005). Bayesian sets. *Advances in Neural Information Processing Systems (NIPS)*.
- Hughes, T., & Ramage, D. (2007). Lexical semantic relatedness with random graph walks. *EMNLP-CoNLL*.
- Jiang, J., & Conrath, D. (1997). Semantic similarity based on corpus statistics and lexical taxonomy. *Proceedings of International Conference on Research in Computational Linguistics*.
- Lin, D. (1998). An information-theoretic definition of similarity. *Proceedings of ICML*.
- Liu, D. C., & Nocedal, J. (1989). On the limited memory method for large scale optimization. *Mathematical Programming B*.
- Pantel, P., Crestan, E., Borkovsky, A., Popescu, A., & Vyas, V. (2009). Web-scale distributional similarity and entity set expansion. *EMNLP*.
- Pasca, M., & Durme, B. V. (2008). Weakly-supervised acquisition of open-domain classes and class attributes from web documents and query logs. *ACL*.
- Roark, B., & Charniak, E. (1998). Noun-phrase co-occurrence statistics for semiautomatic semantic lexicon construction. *ACL-COLING*.
- Snow, R., Jurafsky, D., & Ng, A. (2006). Semantic taxonomy induction from heterogeneous evidence. *ACL*.
- Vyas, V., & Pantel, P. (2009). Semi-automatic entity set refinement. *NAACL/HLT*.
- Vyas, V., Pantel, P., & Crestan, E. (2009). Helping editors choose better seed sets for entity expansion. *CIKM*.
- Wang, R., & Cohen, W. (2007). Language-independent set expansion of named entities using the web. *Seventh IEEE International Conference on Data Mining*.

Learning Better Data Representation using Inference-Driven Metric Learning

Paramveer S. Dhillon **Partha Pratim Talukdar*** **Koby Crammer**
CIS Deptt., Univ. of Penn. Search Labs, Microsoft Research Deptt. of Electrical Engg.
Philadelphia, PA, U.S.A Mountain View, CA, USA The Technion, Haifa, Israel
dhillon@cis.upenn.edu partha@talukdar.net koby@ee.technion.ac.il

Abstract

We initiate a study comparing effectiveness of the transformed spaces learned by recently proposed *supervised*, and *semi-supervised* metric learning algorithms to those generated by previously proposed *unsupervised* dimensionality reduction methods (e.g., PCA). Through a variety of experiments on different real-world datasets, we find IDML-IT, a *semi-supervised* metric learning algorithm to be the most effective.

1 Introduction

Because of the high-dimensional nature of NLP datasets, estimating a large number of parameters (a parameter for each dimension), often from a limited amount of labeled data, is a challenging task for statistical learners. Faced with this challenge, various *unsupervised* dimensionality reduction methods have been developed over the years, e.g., Principal Components Analysis (PCA).

Recently, several supervised metric learning algorithms have been proposed (Davis et al., 2007; Weinberger and Saul, 2009). IDML-IT (Dhillon et al., 2010) is another such method which exploits labeled as well as unlabeled data during metric learning. These methods learn a Mahalanobis distance metric to compute distance between a pair of data instances, which can also be interpreted as learning a transformation of the input data, as we shall see in Section 2.1.

In this paper, we make the following contributions:

Even though different supervised and semi-supervised metric learning algorithms have recently been proposed, effectiveness of the transformed spaces learned by them in NLP

datasets has not been studied before. In this paper, we address that gap: we compare effectiveness of classifiers trained on the transformed spaces learned by metric learning methods to those generated by previously proposed *unsupervised* dimensionality reduction methods. We find IDML-IT, a *semi-supervised* metric learning algorithm to be the most effective.

2 Metric Learning

2.1 Relationship between Metric Learning and Linear Projection

We first establish the well-known equivalence between learning a Mahalanobis distance measure and Euclidean distance in a linearly transformed space of the data (Weinberger and Saul, 2009). Let A be a $d \times d$ positive definite matrix which parameterizes the Mahalanobis distance, $d_A(x_i, x_j)$, between instances x_i and x_j , as shown in Equation 1. Since A is positive definite, we can decompose it as $A = P^\top P$, where P is another matrix of size $d \times d$.

$$\begin{aligned} d_A(x_i, x_j) &= (x_i - x_j)^\top A (x_i - x_j) \quad (1) \\ &= (Px_i - Px_j)^\top (Px_i - Px_j) \\ &= d_{\text{Euclidean}}(Px_i, Px_j) \end{aligned}$$

Hence, computing Mahalanobis distance parameterized by A is equivalent to first projecting the instances into a new space using an appropriate transformation matrix P and then computing Euclidean distance in the linearly transformed space. In this paper, we are interested in learning a better representation of the data (i.e., projection matrix P), and we shall achieve that goal by learning the corresponding Mahalanobis distance parameter A .

We shall now review two recently proposed metric learning algorithms.

* Research carried out while at the University of Pennsylvania, Philadelphia, PA, USA.

2.2 Information-Theoretic Metric Learning (ITML): Supervised

Information-Theoretic Metric Learning (ITML) (Davis et al., 2007) assumes the availability of prior knowledge about inter-instance distances. In this scheme, two instances are considered similar if the Mahalanobis distance between them is upper bounded, i.e., $d_A(x_i, x_j) \leq u$, where u is a non-trivial upper bound. Similarly, two instances are considered dissimilar if the distance between them is larger than certain threshold l , i.e., $d_A(x_i, x_j) \geq l$. Similar instances are represented by set S , while dissimilar instances are represented by set D .

In addition to prior knowledge about inter-instance distances, sometimes prior information about the matrix A , denoted by A_0 , itself may also be available. For example, Euclidean distance (i.e., $A_0 = I$) may work well in some domains. In such cases, we would like the learned matrix A to be as close as possible to the prior matrix A_0 . ITML combines these two types of prior information, i.e., knowledge about inter-instance distances, and prior matrix A_0 , in order to learn the matrix A by solving the optimization problem shown in (2).

$$\begin{aligned} \min_{A \succeq 0} \quad & D_{\text{ld}}(A, A_0) & (2) \\ \text{s.t.} \quad & \text{tr}\{A(x_i - x_j)(x_i - x_j)^\top\} \leq u, & \forall (i, j) \in S \\ & \text{tr}\{A(x_i - x_j)(x_i - x_j)^\top\} \geq l, & \forall (i, j) \in D \end{aligned}$$

where $D_{\text{ld}}(A, A_0) = \text{tr}(AA_0^{-1}) - \log \det(AA_0^{-1}) - n$, is the LogDet divergence.

To handle situations where exactly solving the problem in (2) is not possible, slack variables may be introduced to the ITML objective. To solve this optimization problem, an algorithm involving repeated Bregman projections is presented in (Davis et al., 2007), which we use for the experiments reported in this paper.

2.3 Inference-Driven Metric Learning (IDML): Semi-Supervised

Notations: We first define the necessary notations. Let X be the $d \times n$ matrix of n instances in a d -dimensional space. Out of the n instances, n_l instances are labeled, while the remaining n_u instances are unlabeled, with $n = n_l + n_u$. Let S be a $n \times n$ diagonal matrix with $S_{ii} = 1$ iff instance

x_i is labeled. m is the total number of labels. Y is the $n \times m$ matrix storing training label information, if any. \hat{Y} is the $n \times m$ matrix of estimated label information, i.e., output of any classifier, with \hat{Y}_{il} denoting score of label l at node i .

The ITML metric learning algorithm, which we reviewed in Section 2.2, is supervised in nature, and hence it does not exploit widely available unlabeled data. In this section, we review Inference Driven Metric Learning (IDML) (Algorithm 1) (Dhillon et al., 2010), a recently proposed metric learning framework which combines an existing *supervised* metric learning algorithm (such as ITML) along with *transductive* graph-based label inference to learn a new distance metric from labeled as well as unlabeled data combined. In self-training styled iterations, IDML alternates between metric learning and label inference; with output of label inference used during next round of metric learning, and so on.

IDML starts out with the assumption that existing supervised metric learning algorithms, such as ITML, can learn a better metric if the number of available labeled instances is increased. Since we are focusing on the semi-supervised learning (SSL) setting with n_l labeled and n_u unlabeled instances, the idea is to automatically label the unlabeled instances using a graph based SSL algorithm, and then include instances with low assigned label entropy (i.e., high confidence label assignments) in the next round of metric learning. The number of instances added in each iteration depends on the threshold β^1 . This process is continued until no new instances can be added to the set of labeled instances, which can happen when either all the instances are already exhausted, or when none of the remaining unlabeled instances can be assigned labels with high confidence.

The IDML framework is presented in Algorithm 1. In Line 3, any supervised metric learner, such as ITML, may be used as the METRICLEARNER. Using the distance metric learned in Line 3, a new k-NN graph is constructed in Line 4, whose edge weight matrix is stored in W . In Line 5, GRAPHLABELINF optimizes over the newly constructed graph, the GRF objective (Zhu et al., 2003) shown in (3).

$$\min_{\hat{Y}'} \text{tr}\{\hat{Y}'^\top L \hat{Y}'\}, \text{ s.t. } \hat{S} \hat{Y} = \hat{S} \hat{Y}' \quad (3)$$

where $L = D - W$ is the (unnormalized) Lapla-

¹During the experiments in Section 3, we set $\beta = 0.05$

Algorithm 1: Inference Driven Metric Learning (IDML)

Input: instances X , training labels Y , training instance indicator S , label entropy threshold β , neighborhood size k

Output: Mahalanobis distance parameter A

```
1:  $\hat{Y} \leftarrow Y, \hat{S} \leftarrow S$ 
2: repeat
3:    $A \leftarrow \text{METRICLEARNER}(X, \hat{S}, \hat{Y})$ 
4:    $W \leftarrow \text{CONSTRUCTKNNGRAPH}(X, A, k)$ 
5:    $\hat{Y}' \leftarrow \text{GRAPHLABELINF}(W, \hat{S}, \hat{Y})$ 
6:    $U \leftarrow \text{SELECTLOWENTINST}(\hat{Y}', \hat{S}, \beta)$ 
7:    $\hat{Y} \leftarrow \hat{Y} + U\hat{Y}'$ 
8:    $\hat{S} \leftarrow \hat{S} + U$ 
9: until convergence (i.e.,  $U_{ii} = 0, \forall i$ )
10: return  $A$ 
```

cian, and D is a diagonal matrix with $D_{ii} = \sum_j W_{ij}$. The constraint, $\hat{S}\hat{Y} = \hat{S}'\hat{Y}'$, in (3) makes sure that labels on training instances are not changed during inference. In Line 6, a currently unlabeled instance x_i (i.e., $\hat{S}_{ii} = 0$) is considered a new labeled training instance, i.e., $U_{ii} = 1$, for next round of metric learning if the instance has been assigned labels with high confidence in the current iteration, i.e., if its label distribution has low entropy (i.e., $\text{ENTROPY}(\hat{Y}'_{i,:}) \leq \beta$). Finally in Line 7, training instance label information is updated. This iterative process is continued till no new labeled instance can be added, i.e., when $U_{ii} = 0 \forall i$. IDML returns the learned matrix A which can be used to compute Mahalanobis distance using Equation 1.

3 Experiments

3.1 Setup

Dataset	Dimension	Balanced
Electronics	84816	Yes
Books	139535	Yes
Kitchen	73539	Yes
DVDs	155465	Yes
WebKB	44261	Yes

Table 1: Description of the datasets used in Section 3. All datasets are binary with 1500 total instances in each.

Description of the datasets used during experiments in Section 3 are presented in Table 1. The

first four datasets – Electronics, Books, Kitchen, and DVDs – are from the sentiment domain and previously used in (Blitzer et al., 2007). WebKB is a text classification dataset derived from (Subramanya and Bilmes, 2008). For details regarding features and data pre-processing, we refer the reader to the origin of these datasets cited above. One extra preprocessing that we did was that we only considered features which occurred more 20 times in the entire dataset to make the problem more computationally tractable and also since the infrequently occurring features usually contribute noise. We use classification error (lower is better) as the evaluation metric. We experiment with the following ways of estimating transformation matrix P :

Original²: We set $P = I$, where I is the $d \times d$ identity matrix. Hence, the data is not transformed in this case.

RP: The data is first projected into a lower dimensional space using the Random Projection (RP) method (Bingham and Mannila, 2001). Dimensionality of the target space was set at $d' = \frac{\log n}{\epsilon^2 \log \frac{1}{\epsilon}}$, as prescribed in (Bingham and Mannila, 2001). We use the projection matrix constructed by RP as P . ϵ was set to 0.25 for the experiments in Section 3, which has the effect of projecting the data into a much lower dimensional space (84 for the experiments in this section). This presents an interesting evaluation setting as we already run evaluations in much higher dimensional space (e.g., Original).

PCA: Data instances are first projected into a lower dimensional space using Principal Components Analysis (PCA) (Jolliffe, 2002). Following (Weinberger and Saul, 2009), dimensionality of the projected space was set at 250 for all experiments. In this case, we used the projection matrix generated by PCA as P .

ITML: A is learned by applying ITML (see Section 2.2) on the Original space (above), and then we decompose A as $A = P^\top P$ to obtain P .

²Note that “Original” in the results tables refers to original space with features occurring more than 20 times. We also ran experiments with original set of features (without any thresholding) and the results were worse or comparable to the ones reported in the tables.

Datasets	Original $\mu \pm \sigma$	RP $\mu \pm \sigma$	PCA $\mu \pm \sigma$	ITML $\mu \pm \sigma$	IDML-IT $\mu \pm \sigma$
Electronics	31.3 \pm 0.9	42.5 \pm 1.0	46.4 \pm 2.0	33.0 \pm 1.0	30.7\pm0.7
Books	37.5 \pm 1.1	45.0 \pm 1.1	34.8 \pm 1.4	35.0 \pm 1.1	32.0\pm0.9
Kitchen	33.7 \pm 1.0	43.0 \pm 1.1	34.0 \pm 1.6	30.9 \pm 0.7	29.0\pm1.0
DVDs	39.0 \pm 1.2	47.7 \pm 1.2	36.2 \pm 1.6	37.0 \pm 0.8	33.9\pm1.0
WebKB	31.4 \pm 0.9	33.0 \pm 1.0	27.9 \pm 1.3	28.9 \pm 1.0	25.5\pm1.0

Table 2: Comparison of SVM % classification errors (lower is better), with 50 labeled instances (Sec. 3.2). $n_l=50$. and $n_u = 1450$. All results are averaged over ten trials. All hyperparameters are tuned on a separate random split.

Datasets	Original $\mu \pm \sigma$	RP $\mu \pm \sigma$	PCA $\mu \pm \sigma$	ITML $\mu \pm \sigma$	IDML-IT $\mu \pm \sigma$
Electronics	27.0 \pm 0.9	40.0 \pm 1.0	41.2 \pm 1.0	27.5 \pm 0.8	25.3\pm0.8
Books	31.0 \pm 0.7	42.9 \pm 0.6	31.3 \pm 0.7	29.9 \pm 0.5	27.7\pm0.7
Kitchen	26.3 \pm 0.5	41.9 \pm 0.7	27.0 \pm 0.9	26.1 \pm 0.8	24.8\pm0.9
DVDs	34.7 \pm 0.4	46.8 \pm 0.6	32.9 \pm 0.8	34.0 \pm 0.8	31.8\pm0.9
WebKB	25.7 \pm 0.5	31.1 \pm 0.5	24.9 \pm 0.6	25.6 \pm 0.4	23.9\pm0.4

Table 3: Comparison of SVM % classification errors (lower is better), with 100 labeled instances (Sec. 3.2). $n_l=100$. and $n_u = 1400$. All results are averaged over ten trials. All hyperparameters are tuned on a separate random split.

IDML-IT: A is learned by applying IDML (Algorithm 1) (see Section 2.3) on the Original space (above); with ITML used as METRICLEARNER in IDML (Line 3 in Algorithm 1). In this case, we treat the set of test instances (without their gold labels) as the unlabeled data. In other words, we essentially work in the transductive setting (Vapnik, 2000). Once again, we decompose A as $A = P^\top P$ to obtain P .

We also experimented with the supervised large-margin metric learning algorithm (LMNN) presented in (Weinberger and Saul, 2009). We found ITML to be more effective in practice than LMNN, and hence we report results based on ITML only. Each input instance, x , is now projected into the transformed space as Px . We now train different classifiers on this transformed space. All results are averaged over ten random trials.

3.2 Supervised Classification

We train a SVM classifier, with an RBF kernel, on the transformed space generated by the projection matrix P . SVM hyperparameter, C and RBF kernel bandwidth, were tuned on a separate development split. Experimental results with 50 and 100

labeled instances are shown in Table 2, and Table 3, respectively. From these results, we observe that IDML-IT consistently achieves the best performance across all experimental settings. We also note that in Table 3, performance difference between ITML and IDML-IT in the Electronics and Kitchen domains are statistically significant.

3.3 Semi-Supervised Classification

In this section, we trained the GRF classifier (see Equation 3), a graph-based semi-supervised learning (SSL) algorithm (Zhu et al., 2003), using Gaussian kernel parameterized by $A = P^\top P$ to set edge weights. During graph construction, each node was connected to its k nearest neighbors, with k treated as a hyperparameter and tuned on a separate development set. Experimental results with 50 and 100 labeled instances are shown in Table 4, and Table 5, respectively. As before, we experimented with $n_l = 50$ and $n_l = 100$. Once again, we observe that IDML-IT is the most effective method, with the GRF classifier trained on the data representation learned by IDML-IT achieving best performance in all settings. Here also, we observe that IDML-IT achieves the best performance across all experimental settings.

Datasets	Original $\mu \pm \sigma$	RP $\mu \pm \sigma$	PCA $\mu \pm \sigma$	ITML $\mu \pm \sigma$	IDML-IT $\mu \pm \sigma$
Electronics	47.9 \pm 1.1	49.0 \pm 1.2	43.2 \pm 0.9	34.9 \pm 0.5	34.0\pm0.5
Books	50.0 \pm 1.0	49.4 \pm 1.0	47.9 \pm 0.7	42.1 \pm 0.7	40.6\pm0.7
Kitchen	49.8 \pm 1.1	49.6 \pm 0.9	48.6 \pm 0.8	31.1 \pm 0.5	30.0\pm0.5
DVDs	50.1 \pm 0.5	49.9 \pm 0.7	49.4 \pm 0.6	42.1 \pm 0.4	41.2\pm0.5
WebKB	33.1 \pm 0.4	33.1 \pm 0.3	33.1 \pm 0.3	30.0 \pm 0.4	28.7\pm0.5

Table 4: Comparison of transductive % classification errors (lower is better) over graphs constructed using different methods (see Section 3.3), with $n_l = 50$ and $n_u = 1450$. All results are averaged over ten trials. All hyperparameters are tuned on a separate random split.

Datasets	Original $\mu \pm \sigma$	RP $\mu \pm \sigma$	PCA $\mu \pm \sigma$	ITML $\mu \pm \sigma$	IDML-IT $\mu \pm \sigma$
Electronics	43.5 \pm 0.7	47.2 \pm 0.8	39.1 \pm 0.7	31.3 \pm 0.2	30.8\pm0.3
Books	48.3 \pm 0.5	48.9 \pm 0.3	43.3 \pm 0.4	35.2 \pm 0.5	33.3\pm0.6
Kitchen	45.3 \pm 0.6	48.2 \pm 0.5	41.0 \pm 0.7	30.7 \pm 0.6	29.9\pm0.3
DVDs	48.6 \pm 0.3	49.3 \pm 0.5	45.9 \pm 0.5	42.6 \pm 0.4	41.7\pm0.3
WebKB	33.4 \pm 0.4	33.4 \pm 0.4	33.4 \pm 0.3	30.4 \pm 0.5	28.6\pm0.7

Table 5: Comparison of transductive % classification errors (lower is better) over graphs constructed using different methods (see Section 3.3), with $n_l = 100$ and $n_u = 1400$. All results are averaged over ten trials. All hyperparameters are tuned on a separate random split.

4 Conclusion

In this paper, we compared the effectiveness of the transformed spaces learned by recently proposed *supervised*, and *semi-supervised* metric learning algorithms to those generated by previously proposed *unsupervised* dimensionality reduction methods (e.g., PCA). To the best of our knowledge, this is the first study of its kind involving NLP datasets. Through a variety of experiments on different real-world NLP datasets, we demonstrated that supervised as well as semi-supervised classifiers trained on the space learned by IDML-IT consistently result in the lowest classification errors. Encouraged by these early results, we plan to explore further the applicability of IDML-IT in other NLP tasks (e.g., entity classification, word sense disambiguation, polarity lexicon induction, etc.) where better representation of the data is a pre-requisite for effective learning.

Acknowledgments

Thanks to Kuzman Ganchev for providing detailed feedback on a draft of this paper. This work was supported in part by NSF IIS-0447972 and DARPA HRO1107-1-0029.

References

- E. Bingham and H. Mannila. 2001. Random projection in dimensionality reduction: applications to image and text data. In *ACM SIGKDD*.
- J. Blitzer, M. Dredze, and F. Pereira. 2007. Biographies, bollywood, boom-boxes and blenders: Domain adaptation for sentiment classification. In *ACL*.
- J.V. Davis, B. Kulis, P. Jain, S. Sra, and I.S. Dhillon. 2007. Information-theoretic metric learning. In *ICML*.
- P. S. Dhillon, P. P. Talukdar, and K. Crammer. 2010. Inference-driven metric learning for graph construction. Technical Report MS-CIS-10-18, CIS Department, University of Pennsylvania, May.
- IT Jolliffe. 2002. *Principal component analysis*. Springer verlag.
- A. Subramanya and J. Bilmes. 2008. Soft-Supervised Learning for Text Classification. In *EMNLP*.
- V.N. Vapnik. 2000. *The nature of statistical learning theory*. Springer Verlag.
- K.Q. Weinberger and L.K. Saul. 2009. Distance metric learning for large margin nearest neighbor classification. *The Journal of Machine Learning Research*.
- X. Zhu, Z. Ghahramani, and J. Lafferty. 2003. Semi-supervised learning using Gaussian fields and harmonic functions. In *ICML*.

Wrapping up a Summary: from Representation to Generation

Josef Steinberger and Marco Turchi and
Mijail Kabadjov and Ralf Steinberger

EC Joint Research Centre
21027, Ispra (VA), Italy

{Josef.Steinberger, Marco.Turchi,
Mijail.Kabadjov, Ralf.Steinberger}
@jrc.ec.europa.eu

Nello Cristianini

University of Bristol,
Bristol, BS8 1UB, UK

nello@support-vector.net

Abstract

The main focus of this work is to investigate robust ways for generating summaries from summary representations without recurring to simple sentence extraction and aiming at more human-like summaries. This is motivated by empirical evidence from TAC 2009 data showing that human summaries contain on average more and shorter sentences than the system summaries. We report encouraging preliminary results comparable to those attained by participating systems at TAC 2009.

1 Introduction

In this paper we adopt the general framework for summarization put forward by Spärck-Jones (1999) – which views summarization as a three-fold process: interpretation, transformation and generation – and attempt to provide a clean instantiation for each processing phase, with a particular emphasis on the last, summary-generation phase often omitted or over-simplified in the mainstream work on summarization.

The advantages of looking at the summarization problem in terms of distinct processing phases are numerous. It not only serves as a common ground for comparing different systems and understanding better the underlying logic and assumptions, but it also provides a neat framework for developing systems based on clean and extendable designs. For instance, Gong and Liu (2002) proposed a method based on Latent Semantic Analysis (LSA) and later J. Steinberger et al. (2007) showed that solely by enhancing the first source *interpretation* phase, one is already able to produce better summaries.

There has been limited work on the last summary generation phase due to the fact that it is unarguably a very challenging problem. The vast

amount of approaches assume simple sentence selection, a type of extractive summarization, where often the summary representation and the end summary are, indeed, conflated.

The main focus of this work is, thus, to investigate robust ways for generating summaries from summary representations without recurring to simple sentence extraction and aiming at more human-like summaries. This decision is also motivated by empirical evidence from TAC 2009 data (see table 1) showing that human summaries contain on average more and shorter sentences than the system summaries. The intuition behind this is that, by containing more sentences, a summary is able to capture more of the important content from the source.

Our initial experimental results show that our approach is feasible, since it produces summaries, which when evaluated against the TAC 2009 data¹ yield ROUGE scores (Lin and Hovy, 2003) comparable to the participating systems in the Summarization task at TAC 2009. Taking into account that our approach is completely unsupervised and language-independent, we find our preliminary results encouraging.

The remainder of the paper is organised as follows: in the next section we briefly survey the related work, in §3 we describe our approach to summarization, in §4 we explain how we tackle the generation step, in §5 we present and discuss our experimental results and towards the end we conclude and give pointers to future work.

2 Related Work

There is a large body of literature on summarization (Hovy, 2005; Erkan and Radev, 2004; Kupiec et al., 1995). The most closely related work to the approach presented hereby is work on summarization attempting to go beyond simple sentence ex-

¹<http://www.nist.gov/tac/>

traction and to a lesser degree work on sentence compression. We survey below work along these lines.

Although our approach is related to sentence compression (Knight and Marcu, 2002; Clarke and Lapata, 2008), it is subtly different. Firstly, we reduce the number of terms to be used in the summary at a global level, not at a local per-sentence level. Secondly, we directly exploit the resulting structures from the SVD making the last generation step fully aware of previous processing stages, as opposed to tackling the problem of sentence compression in isolation.

A similar approach to our sentence reconstruction method has been developed by Quirk et al. (2004) for paraphrase generation. In their work, training and test sets contain sentence pairs that are composed of two different proper English sentences and a paraphrase of a source sentence is generated by finding the optimal path through a paraphrases lattice.

Finally, it is worth mentioning that we are aware of the ‘capsule overview’ summaries proposed by Boguraev and Kennedy (1997) which is similar to our TSR (see below), however, as opposed to their emphasis on a suitable browsing interface rather than producing a readable summary, we precisely attempt the latter.

3 Three-fold Summarization: Interpretation, Transformation and Generation

We chose the LSA paradigm for summarization, since it provides a clear and direct instantiation of Spärck-Jones’ three-stage framework.

In LSA-based summarization the interpretation phase takes the form of building a term-by-sentence matrix $A = [A_1, A_2, \dots, A_n]$, where each column $A_j = [a_{1j}, a_{2j}, \dots, a_{nj}]^T$ represents the weighted term-frequency vector of sentence j in a given set of documents. We adopt the same weighting scheme as the one described in (Steinberger et al., 2007), as well as their more general definition of term entailing not only unigrams and bigrams, but also named entities.

The transformation phase is done by applying singular value decomposition (SVD) to the initial term-by-sentence matrix defined as $A = U\Sigma V^T$.

The generation phase is where our main contribution comes in. At this point we depart from standard LSA-based approaches and aim at produc-

ing a succinct summary representation comprised only of salient terms – Term Summary Representation (TSR). Then this TSR is passed on to another module which attempts to produce complete sentences. The module for sentence reconstruction is described in detail in section 4, in what follows we explain the method for producing a TSR.

3.1 Term Summary Representation

To explain how a term summary representation (TSR) is produced, we first need to define two concepts: *salience score* of a given term and *salience threshold*. Salience score for each term in matrix A is given by the magnitude of the corresponding vector in the matrix resulting from the dot product of the matrix of left singular vectors with the diagonal matrix of singular values. More formally, let $T = U \cdot \Sigma$ and then for each term i , the salience score is given by $|\vec{T}_i|$. Salience threshold is equal to the salience score of the top k^{th} term, when all terms are sorted in descending order on the basis of their salience scores and a cutoff is defined as a percentage (e.g., top 15%). In other words, if the total number of terms is n , then $100 * k/n$ must be equal to the percentage cutoff specified.

The generation of a TSR is performed in two steps. First, an initial pool of sentences is selected by using the same technique as in (Steinberger and Ježek, 2009) which exploits the dot product of the diagonal matrix of singular values with the right singular vectors: $\Sigma \cdot V^T$.² This initial pool of sentences is the output of standard LSA approaches.

Second, the terms from the source matrix A are identified in the initial pool of sentences and those terms whose *salience score* is above the *salience threshold* are copied across to the TSR. Thus, the TSR is formed by the most (globally) salient terms from each one of the sentences. For example:

- **Extracted Sentence:** “Irish Prime Minister Bertie Ahern admitted on Tuesday that he had held a series of private one-on-one meetings on the Northern Ireland peace process with Sinn Fein leader Gerry Adams, but denied they had been secret in any way.”
- **TSR Sentence at 10%:** “Irish Prime Minister Bertie Ahern Tuesday had held one-on-one meetings Northern Ireland peace process Sinn Fein leader Gerry Adams”³

²Due to space constraints, full details on that step are omitted here, see (Steinberger and Ježek, 2009).

³The TSR sentence is stemmed just before feeding it to the reconstruction module discussed in the next section.

Average number of:	Human Summaries	System Summaries	At 100%	At 15%	At 10%	At 5%	At 1%
Sentences/summary	6.17	3.82	3.8	3.95	4.39	5.18	12.58
Words/sentence	15.96	25.01	26.24	25.1	22.61	19.08	7.55
Words/summary	98.46	95.59	99.59	99.25	99.18	98.86	94.96

Table 1: Summary statistics on TAC’09 data (initial summaries).

Metric	LSA _{extract}	At 100%	At 15%	At 10%	At 5%	At 1%
ROUGE-1	0.371	0.361	0.362	0.365	0.372	0.298
ROUGE-2	0.096	0.08	0.081	0.083	0.083	0.083
ROUGE-SU4	0.131	0.125	0.126	0.128	0.131	0.104

Table 2: Summarization results on TAC’09 data (initial summaries).

4 Noisy-channel model for sentence reconstruction

This section describes a probabilistic approach to the reconstruction problem. We adopt the noisy-channel framework that has been widely used in a number of other NLP applications. Our interpretation of the noisy channel consists of looking at a stemmed string without stopwords and imagining that it was originally a long string and that someone removed or stemmed some text from it. In our framework, reconstruction consists of identifying the original long string.

To model our interpretation of the noisy channel, we make use of one of the most popular classes of SMT systems: the Phrase Based Model (PBM) (Zens et al., 2002; Och and Ney, 2001; Koehn et al., 2003). It is an extension of the noisy-channel model and was introduced by Brown et al. (1994), using phrases rather than words. In PBM, a source sentence f is segmented into a sequence of I phrases $f^I = [f_1, f_2, \dots, f_I]$ and the same is done for the target sentence e , where the notion of phrase is not related to any grammatical assumption; a phrase is an n-gram. The best translation e_{best} of f is obtained by:

$$e_{best} = \arg \max_e p(e|f) = \arg \max_e \prod_{i=1}^I \phi(f_i|e_i)^{\lambda_\phi}$$

$$d(a_i - b_{i-1})^{\lambda_d} \prod_{i=1}^{|e|} p_{LM}(e_i|e_1 \dots e_{i-1})^{\lambda_{LM}}$$

where $\phi(f_i|e_i)$ is the probability of translating a phrase e_i into a phrase f_i . $d(a_i - b_{i-1})$ is the distance-based reordering model that drives the system to penalize substantial reorderings of words during translation, while still allowing some flexibility. In the reordering model, a_i denotes the

start position of the source phrase that was translated into the i^{th} target phrase, and b_{i-1} denotes the end position of the source phrase translated into the $(i-1)^{\text{th}}$ target phrase. $p_{LM}(e_i|e_1 \dots e_{i-1})$ is the language model probability that is based on the Markov chain assumption. It assigns a higher probability to fluent/grammatical sentences. λ_ϕ , λ_{LM} and λ_d are used to give a different weight to each element (for more details see (Koehn et al., 2003)).

In our reconstruction problem, the difference between the source and target sentences is not in terms of languages, but in terms of forms. In fact, our source sentence f is a stemmed sentence without stopwords, while the target sentence e is a complete English sentence. “Translate” means to reconstruct the most probable sentence e given f inserting new words and reproducing the inflected surface forms of the source words.

4.1 Training of the model

In Statistical Machine Translation, a PBM system is trained using parallel sentences, where each sentence in a language is paired with another sentence in a different language and one is the translation of the other.

In the reconstruction problem, we use a set, S_1 of 2,487,414 English sentences extracted from the news. This set is duplicated, S_2 , and for each sentence in S_2 , stopwords are removed and the remaining words are stemmed using Porter’s stemmer (Porter, 1980). Our stopword list contains 488 words. Verbs are not included in this list, because they are relevant for the reconstruction task. To optimize the lambda parameters, we select 2,000 pairs as development set.

An example of training sentence pair is:

- Source Sentence: “royal mail ha doubl profit 321 million huge fall number letter post”
- Target Sentence: “royal mail has doubled its profits to 321 million despite a huge fall in the number of letters being posted”

In this work we use Moses (Koehn et al., 2007), a complete phrase-based translation toolkit for academic purposes. It provides all the state-of-the-art components needed to create a phrase-based machine translation system. It contains different modules to preprocess data, train the Language Models and the Translation Models.

5 Experimental Results

For our experiments we made use of the TAC 2009 data which conveniently contains human-produced summaries against which we could evaluate the output of our system (NIST, 2009).

To begin our inquiry we carried out a phase of exploratory data analysis, in which we measured the average number of sentences per summary, words per sentence and words per summary in human vs. system summaries in the TAC 2009 data. Additionally, we also measured these statistics of summaries produced by our system at five different percentage cutoffs: 100%, 15%, 10%, 5% and 1%.⁴ The results from this exploration are summarised in table 1. The most notable thing is that human summaries contain on average more and shorter sentences than the system summaries (see 2nd and 3rd column from left to right). Secondly, we note that as the percentage cutoff decreases (from 4th column rightwards) the characteristics of the summaries produced by our system are increasingly more similar to those of the human summaries. In other words, within the 100-word window imposed by the TAC guidelines, our system is able to fit more (and hence shorter) sentences as we decrease the percentage cutoff.

Summarization performance results are shown in table 2. We used the standard ROUGE evaluation (Lin and Hovy, 2003) which has been also used for TAC. We include the usual ROUGE metrics: R_1 is the maximum number of co-occurring unigrams, R_2 is the maximum number of co-occurring bigrams and R_{SU4} is the skip bigram measure with the addition of unigrams as counting

⁴Recall from section §3 that the salience threshold is a function of the percentage cutoff.

unit. The last five columns of table 2 (from left to right) correspond to summaries produced by our system at various percentage cutoffs. The 2nd column, $LSA_{extract}$, corresponds to the performance of our system at producing summaries by sentence extraction only.⁵

In the light of the above, the decrease in performance from column $LSA_{extract}$ to column ‘At 100%’ can be regarded as reconstruction error.⁶ Then, as we decrease the percentage cutoff (from 4th column rightwards) we are increasingly covering more of the content comprised by the human summaries (as far as the ROUGE metrics are able to gauge this, of course). In other words, the improvement of content coverage makes up for the reconstruction error, and at 5% cutoff we already obtain ROUGE scores comparable to $LSA_{extract}$. This suggests that if we improve the quality of our sentence reconstruction we would potentially end up with a better performing system than a typical LSA system based on sentence selection. Hence, we find these results very encouraging.

Finally, we admittedly note that by applying a percentage cutoff on the initial term set and further performing the sentence reconstruction we gain in content coverage, to a certain extent, on the expense of sentence readability.

6 Conclusion

In this paper we proposed a novel approach to summary generation from summary representation based on the LSA summarization framework and on a machine-translation-inspired technique for sentence reconstruction.

Our preliminary results show that our approach is feasible, since it produces summaries which resemble better human summaries in terms of the average number of sentences per summary and yield ROUGE scores comparable to the participating systems in the Summarization task at TAC 2009. Bearing in mind that our approach is completely unsupervised and language-independent, we find our results promising.

In future work we plan on working towards improving the quality of our sentence reconstruction step in order to produce better and more readable sentences.

⁵These are, effectively, what we called initial pool of sentences in section 3, before the TSR generation.

⁶The only difference between the two types of summaries is the reconstruction step, since we are including 100% of the terms.

References

- B. Boguraev and C. Kennedy. 1997. Saliency-based content characterisation of text documents. In I. Mani, editor, *Proceedings of the Workshop on Intelligent and Scalable Text Summarization at the Annual Joint Meeting of the ACL/EACL*, Madrid.
- P. Brown, S. Della Pietra, V. Della Pietra, and R. Mercer. 1994. The mathematics of statistical machine translation: Parameter estimation. *Computational Linguistics*, 19(2):263–311.
- J. Clarke and M. Lapata. 2008. Global inference for sentence compression: An integer linear programming approach. *Journal of Artificial Intelligence Research*, 31:273–318.
- G. Erkan and D. Radev. 2004. LexRank: Graph-based centrality as saliency in text summarization. *Journal of Artificial Intelligence Research (JAIR)*.
- Y. Gong and X. Liu. 2002. Generic text summarization using relevance measure and latent semantic analysis. In *Proceedings of ACM SIGIR*, New Orleans, US.
- E. Hovy. 2005. Automated text summarization. In Ruslan Mitkov, editor, *The Oxford Handbook of Computational Linguistics*, pages 583–598. Oxford University Press, Oxford, UK.
- K. Knight and D. Marcu. 2002. Summarization beyond sentence extraction: A probabilistic approach to sentence compression. *Artificial Intelligence*, 139(1):91–107.
- P. Koehn, F. Och, and D. Marcu. 2003. Statistical phrase-based translation. In *Proceedings of NAACL '03*, pages 48–54, Morristown, NJ, USA.
- P. Koehn, H. Hoang, A. Birch, C. Callison-Burch, M. Federico, N. Bertoldi, B. Cowan, W. Shen, C. Moran, R. Zens, C. Dyer, O. Bojar, A. Constantin, and E. Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In *Proceedings of ACL '07, demonstration session*.
- J. Kupiec, J. Pedersen, and F. Chen. 1995. A trainable document summarizer. In *Proceedings of the ACM SIGIR*, pages 68–73, Seattle, Washington.
- C. Lin and E. Hovy. 2003. Automatic evaluation of summaries using n-gram co-occurrence statistics. In *Proceedings of HLT-NAACL*, Edmonton, Canada.
- NIST, editor. 2009. *Proceeding of the Text Analysis Conference*, Gaithersburg, MD, November.
- F. Och and H. Ney. 2001. Discriminative training and maximum entropy models for statistical machine translation. In *Proceedings of ACL '02*, pages 295–302, Morristown, NJ, USA.
- M. Porter. 1980. An algorithm for suffix stripping. *Program*, 14(3):130–137.
- C. Quirk, C. Brockett, and W. Dolan. 2004. Monolingual machine translation for paraphrase generation. In *Proceedings of EMNLP*, volume 149. Barcelona, Spain.
- K. Spärck-Jones. 1999. Automatic summarising: Factors and directions. In I. Mani and M. Maybury, editors, *Advances in Automatic Text Summarization*. MIT Press.
- J. Steinberger and K. Ježek. 2009. Update summarization based on novel topic distribution. In *Proceedings of the 9th ACM DocEng, Munich, Germany*.
- J. Steinberger, M. Poesio, M. Kabadjov, and K. Ježek. 2007. Two uses of anaphora resolution in summarization. *Information Processing and Management*, 43(6):1663–1680. Special Issue on Text Summarization (Donna Harman, ed.).
- R. Zens, F. J. Och, and H. Ney. 2002. Phrase-based statistical machine translation. In *Proceedings of KI '02*, pages 18–32, London, UK. Springer-Verlag.

Author Index

- Aharonson, Vered, 331
Alumäe, Tanel, 301
Ambati, Vamshi, 365
Asahara, Masayuki, 98
- Baker, Collin, 68
Becker, Israela, 331
Beigman Klebanov, Beata, 253
Beigman, Eyal, 253
Ben Aharon, Roni, 241
Benajiba, Yassine, 281
Bienenstock, Elie, 215
Blackwood, Graeme, 27
Blunsom, Phil, 225
Bojar, Ondřej, 86
Bouma, Gerlof, 109
Bui, Trung H., 307
Buttler, David, 156
Byrne, William, 27
- Campbell, Gwendolyn, 43
Cao, Hailong, 17
Carbonell, Jaime, 365
Cardie, Claire, 156, 269, 336
Carpuat, Marine, 178
Casacuberta, Francisco, 173
Chang, Jason S., 115
Chang, Yu-Chia, 115
Charniak, Eugene, 33
Chiang, David, 209
Chodorow, Martin, 353
Choi, Yejin, 269, 336
Cohn, Trevor, 225
Crammer, Koby, 377
Cristianini, Nello, 382
Cui, Lei, 6
- Dagan, Ido, 241
Danescu-Niculescu-Mizil, Cristian, 247
de Gispert, Adrià, 27
Deng, Yonggang, 22
Dhillon, Paramveer S., 377
Diab, Mona, 281
Diermeier, Daniel, 253
Dixon, Paul, 275
- Dobrinkat, Marcus, 80
Dzikovska, Myroslava O., 43
- Elsner, Micha, 33
Erk, Katrin, 92
Etzioni, Oren, 286
- Fader, Anthony, 286
Faruquie, Tanveer A, 126
Fellbaum, Christiane, 68
Finlayson, Mark, 49
Foster, Jennifer, 353
Fujino, Akinori, 137
Fujita, Sanae, 162
Furui, Sadaoki, 236, 275
- Ganchev, Kuzman, 194
Gilbert, Nathan, 156
Gillenwater, Jennifer, 194
González Rubio, Jesús, 173
Goudbeek, Martijn, 55
Graça, João, 194
Gurevych, Iryna, 263
- Habash, Nizar, 178
Haghighi, Aria, 291
Hasegawa, Takaaki, 325
Heie, Matthias H., 236
Hervas, Raquel, 49
Hou, Yuexian, 120
Hwang, Young-Sook, 142
Hysom, David, 156
- Ide, Nancy, 68
Iwata, Tomoharu, 184
- Jakob, Niklas, 263
Jeong, Minwoo, 151
Johnson, Mark, 215
Ju, Yun-Cheng, 313
- K, Hima Prasad, 126
Kabadjov, Mijail, 382
Kato, Yoshihide, 74
Keller, Frank, 60
Kettunen, Kimmo, 80

Kikui, Genichiro, 325
Kipersztok, Oscar, 371
Kitagawa, Kotaro, 189
Klein, Dan, 200, 291, 348
Koller, Daphne, 371
Kos, Kamil, 86
Kovashka, Adriana, 38
Krahmer, Emiel, 55
Kulick, Seth, 342
Kurimo, Mikko, 301

Lall, Ashwin, 231
Lamar, Michael, 215
Lee, Lillian, 247
Lewis, William, 220
Li, Decong, 296
Li, Mu, 6
Li, Sujian, 296
Li, Wenjie, 296
Li, Xiao-Li, 359
Liu, Bing, 359
Liu, Qun, 12, 142
Liu, Yang, 12
Lv, Yajuan, 12, 142

Mareček, David, 86
Maron, Yariv, 215
Marton, Yuval, 178
Matsubara, Shigeki, 74
Matsumoto, Yuji, 98
Matsuo, Yoshihiro, 325
Maxwell, Tamsin, 120
Mi, Haitao, 12
Mitamura, Teruko, 115
Mochihashi, Daichi, 184
Mohania, Mukesh, 126
Mooney, Raymond, 38
Moore, Johanna D., 43
Moore, Robert C., 220
Murthy, Karin, 126

Nagata, Masaaki, 137, 162
Navaretta, Costanza, 318
Ng, See-Kiong, 359
Nishikawa, Hitoshi, 325

Onishi, Takashi, 1
Ortiz Martínez, Daniel, 173

Pado, Sebastian, 92
Paek, Tim, 313
Paggio, Patrizia, 318
Pal, Christopher, 258

Passonneau, Rebecca, 68
Pauls, Adam, 200, 209, 348
Pereira, Fernando, 194
Peters, Stanley, 307
Power, Richard, 132

Qu, Weiguang, 296
Quirk, Chris, 200

Raghavan, Sindhu, 38
Riloff, Ellen, 156
Rosso, Paolo, 281

Sawada, Hiroshi, 184
Søgaard, Anders, 205
Sharif Razavian, Narges, 147
Shindo, Hiroyuki, 137
Soderland, Stephen, 286
Song, Dawei, 120
Steinberger, Josef, 382
Steinberger, Ralf, 382
Steinhauser, Natalie, 43
Stoyanov, Veselin, 156
Su, Jinsong, 12
Subramaniam, L Venkata, 126
Sumita, Eiichiro, 1, 17
Sun, Weiwei, 103, 168
Szpektor, Idan, 241

Taira, Hirotoshi, 162
Talukdar, Partha Pratim, 377
Tanaka-Ishii, Kumiko, 189
Tapiovaara, Tero, 80
Taskar, Ben, 194
Tetreault, Joel, 353
Titov, Ivan, 151
Turchi, Marco, 382

Utiyama, Masao, 1

Van Durme, Benjamin, 231
Vaswani, Ashish, 209
Väyrynen, Jaakko, 80
Vickrey, David, 371
Vogel, Stephan, 147, 365

Wang, Wei, 296
Wang, Zhiyang, 142
Watanabe, Yotaro, 98
Wei, Bin, 258
Whittaker, Edward W. D., 236
Wu, Jian-Cheng, 115

Xiang, Bing, 22

Yan, Tingxu, 120

Yang, Dong, 275

Yessenalina, Ainur, 336

Zhang, Dongdong, 6

Zhang, Lei, 359

Zhang, Peng, 120

Zhao, Tiejun, 6

Zhou, Bowen, 22

Zhou, Ming, 6

Zitouni, Imed, 281

ACL 2010

**48th Annual Meeting of the
Association for Computational Linguistics**

Proceedings of the Student Research Workshop

13 July 2010
Uppsala University
Uppsala, Sweden

Production and Manufacturing by
Taberg Media Group AB
Box 94, 562 02 Taberg
Sweden

©2010 The Association for Computational Linguistics

Order copies of this and other ACL proceedings from:

Association for Computational Linguistics (ACL)
209 N. Eighth Street
Stroudsburg, PA 18360
USA
Tel: +1-570-476-8006
Fax: +1-570-476-0860
acl@aclweb.org

Introduction

We are pleased to present the papers accepted for presentation at the Student Research Workshop of the 48th Annual Meeting of the Association for Computational Linguistics held in Uppsala, Sweden, July 11-16, 2010. The Student Research Workshop, an established tradition at the annual meetings of ACL, offers students the opportunity to present their work in a setting embedded in the main conference and features ACL's efforts to invest in young researchers who will be a part of the research community in Computational Linguistics, Natural Language Processing and related fields in the near future. The workshop aimed at enabling students to exchange ideas with other researchers and experts and to receive useful feedback and directions for future research in an early stage of their work.

We are indebted to many people who contributed to the organization of the ACL 2010 Student Research Workshop. A total of 54 students and senior researchers agreed to serve on the program committee, of which 21 were from North America, 24 were from Europe, 2 were from Middle East, and 7 were from Asia/Pacific regions. The committee members, who represented various areas of NLP and computational linguistics research, well understood the spirit of the Student Research Workshop and provided high-quality, constructive, and elaborate reviews to all students who submitted their work. We are thankful to the members of the program committee for the time they spent reading and reviewing the papers.

We received 37 submissions from all over the world. Each submission was assigned to 3 reviewers, at least one of which was a senior researcher. We accepted 7 submissions for oral and poster presentation and 12 submissions for only poster presentation during a parallel session with the main conference. The presentation format was assigned based on the suggestions of the reviewers about how the work could be presented best, and does not indicate a quality difference among accepted papers. We would like to thank all students who contributed to the success of this years Student Research Workshop by submitting their research papers from a wide range of topics.

We are very grateful to Tomek Strzalkowski and Marketa Lopatkova, our tireless faculty advisors, for their advice, constant support, and obtaining funding. We also would like to thank the conference organizers of ACL'2010: the general chair, Jan Hajic, the program chairs, Sandra Carberry and Stephen Clark, the publications chairs Jing-Shin Chang and Philipp Koehn, and the local organization committee in particular Joakim Nivre, and Priscilla Rasmussen.

Finally, we wish to thank the U.S. National Science Foundation (NSF), The Don and Betty Walker Student Scholarship Fund of the Association for Computational Linguistics, and European Chapter of the ACL (EACL) for generously sponsoring this workshop by offering grants to assist presenters in covering their registration, accommodation, and travel costs.

The ACL 2010 Student Research Workshop Co-Chairs
Seniz Demir, Nils Reiter, Jan Raab

Faculty Advisors:

Marketa Lopatkova
Charles University
Prague, Czech Republic

Tomek Strzalkowski
State University of New York (SUNY)
Albany, NY, USA

Chairs:

Seniz Demir
University of Delaware
Newark, DE, USA

Jan Raab
Charles University
Prague, Czech Republic

Nils Reiter
Heidelberg University
Germany

Program Committee:

Eneko Agirre, University of the Basque Country, Spain
Regina Barzilay, MIT, USA
Fadi Biadsy, Columbia University, USA
Johan Bos, Italy
Razvan Bunescu, Ohio University, USA
Richard Burns, University of Delaware, USA
Aoife Cahill, University of Stuttgart, Germany
Bethold Crysmann, Saarland University, Germany
Micha Elsner, Brown University, USA
Alex Gruenstein, Google, USA
Tunga Gungor, Bogazici University, Turkey
Aria Haghighi, UC Berkeley, USA
Dilek Hakkani-Tur, International Computer Science Institute (ICSI) - Berkeley, USA
Ilana Heintz, Ohio State University, USA
Julia Hockenmaier, University of Illinois at Urbana-Champaign, USA
Alexander Koller, Saarland University, Germany
Udo Kruschwitz, University of Essex, UK
Cheongjae Lee, Pohang University of Science and Engineering, Korea
Wei-Hao Lin, Google, USA
Yang Liu, University of Texas at Dallas, USA

Daniel Marcu, Information Sciences Institute (ISI) - USC, USA
Philippe Martin, University of Toronto, Canada
Sameer Maskey, IBM, USA
Diana Maynard, University of Sheffield, UK
David McDonald, BBN Technologies, USA
Timothy Miller, University of Minnesota Twin Cities, USA
Jeremy Nicholson, University of Melbourne, Australia
Kemal Oflazer, Carnegie Mellon University in Qatar, Qatar
Rainer Osswald, Distance University of Hagen, Germany
You Ouyang, Hong Kong Polytechnic University, Hong Kong
Sebastian Pado, University of Stuttgart, Germany
Katerina Pastra, Institute for Language and Speech Processing, Greece
Paul Piwek, The Open University, UK
Jan Ptacek, Charles University in Prague, Czech Republic
Verena Rieser, University of Edinburgh, UK
Jan Romportl, University of West Bohemia, Czech Republic
Helmut Schmid, University of Stuttgart, Germany
Sabine Schulte im Walde, University of Stuttgart, Germany
Samira Shaikh, University at Albany - SUNY, NY, USA
Drahomira Spoustova, Charles University in Prague, Czech Republic
Xu Sun, University of Tokyo, Japan
Idan Szpektor, Bar Ilan University, Israel
Ahmet Cuneyd Tantug, Istanbul Technical University, Turkey
Keith Trnka, University of Delaware, USA
Ming-Feng Tsai, National University of Singapore, Singapore
Bonnie Webber, University of Edinburgh, UK
Furu Wei, IBM, China
Michael White, Ohio State University, USA
Yorick Wilks, University of Sheffield, UK
Emine Yilmaz, Microsoft, UK
Zdenek Zabokrtsky, Charles University in Prague, Czech Republic
Daniel Zeman, Charles University in Prague, Czech Republic
Guodong Zhou, Soochow University, China
Xiaodan Zhu, University of Toronto, Canada

Table of Contents

<i>Non-Cooperation in Dialogue</i>	
Brian Plüss	1
<i>Towards Relational POMDPs for Adaptive Dialogue Management</i>	
Pierre Lison	7
<i>WSD as a Distributed Constraint Optimization Problem</i>	
Siva Reddy and Abhilash Inumella	13
<i>A Probabilistic Generative Model for an Intermediate Constituency-Dependency Representation</i>	
Federico Sangati	19
<i>Sentiment Translation through Lexicon Induction</i>	
Christian Scheible	25
<i>Unsupervised Search for the Optimal Segmentation for Statistical Machine Translation</i>	
Coşkun Mermer and Ahmet Afşın Akın	31
<i>How Spoken Language Corpora Can Refine Current Speech Motor Training Methodologies</i>	
Daniil Umanski and Federico Sangati	37
<i>Mood Patterns and Affective Lexicon Access in Weblogs</i>	
Thin Nguyen	43
<i>Growing Related Words from Seed via User Behaviors: A Re-Ranking Based Approach</i>	
Yabin Zheng, Zhiyuan Liu and Lixing Xie	49
<i>Transition-Based Parsing with Confidence-Weighted Classification</i>	
Martin Haulrich	55
<i>Expanding Verb Coverage in Cyc with VerbNet</i>	
Clifton McFate	61
<i>A Framework for Figurative Language Detection Based on Sense Differentiation</i>	
Daria Bogdanova	67
<i>Automatic Selectional Preference Acquisition for Latin Verbs</i>	
Barbara McGillivray	73
<i>Edit Tree Distance Alignments for Semantic Role Labelling</i>	
Hector-Hugo Franco-Penya	79
<i>Automatic Sanskrit Segmentizer Using Finite State Transducers</i>	
Vipul Mittal	85
<i>Adapting Self-Training for Semantic Role Labeling</i>	
Rasoul Samad Zadeh Kaljahi	91
<i>Weakly Supervised Learning of Presupposition Relations between Verbs</i>	
Galina Tremper	97
<i>Importance of Linguistic Constraints in Statistical Dependency Parsing</i>	
Bharat Ram Ambati	103

The Use of Formal Language Models in the Typology of the Morphology of Amerindian Languages

Andres Osvaldo Porta..... 109

Conference Program

Tuesday, July 13, 2010

Oral Session (10:30–11:45)

Chaired by Jan Raab

- 10:30–10:40 *Non-Cooperation in Dialogue*
Brian Plüss
- 10:40–10:50 *Towards Relational POMDPs for Adaptive Dialogue Management*
Pierre Lison
- 10:50–11:00 *WSD as a Distributed Constraint Optimization Problem*
Siva Reddy and Abhilash Inumella
- 11:00–11:10 *A Probabilistic Generative Model for an Intermediate Constituency-Dependency Representation*
Federico Sangati
- 11:10–11:20 *Sentiment Translation through Lexicon Induction*
Christian Scheible
- 11:20–11:30 *Unsupervised Search for the Optimal Segmentation for Statistical Machine Translation*
Coşkun Mermer and Ahmet Afşın Akın
- 11:30–11:40 *How Spoken Language Corpora Can Refine Current Speech Motor Training Methodologies*
Daniil Umanski and Federico Sangati
- 11:55–13:15 Each talk is accompanied by a poster presentation

Tuesday, July 13, 2010 (continued)

Poster Session (11:55–13:15)

Mood Patterns and Affective Lexicon Access in Weblogs

Thin Nguyen

Growing Related Words from Seed via User Behaviors: A Re-Ranking Based Approach

Yabin Zheng, Zhiyuan Liu and Lixing Xie

Transition-Based Parsing with Confidence-Weighted Classification

Martin Haulrich

Expanding Verb Coverage in Cyc with VerbNet

Clifton McFate

A Framework for Figurative Language Detection Based on Sense Differentiation

Daria Bogdanova

Automatic Selectional Preference Acquisition for Latin Verbs

Barbara McGillivray

Edit Tree Distance Alignments for Semantic Role Labelling

Hector-Hugo Franco-Penya

Automatic Sanskrit Segmentizer Using Finite State Transducers

Vipul Mittal

Adapting Self-Training for Semantic Role Labeling

Rasoul Samad Zadeh Kaljahi

Weakly Supervised Learning of Presupposition Relations between Verbs

Galina Tremper

Importance of Linguistic Constraints in Statistical Dependency Parsing

Bharat Ram Ambati

The Use of Formal Language Models in the Typology of the Morphology of Amerindian Languages

Andres Osvaldo Porta

Non-Cooperation in Dialogue

Brian Plüss

Centre for Research in Computing

The Open University

Milton Keynes, UK

b.pluss@open.ac.uk

Abstract

This paper presents ongoing research on computational models for non-cooperative dialogue. We start by analysing different levels of cooperation in conversation. Then, inspired by findings from an empirical study, we propose a technique for measuring non-cooperation in political interviews. Finally, we describe a research programme towards obtaining a suitable model and discuss previous accounts for conflictive dialogue, identifying the differences with our work.

1 Introduction

Most approaches to modeling conversation are based on a strong notion of cooperation between the dialogue participants (DPs). Traditional models using intentions (Cohen and Levesque, 1991), dialogue games (Power, 1979), shared plans (Grosz and Sidner, 1990) or collaborative problem-solving (Blaylock and Allen, 2005) explain dialogue situations in which DPs recognise each other's intentions and, at least to a certain extent, accept each other's goals when deciding on their actions. These assumptions are theoretically grounded, as most work in linguistics has considered situations in which DPs share a common goal and cooperate to achieve it by means of conversation (Grice, 1975; Clark and Schaefer, 1989). They are also practically sound: dialogue models are usually implemented in the form of dialogue systems, built for the purpose of providing a service to their users (e.g., TRAINS (Allen and Schubert, 1991)). In this scenario, failure to cooperate, either on the side of the system or of the user, is against the premises on which the system is conceived and used.

In everyday conversation, however, a great many situations escape the arguments above. Con-

sider the following example¹:

- (1) PAXMAN [1]: (interrupting) Did you threaten to overrule him?
HOWARD [2]: I, I, was not entitled to *instruct* Derek Lewis, and I did not instruct him.
PAXMAN [3]: Did you threaten to overrule him?
HOWARD [4]: The truth of the matter is that Mr. Marriott was not suspended. I...
PAXMAN [5]: (overlapping) Did you threaten to overrule him?
HOWARD [6]: ... did not overrule Derek Lewis.
PAXMAN [7]: Did you *threaten* to overrule him?
HOWARD [8]: I took advice on what I could or could not do...
PAXMAN [9]: (overlapping) Did you threaten to overrule him, Mr. Howard?
HOWARD[10]: ... and I acted scrupulously in accordance with that advice, I did *not* overrule Derek Lewis...
PAXMAN[11]: (overlapping) Did you threaten to overrule him?
HOWARD[12]: ... Mr. Marriott was *not* suspended.
PAXMAN[13]: Did you threaten to overrule him?
HOWARD[14]: (pauses) I have accounted for my decision to dismiss Derek Lewis...
PAXMAN[15]: (overlapping) Did you threaten to overrule him?
HOWARD[16]: ... in great detail, before the House of Commons.
PAXMAN[17]: I note that you're not answering the question of whether you *threatened* to overrule him.

(*Newsnight*, BBC, 1997)

We take it for granted that, at some level, Paxman and Howard are sharing a goal, for otherwise they would not be having an interview. Still, the exchange is clearly conflictive, to the point that their behaviour compromises the flow of the conversation.

Heritage (1998) analyses the distinctive roles of DPs in news interviews:

¹BBC presenter Jeremy Paxman questions former UK Home Secretary Michael Howard with respect to a meeting in 1995 between Howard and the head of the Prison Service, Derek Lewis, about the dismissal of the governor of Parkhurst Prison, John Marriott, due to repeated security failures. The case was given considerable attention in the media, as a result of accusations by Lewis that Howard had instructed him, thus exceeding the powers of his office.

“the participants -IRs [=interviewers] and IEs [=interviewees]- exclude themselves from a wide variety of actions that they are normally free to do in the give and take of ordinary conversation. If IRs restrict themselves to asking questions, then they cannot - at least overtly - express opinions, or argue with, debate or criticize the interviewees’ positions nor, conversely, agree with, support or defend them. Correspondingly, if IEs restrict themselves to answers (or responses) to questions, then they cannot ask questions (of IRs or other IEs), nor make unsolicited comments on previous remarks, initiate changes of topic, or divert the discussion into criticisms of the IR or the broadcasting organization.”

(Heritage, 1998, p.8)

Now, consider the fragment below²:

- (2) PAXMAN[1]: Can you clear up whether or not you did threaten to overrule Derek Lewis when you were Home Secretary?
 HOWARD[2]: Oh, come on, Jeremy, you are really going to go back over that again? As...
 PAXMAN[3]: (overlapping) You’ve had seven years to think about it!
 HOWARD[4]: (overlapping)...as, as it happens, I didn’t. Are you satisfied now?
 PAXMAN[5]: Thank you. Why didn’t you say that at the time?
 HOWARD[6]: I, well, we’ve been over this many, many times. I, I, I knew that everyone was crawling over every syllable I said about that, and I wanted to check very carefully what I said before answering your question.

(*Newsnight*, BBC, 2004)

On this occasion, Howard provides an answer almost immediately and the flow of the conversation contrasts noticeably with that in (1). The investigation reported in this article aims at shedding light on the nature of non-cooperation in dialogue, by capturing the intuitions that allow us to differentiate between both conversations in terms of participant behaviour.

Dialogue games supporters could say that there is a game that describes the interaction in the first example. While this might be true, such an approach would force us, in the limit, to define one game for each possible conversation that would not fit a certain standard. Walton and Krabbe (1995) attempt a game-based approach in their study of natural argumentation. They claim that a rigorous model of conversational interaction is useful, but accept that most of the huge variety of everyday conversation escapes it. Dialogue games are based on strict rules that capture typical dialogue situations while leaving out considerable detail. As example (1) shows, DPs behaviour can

²This exchange took place seven years after (1), when public awareness of the 1995 affair had dissipated.

divert from the typical case in unexpected ways, falling outside the characterisation³.

Nevertheless, the rules and patterns captured by game models are useful, as they describe the expected behaviour of the DPs under a certain conversational scenario. In our research, we aim at reconciling two worlds, using the insights from dialogue games to provide a description of expected behaviour in the form of social obligations, but looking at naturally occurring cases that deviate from the norm. This, in turn, calls for a technique to measure non-cooperation in dialogue and in this paper we provide one that is theoretically sound and supported by empirical evidence.

The following section discusses levels of cooperation in dialogue; Section 3 presents an empirical study and a practical measure of non-cooperation in political interviews; in Section 4 we discuss related work, our working hypothesis and a methodology; and Section 5 has the conclusions.

2 Linguistic and Non-Linguistic Cooperation

Cooperation in dialogue can happen at different levels. In most cases, conversation supports a social activity that constrains the behaviour acceptable or expected from the participants. In addition, conversational behaviour determines how cooperatively participants engage in a social activity. However, cooperation at the conversational level does not necessarily translate to the social level. Consider, for instance, a witness under interrogation in a U.S. trial refusing to answer a question by appealing to the Fifth Amendment of the Constitution⁴. Such behaviour will be accepted in the conversational setting as established by law, although it is not cooperative in relation with the goals of the trial. Non-cooperation at the conversational level, on the other hand, usually results in lack of cooperation at the social level. Take as an example, the same witness remaining silent, rather than answering or appealing to the Fifth Amendment.

To illustrate further, consider a fictional alternative to (1), where Howard replies by saying “I will not answer that question, as it is not relevant to whether I exceeded the powers of my office”.

³Consider, for instance, Ginzburg’s QUD model (Ginzburg, 1996) when applied to dialogue (1), in which Howard repeatedly fails to either accept or reject Paxman’s question.

⁴“No person shall (. . .) be compelled in any criminal case to be a witness against himself”.

This is not cooperative for the interview, but it is so at the linguistic level. It would help in preserving the flow of the conversation, e.g., by triggering a sub-dialogue to solve the disagreement.

The distinction between linguistic and non-linguistic (also called task-related, high-level or social) cooperation has been addressed before. Attardo (1997) revisits Gricean pragmatics, relating non-linguistic cooperation to participants' behaviour towards realising task-related goals, and linguistic cooperation to assumptions on their respective behaviour in order to encode and decode intended meaning. From a computational perspective, Bunt (1994) relies on a similar distinction for defining dialogue acts. Also, Traum and Allen (1994) introduce discourse obligations as an alternative to joint intentions and shared plans, to allow for models of dialogues in which participants do not share the same high-level goals and where behaviour is also determined by "a sense of obligation to behave within limits set by the society" (Traum and Allen, 1994, p.2).

Walton and Krabbe (1995) proposed a typology of dialogue based on the initial situation triggering the exchange and participants' shared aims and individual goals. Based on their work, Reed and Long (1997) distinguish cases where participants follow a common set of dialogue rules and stay within a mutually acknowledged framework from a stronger notion in which their individual goals are in the same direction. Borrowing from the latter, in the rest of the paper, we will speak of *collaboration* when DPs share the same task-level goals, and use *cooperation* when participants follow the conversational obligations imposed by the social activity (i.e., linguistic cooperation as discussed above). We will not deal with collaboration here, though, as our focus is on non-cooperation.

3 An Empirical Study

In this section, we describe an empirical pilot study aimed at identifying a set of features that distinguish cooperative from non-cooperative conversational behaviour and at establishing a suitable domain in which to focus our work.

3.1 The Corpus

We collected the transcripts of 10 adversarial dialogues: 4 political interviews, 2 entertainment interviews, 1 parliamentary inquiry, 1 courtroom confrontation, 1 courtroom interrogation and 1

dispute. The corpus includes 2 collaborative political interviews for result comparison and is nearly 14,500 words long⁵.

In a first analysis, we identified those surface features that characterised each conversation as conflictive: e.g., interruptions, short turns, unfinished adjacency pairs, verbatim repetition. Next, looking for a better understanding, we performed an in-depth case study of one of the examples, approaching the analysis from different angles.

By studying, e.g., the observance of turn-taking rules, the implicatures of the participants and, more extensively, how the case fitted within the normative framework proposed by Walton and Krabbe (1995), we were able to better identify the nature of non-cooperative features present in the dialogue and establish a formalisable framework for approaching non-cooperative dialogue.

As for the domain, the wealth of interesting conversational situations that arise in political interviews make a suitable context for this research. In the English-speaking world, journalists are well-known for their incisive approach to public servants. At the same time, politicians are usually well trained to deliver a set of key messages when speaking in public, and to avoid issues unfavorable to their image. We will only consider naturally occurring (i.e. non-scripted) two-party interviews.

3.2 Degrees of Non-Cooperation

Based on the analysis described above, we propose a technique for measuring non-cooperation in political interviews using a set of non-cooperative features (NCFs). The number of occurrences of these features will determine the degree of non-cooperation (DNC) of an exchange.

We grouped NCFs following three aspects of conversation: turn-taking, grounding and speech acts (see Table 1 for a complete list).

Turn-taking rules (Sacks et al., 1974) establish that speakers make their contributions at adequate places and in particular ways. Interlocutors in a political interview are expected to respect transition-relevance places, openings and closings according to social conventions. Failing to do so (e.g., by interrupting each other) constitutes a non-cooperative feature.

Grounding (Clark and Schaefer, 1989) refers to participants' acknowledgement of each other's

⁵These resources are available at <http://www.open.ac.uk/blogs/brianpluss/pilot-study/>.

Turn-Taking	For both speakers: <ul style="list-style-type: none"> • interrupting • overlapping • ending the exchange abruptly
Grounding	Interviewer fails to either: <ul style="list-style-type: none"> • ask next relevant question • move to next topical issue • state irrelevance of answer Interviewee fails to either: <ul style="list-style-type: none"> • give relevant answer • reject question
Speech Acts	Interviewer either: <ul style="list-style-type: none"> • expresses personal opinion • argues, debates with or criticises interviewee’s position subjectively • agrees with, supports or defends interviewee’s position subjectively Interviewee either: <ul style="list-style-type: none"> • asks (non-CR) question • makes irrelevant comment • initiates change of topic • criticises interviewer

Table 1: NCFs for political interviews

contributions by providing evidence of understanding (e.g, continued attention, relevant next contribution). In political interviews a question is acknowledged by rejecting it or by providing a direct answer. Likewise, answers are acknowledged by rejecting their relevance, by asking a next relevant question or by moving on to a new topical issue. Failing to provide sufficient evidence of understanding is also a non-cooperative feature.

Speech Act theory (Searle, 1979) classifies utterances according to their associated force and propositional content. Going back to Heritage’s comment, in a political interview participants can fail to restrict their speech acts to the force and content expected for their role. Non-cooperative features related to speech acts include the interviewer expressing a personal opinion or criticising subjectively the interviewee’s positions and the interviewee asking questions (except for clarification requests) or making irrelevant comments.

We define the degree of non-cooperation (DNC) of a dialogue as the proportion of utterances with one of more occurrences of these non-cooperative features⁶. Furthermore, the DNC could be thus computed for the whole conversation and also for each participant, by counting only occurrences of features and utterances from each DP.

As an example, consider an extended fragment

⁶At this stage, all NCFs are weighted equally. This is a simplifying assumption we will remove in the future so that, e.g., an interviewee attempting a change of topic has a stronger impact on the DNC than, say, one interrupting.

of (1) annotated with non-cooperative features (**O**: overlap; **GF**: grounding failure; **UC**: unsolicited comment; **I**: interruption; **TC**: topic change):

- (3) P [11] : Uir.1 (overlapping) Did you threaten to **O**
overrule him?
H[12] : Uie.1 ... Mr. Marriot was *not* suspended. **GF**
P [13] : Uir.2 Did you threaten to overrule him? **GF**
H[14] : Uie.2 (pauses) I have accounted for my decision to dismiss Derek Lewis. . .
P [15] : Uir.3 (overlapping) Did you threaten to **O**
overrule him?
H[16] : Uie.2 ... in great detail before the House of **UC**
Commons.
P [17] : Uir.4 I note that you’re not answering the question whether you *threatened* to overrule him.
H[18] : Uie.3 Well, the important aspect of this **GF** which it’s very clear to bear in mind. . .
P [19] : Uir.5 (interrupting) I’m sorry, I’m going to **I**
be frightfully rude but. . .
H[20] : Uie.4 Yes, you can. . .
P [21] : Uir.6 (overlapping) I’m sorry. . . **O**
H[22] : Uie.4 (overlapping) ...you can put the **O**
question and I will give you, I will give you an answer.
P [23] : Uir.7 ...it’s a straight yes-or-no question and a straight yes-or-no answer:
Uir.8 did you threaten to overrule him?
H[24] : Uie.5 I discussed the matter with Derek Lewis.
Uie.6 I gave him the benefit of my opinion.
Uie.7 I gave him the benefit of my opin- **UC**
ion in strong language, but I did not instruct him because I was not, er, entitled to instruct him.
Uie.8 I was entitled to express my opinion **UC** and that is what I did.
P [25] : Uir.9 With respect, that is not answering the question of whether you threatened to overrule him.
H[26] : Uie.9 It’s dealing with the relevant point **TC** which was what I was entitled to do and what I was not entitled to do,
Uie.10 and I have dealt with this in detail **UC** before the House of Commons and before the select committee.

Table 2 summarises non-cooperative features, utterances and the degree of non-cooperation for each participant and for the whole fragment.

	P (ir)	H (ie)	Fragment
Interruptions	1	0	1
Overlaps	3	1	4
Grounding Failure	1	2	3
Unsolicited Comments	0	4	4
Topic Change	0	1	1
Total NCFs	5	8	13
Utterances	9	10	19
DNC	0.56	0.80	0.68

Table 2: Computing the DNC for dialogue (3)

The DNC was computed for all the political interviews in the corpus. Table 3 shows the val-

	Dialogue	Utterances	NCF	DNC
Adversarial Dialogues	1. Paxman v. Howard	54	30	0.56
	Paxman (IR)	24	13	0.54
	Howard (IE)	30	17	0.57
	2. Paxman v. Galloway	48	15	0.31
	Paxman (IR)	19	7	0.37
	Galloway (IE)	29	8	0.28
	4. O'Reilly v. Hartman	36	9	0.25
	O'Reilly (IR)	15	4	0.27
	Hartman (IE)	21	5	0.24
	8. Rather v. Bush (Turns 107-133)	40	19	0.48
Rather (IR)	18	8	0.44	
Bush (IE)	22	11	0.5	
Cooperative Dialogues	7. Keating v. Thatcher	57	0	0
	Keating (IR)	12	0	0
	Thatcher (IE)	45	0	0
	8. Brodie v. Blair (Turns 7-20)	31	2	0.06
	Brodie (IR)	9	1	0.11
	Blair (IE)	22	1	0.05

Table 3: DNC of political interviews in the corpus

ues obtained. Adversarial interviews have a large number of NCFs, thus a high value for the DNC. On the other hand, collaborative exchanges have low occurrence of NCFs (or none at all)⁷.

4 Discussion

There have been previous approaches to modeling dialogue on the basis that participants are not always fully cooperative. Jameson (1989) presents an extensive study for modeling bias, individual goals, projected image and belief ascription in conversation. User-model approaches are flexible to account for intricate situations but, as noted by Taylor et al. (1996), can lead to problems like infinite regress in nested beliefs. Taylor (1994) addressed non-cooperative dialogue behaviour by implementing CYNIC, a dialogue system able to generate and recognise deception; a notion of non-cooperation weaker than the one we address.

More recently, Traum (2008) brought attention to the need for computational accounts of dialogue situations in which a broader notion of cooperation is not assumed: e.g., intelligent tutoring systems, bargaining agents, role-playing training

⁷These results and the validity of DNC measure need further evaluation. We are currently performing two studies: one to determine inter-annotator agreement of the coding scheme for NCFs, and another to test how NCFs correlate to human judgements of non-cooperative conversational behaviour.

agents⁸. Traum's work on conflictive dialogue is mainly aimed at creating virtual humans with abilities to engage in adversarial dialogue. Traum et al. (2008) present a model of conversation strategies for negotiation, that includes variables representing trust, politeness and emotions, and a set of conversational strategies. Despite being adversarial in nature, the conversational scenarios are modeled by means of rules, that are followed by the interlocutors, according to the values of some of the variables. Hence, the dialogues are adversarial, but cooperative under our characterisation of linguistic non-cooperation, and it is not clear how effectively the model accounts for cases in which participants fail to follow the rules of a scenario.

4.1 Working Hypothesis

Finding a suitable model of non-cooperative dialogue involves bridging the gap between the theoretical aspects mentioned so far and the evidence in the empirical data of the previous section. Following Traum and Allen (1994), we base on the hypothesis that non-cooperative features result from decisions that participants make during the conversation, by considering the obligations imposed by the social activity and their individual goals, with an adequate configuration of the priorities for goals and obligations.

Thus, a participant with high priorities for individual goals might compromise the workings of a conversation by choosing contributions that go against the norms of the social activity. On the other hand, participants with higher priorities associated with obligations will favour contributions consistent with the rules of the social activity.

4.2 Research Methodology

For the next steps of the project, we will construct a model based on the hypothesis and test it by means of simulation⁹.

The construction of the model is a formalization of the working hypothesis, including rules for political interviews, goals, obligations, priorities and a dialogue management component. At the

⁸Traum also provides a list of "behaviours of interest", along the lines of the NCFs we identified above: e.g., unilateral topic shifts or topic maintenance, unhelpful criticism, withholding of information, lying, deception, antagonism.

⁹The use of simulation in dialogue modeling was pioneered by Power (1979). It suits our project better than alternatives (e.g., Wizard-of-Oz, dialogue systems), by making it easier to introduce modifications, do re-runs, and generate a large number of cases with different parameter settings.

moment of writing, we are investigating the line of research on obligation-driven dialogue modeling, initiated by Traum and Allen (1994) and developed further by Poesio and Traum (1998) and Kreutel and Matheson (2003).

For the simulation, DPs will be autonomous conversational agents with a cognitive state consisting of goals, a notion of their expected behaviour in a political interview, priorities, and some knowledge of the world. We are currently implementing a prototype based on EDIS (Matheson et al., 2000).

5 Conclusions

In this paper we presented an attempt to shed light on non-cooperation in dialogue by proposing a practical measure of the degree of linguistic non-cooperation in political interviews and a methodology towards a suitable computational model.

Acknowledgments

We would like to thank the NLG group at The Open University (especially Paul Piwek, Richard Power and Sandra Williams) for helpful discussion and comments on previous versions of this paper; and three anonymous reviewers for thoughtful feedback and suggestions.

References

- J.F. Allen and L.K. Schubert. 1991. The TRAINS project. TRAINS Technical Note 91-1. *Computer Science Dept. University of Rochester*.
- S. Attardo. 1997. Locutionary and perlocutionary cooperation: The perlocutionary cooperative principle. *Journal of Pragmatics*, 27(6):753–779.
- N. Blaylock and J. Allen. 2005. A collaborative problem-solving model of dialogue. In *Proceedings of the 6th SIGdial Workshop on Discourse and Dialogue*, pages 200–211, Lisbon, Portugal.
- Harry Bunt. 1994. Context and dialogue control. *THINK Quarterly*, 3.
- H.H. Clark and E.F. Schaefer. 1989. Contributing to discourse. *Cognitive science*, 13(2):259–294.
- P.R. Cohen and H.J. Levesque. 1991. Confirmations and joint action. In *Proceedings of the 12th International Joint Conference on Artificial Intelligence*, pages 951–957.
- J. Ginzburg. 1996. Interrogatives: Questions, facts and dialogue. *The handbook of contemporary semantic theory*, 5:359–423.
- H. P. Grice. 1975. Logic and conversation. *Syntax and Semantics*, 3:41–58.
- B.J. Grosz and C.L. Sidner. 1990. Plans for discourse. *Intentions in communication*, pages 417–444.
- J. Heritage. 1998. Conversation analysis and institutional talk. Analyzing distinctive turn-taking systems. In *Proceedings of the 6th International Congress of IADA, Tubingen, Niemeyer*.
- A. Jameson. 1989. But what will the listener think? Belief ascription and image maintenance in dialog. *User Models in Dialog Systems. Springer-Verlag*, pages 255–312.
- J. Kreutel and C. Matheson. 2003. Incremental information state updates in an obligation-driven dialogue model. *Logic Journal of IGPL*, 11(4):485.
- C. Matheson, M. Poesio, and D. Traum. 2000. Modelling grounding and discourse obligations using update rules. In *Proceedings of the 1st NAACL conference*, pages 1–8, San Francisco, CA, USA.
- M. Poesio and D. Traum. 1998. Towards an axiomatization of dialogue acts. In *Proceedings of the Twente Workshop on the Formal Semantics and Pragmatics of Dialogues*, pages 207–222.
- R. Power. 1979. The organisation of purposeful dialogues. *Linguistics*, 17:107–152.
- C. Reed and D. Long. 1997. Collaboration, cooperation and dialogue classification. *Working Notes of the IJCAI97 Workshop on Collaboration, Cooperation and Conflict in Dialog Systems, IJCAI 97*, pages 73–78.
- H. Sacks, E.A. Schegloff, and G. Jefferson. 1974. A simplest systematics for the organization of turn-taking for conversation. *Language*, pages 696–735.
- J.R. Searle. 1979. A Taxonomy of Illocutionary Acts. *Expression and meaning: studies in the theory of speech acts*, pages 1–29.
- J. A. Taylor, J. Carletta, and C. Mellish. 1996. Requirements for belief models in cooperative dialogue. *User Modeling and User-Adapted Interaction*, 6(1):23–68.
- J.A. Taylor. 1994. *A multi-agent planner for modelling dialogue*. Ph.D. Thesis, School of Cognitive and Computing Sciences, University of Sussex.
- D.R. Traum and J.F. Allen. 1994. Discourse obligations in dialogue processing. In *Proceedings of the 32nd annual meeting of ACL*, pages 1–8. Morristown, NJ, USA.
- D. Traum, W. Swartout, J. Gratch, and S. Marsella. 2008. A virtual human dialogue model for non-team interaction. *Recent Trends in Discourse and Dialogue. Springer*.
- D. Traum. 2008. Extended Abstract: Computational Models of Non-cooperative dialogue. In *Proceedings of LONDIAL 2008, the 12th Workshop on the Semantics and Pragmatics of Dialogue*, pages 11–14, London, UK.
- D. Walton and E. Krabbe. 1995. *Commitment in dialogue: Basic concepts of interpersonal reasoning*. State University of New York Press.

Towards Relational POMDPs for Adaptive Dialogue Management

Pierre Lison

Language Technology Lab
German Research Centre for Artificial Intelligence (DFKI GmbH)
Saarbrücken, Germany

Abstract

Open-ended spoken interactions are typically characterised by both structural complexity and high levels of uncertainty, making dialogue management in such settings a particularly challenging problem. Traditional approaches have focused on providing theoretical accounts for either the uncertainty or the complexity of spoken dialogue, but rarely considered the two issues simultaneously. This paper describes ongoing work on a new approach to dialogue management which attempts to fill this gap. We represent the interaction as a Partially Observable Markov Decision Process (POMDP) over a rich state space incorporating both dialogue, user, and environment models. The tractability of the resulting POMDP can be preserved using a mechanism for dynamically constraining the action space based on prior knowledge over locally relevant dialogue structures. These constraints are encoded in a small set of general rules expressed as a Markov Logic network. The first-order expressivity of Markov Logic enables us to leverage the rich relational structure of the problem and efficiently abstract over large regions of the state and action spaces.

1 Introduction

The development of spoken dialogue systems for rich, open-ended interactions raises a number of challenges, one of which is dialogue management. The role of dialogue management is to determine which communicative actions to take (i.e. what to say) given a goal and particular observations about the interaction and the current situation.

Dialogue managers have to face several issues. First, spoken dialogue systems must usually deal

with high levels of noise and uncertainty. These uncertainties may arise from speech recognition errors, limited grammar coverage, or from various linguistic and pragmatic ambiguities.

Second, open-ended dialogue is characteristically complex, and exhibits rich relational structures. Natural interactions should be adaptive to a variety of factors dependent on the interaction history, the general context, and the user preferences. As a consequence, the state space necessary to model the dynamics of the environment tends to be large and sparsely populated.

These two problems have typically been addressed separately in the literature. On the one hand, the issue of uncertainty in speech understanding is usually dealt using a range of probabilistic models combined with decision-theoretic planning. Among these, *Partially Observable Markov Decision Process* (POMDP) models have recently emerged as a unifying mathematical framework for dialogue management (Williams and Young, 2007; Lemon and Pietquin, 2007). POMDPs provide an explicit account for a wide range of uncertainties related to partial observability (noisy, incomplete spoken inputs) and stochastic action effects (the world may evolve in unpredictable ways after executing an action).

On the other hand, structural complexity is typically addressed with logic-based approaches. Some investigated topics in this paradigm are pragmatic interpretation (Thomason et al., 2006), dialogue structure (Asher and Lascarides, 2003), or collaborative planning (Kruijff et al., 2008). These approaches are able to model sophisticated dialogue behaviours, but at the expense of robustness and adaptivity. They generally assume complete observability and provide only a very limited account (if any) of uncertainties.

We are currently developing a hybrid approach which *simultaneously* tackles the uncertainty and complexity of dialogue management, based on a

POMDP framework. We present here our ongoing work on this issue. In this paper, we more specifically describe a new mechanism for dynamically constraining the space of possible actions available at a given time. Our aim is to use such mechanism to significantly reduce the search space and therefore make the planning problem globally more tractable. This is performed in two consecutive steps. We first structure the state space using *Markov Logic Networks*, a first-order probabilistic language. Prior pragmatic knowledge about dialogue structure is then exploited to derive the set of dialogue actions which are locally admissible or relevant, and prune all irrelevant ones. The first-order expressivity of Markov Logic Networks allows us to easily specify the constraints via a small set of general rules which abstract over large regions of the state and action spaces.

Our long-term goal is to develop a unified framework for adaptive dialogue management in rich, open-ended interactional settings.

This paper is structured as follows. Section 2 lays down the formal foundations of our work, by describing dialogue management as a POMDP problem. We then describe in Section 3 our approach to POMDP planning with control knowledge using Markov Logic rules. Section 4 discusses some further aspects of our approach and its relation to existing work, followed by the conclusion in Section 5.

2 Background

2.1 Partially Observable Markov Decision Processes (POMDPs)

POMDPs are a mathematical model for sequential decision-making in partially observable environments. It provides a powerful framework for control problems which combine partial observability, uncertain action effects, incomplete knowledge of the environment dynamics and multiple, potentially conflicting objectives.

Via reinforcement learning, it is possible to automatically *learn* near-optimal action policies given a POMDP model combined with real or simulated user data (Schatzmann et al., 2007).

2.1.1 Formal definition

A POMDP is a tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{Z}, T, \Omega, R \rangle$, where:

- \mathcal{S} is the **state space**, which is the model of the world from the agent’s viewpoint. It is defined as a set of mutually exclusive states.

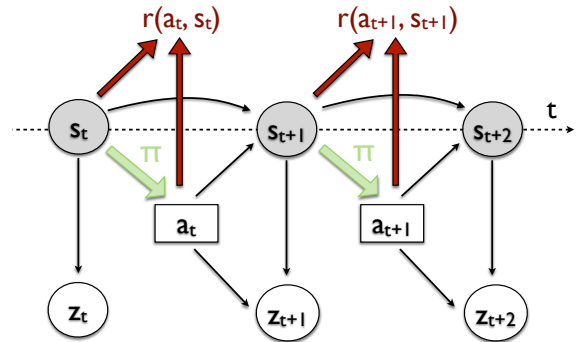


Figure 1: *Bayesian decision network* corresponding to the POMDP model. Hidden variables are greyed. Actions are represented as rectangles to stress that they are system actions rather than observed variables. Arcs into circular nodes express influence, whereas arcs into squared nodes are informational. For readability, only one state is shown at each time step, but it should be noted that the policy π is function of the full belief state rather than a single (unobservable) state.

- \mathcal{A} is the **action space**: the set of possible actions at the disposal of the agent.
- \mathcal{Z} is the **observation space**: the set of observations which can be captured by the agent. They correspond to features of the environment which can be directly perceived by the agent’s sensors.
- T is the **transition function**, defined as $T : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$, where $T(s, a, s') = P(s' | s, a)$ is the probability of reaching state s' from state s if action a is performed.
- Ω is the **observation function**, defined as $\Omega : \mathcal{Z} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$, with $\Omega(z, a, s') = P(z | a, s')$, i.e. the probability of observing z after performing a and being now in state s' .
- R is the **reward function**, defined as $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathfrak{R}$, $R(s, a)$ encodes the utility for the agent to perform the action a while in state s . It is therefore a model for the goals or preferences of the agent.

A graphical illustration of a POMDP model as a Bayesian decision network is provided in Fig. 1.

In addition, a POMDP can include additional parameters such as the horizon of the agent (num-

ber of look-ahead steps), and the discount factor (weighting scheme for non-immediate rewards).

2.1.2 Beliefs and belief update

A key idea of POMDP is the assumption that the state of the world is not directly accessible, and can only be inferred via observation. Such uncertainty is expressed in the **belief state** b , which is a probability distribution over possible states, that is: $b : \mathcal{S} \rightarrow [0, 1]$. The belief state for a state space of cardinality n is therefore represented in a real-valued simplex of dimension $(n-1)$.

This belief state is dynamically updated before executing each action. The belief state update operates as follows. At a given time step t , the agent is in some unobserved state $s_t = s \in \mathcal{S}$. The probability of being in state s at time t is written as $b_t(s)$. Based on the current belief state b_t , the agent selects an action a_t , receives a reward $R(s, a_t)$ and transitions to a new (unobserved) state $s_{t+1} = s'$, where s_{t+1} depends only on s_t and a_t . The agent then receives a new observation o_{t+1} which is dependent on s_{t+1} and a_t .

Finally, the belief distribution b_t is updated, based on o_{t+1} and a_t as follows¹.

$$b_{t+1}(s') = P(s'|o_{t+1}, a_t, b_t) \quad (1)$$

$$= \frac{P(o_{t+1}|s', a_t, b_t)P(s'|a_t, b_t)}{P(o_{t+1}|a_t, b_t)} \quad (2)$$

$$= \frac{P(o_{t+1}|s', a_t) \sum_{s \in \mathcal{S}} P(s'|a_t, s)P(s|a_t, b_t)}{P(o_{t+1}|a_t, b_t)} \quad (3)$$

$$= \alpha \Omega(o_{t+1}, s', a_t) \sum_{s \in \mathcal{S}} T(s, a_t, s')b_t(s) \quad (4)$$

where α is a normalisation constant. An initial belief state b_0 must be specified at runtime as a POMDP parameter when initialising the system.

2.1.3 POMDP policies

Given a POMDP model $\langle \mathcal{S}, \mathcal{A}, \mathcal{Z}, T, Z, R \rangle$, the agent should execute at each time-step the action which maximises its expected cumulative reward over the horizon. The function $\pi : \mathcal{B} \rightarrow \mathcal{A}$ defines a *policy*, which determines the action to perform for each point of the belief space.

The expected reward for policy π starting from belief b is defined as:

$$J^\pi(b) = E \left[\sum_{t=0}^h \gamma^t R(s_t, a_t) \mid b, \pi \right] \quad (5)$$

¹As a notational shorthand, we write $P(s_t=s)$ as $P(s)$ and $P(s_{t+1}=s')$ as $P(s')$.

The optimal policy π^* is then obtained by optimizing the long-term reward, starting from b_0 :

$$\pi^* = \operatorname{argmax}_{\pi} J^\pi(b_0) \quad (6)$$

The optimal policy π^* yields the highest expected reward value for each possible belief state. This value is compactly represented by the optimal value function, noted V^* , which is a solution to the Bellman optimality equation (Bellman, 1957).

Numerous algorithms for (offline) policy optimisation and (online) planning are available. For large spaces, exact optimisation is impossible and approximate methods must be used, see for instance grid-based (Thomson and Young, 2009) or point-based (Pineau et al., 2006) techniques.

2.2 POMDP-based dialogue management

Dialogue management can be easily cast as a POMDP problem, with the *state space* being a compact representation of the interaction, the *action space* being a set of dialogue moves, the *observation space* representing speech recognition hypotheses, the *transition function* defining the dynamics of the interaction (which user reaction is to be expected after a particular dialogue move), and the *observation function* describing a “sensor model” between observed speech recognition hypotheses and actual utterances. Finally, the *reward function* encodes the utility of dialogue policies – it typically assigns a big positive reward if a long-term goal has been reached (e.g. the retrieval of some important information), and small negative rewards for minor “inconveniences” (e.g. prompting the user to repeat or asking for confirmations).

Our long-term aim is to apply such POMDP framework to a rich dialogue domain for human-robot interaction (Kruijff et al., 2010). These interactions are typically open-ended, relatively long, include high levels of noise, and require complex state and action spaces. Furthermore, the dialogue system also needs to be *adaptive* to its user (attributed beliefs and intentions, attitude, attentional state) and to the current situation (currently perceived entities and events).

As a consequence, the state space must be expanded to include these knowledge sources. Belief monitoring is then used to continuously update the belief state based on perceptual inputs (see also (Bohus and Horvitz, 2009) for an overview of techniques to extract such information). These requirements can only be fulfilled if we address the

“curse of dimensionality” characteristic of traditional POMDP models. The next section provides a tentative answer.

3 Approach

3.1 Control knowledge

Classical approaches to POMDP planning operate directly on the full action space and select the next action to perform based on the maximisation of the expected cumulative reward over the specified horizon. Such approaches can be used in small-scale domains with a limited action space, but quickly become intractable for larger ones, as the planning time increases exponentially with the size of the action space. Significant planning time is therefore spent on actions which should be directly discarded as irrelevant². Dismissing these actions *before* planning could therefore provide important computational gains.

Instead of a direct policy optimisation over the full action space, our approach formalises action selection as a *two-step* process. As a first step, a set of *relevant dialogue moves* is constructed from the full action space. The POMDP planner then computes the optimal (highest-reward) action on this reduced action space in a second step.

Such an approach is able to significantly reduce the dimensionality of the dialogue management problem by taking advantage of prior knowledge about the expected relational structure of spoken dialogue. This prior knowledge is to be encoded in a set of general rules describing the admissible dialogue moves in a particular situation.

How can we express such rules? POMDPs are usually modeled with Bayesian networks which are inherently propositional. Encoding such rules in a propositional framework requires a distinct rule for every possible state and action instance. This is not a feasible approach. We therefore need a first order (probabilistic) language able to express generalities over large regions of the state action spaces. Markov Logic is such a language.

3.2 Markov Logic Networks (MLNs)

Markov Logic combines first-order logic and probabilistic graphical models in a unified representation (Richardson and Domingos, 2006). A

²For instance, an agent hearing a user command such as “Please take the mug on your left” might spend a lot of planning time calculating the expected future reward of dialogue moves such as “Is the box green?” or “Your name is John”, which are irrelevant to the situation.

Markov Logic Network L is a set of pairs (F_i, w_i) , where F_i is a formula in first-order logic and w_i is a real number representing the formula weight.

A Markov Logic Network L can be seen as a *template* for constructing markov networks³. To construct a markov network from L , one has to provide an additional set of constants $C = \{c_1, c_2, \dots, c_{|C|}\}$. The resulting markov network is called a *ground markov network* and is written $M_{L,C}$. The ground markov network contains one feature for each possible grounding of a first-order formula in L , with the corresponding weight. The technical details of the construction of $M_{L,C}$ from the two sets L and C is explained in several papers, see e.g. (Richardson and Domingos, 2006).

Once the markov network $M_{L,C}$ is constructed, it can be exploited to perform *inference* over arbitrary queries. Efficient probabilistic inference algorithms such as Markov Chain Monte Carlo (MCMC) or other sampling techniques can then be used to this end (Poon and Domingos, 2006).

3.3 States and actions as relational structures

The specification of Markov Logic rules applying over complete regions of the state and action spaces (instead of over single instances) requires an explicit relational structure over these spaces.

This is realised by factoring the state and action spaces into a set of distinct, conditionally independent features. A state s can be expanded into a tuple $\langle f_1, f_2, \dots, f_n \rangle$, where each sub-state f_i is assigned a value from a set $\{v_1, v_2, \dots, v_m\}$. Such structure can be expressed in first-order logic with a binary predicate $f_i(s, v_j)$ for each sub-state f_i , where v_j is the value of the sub-state f_i in s . The same type of structure can be defined over actions. This factoring leads to a relational structure of arbitrary complexity, compactly represented by a set of unary and binary predicates.

For instance, (Young et al., 2010) factors each dialogue state into three independent parts $s = \langle s_u, a_u, s_d \rangle$, where s_u is the user goal, a_u the last user move, and s_d the dialogue history. These can be expressed in Markov Logic with predicates such as $\text{UserGoal}(s, s_u)$, $\text{LastUserMove}(s, a_u)$, or $\text{History}(s, s_d)$.

³Markov networks are undirected graphical models.

3.4 Relevant action space

For a given state s , the relevant action space $RelMoves(\mathcal{A}, s)$ is defined as:

$$\{a_m : a_m \in \mathcal{A} \wedge \text{RelevantMove}(a_m, s)\} \quad (7)$$

The truth-value of the predicate $\text{RelevantMove}(a_m, s)$ is determined using a set of Markov Logic rules dependent on both the state s and the action a_m . For a given state s , the relevant action space is constructed via probabilistic inference, by estimating the probability $P(\text{RelevantMove}(a_m, s))$ for each action a_m , and selecting the subset of actions for which the probability is above a given threshold.

Eq. 8 provides a simple example of such Markov Logic rule:

$$\text{LastUserMove}(s, a_u) \wedge \text{PolarQuestion}(a_u) \wedge \text{YesNoAnswer}(a_m) \rightarrow \text{RelevantMove}(a_m, s) \quad (8)$$

It defines an admissible dialogue move for a situation where the user asks a polar question to the agent (e.g. “do you see my hand?”). The rule specifies that, if a state s contains a_u as last user move, and if a_u is a polar question, then an answer a_m of type yes-no is a relevant dialogue move for the agent. This rule is (implicitly) universally quantified over s , a_u and a_m .

Each of these Markov Logic rules has a weight attached to it, expressing the strength of the implication. A rule with infinite weight and satisfied premises will lead to a relevant move with probability 1. Softer weights can be used to describe moves which are less relevant but still possible in a particular context. These weights can either be encoded by hand or learned from data (how to perform this efficiently remains an open question).

3.5 Rules application on POMDP belief state

The previous section assumed that the state s is known. But the real state of a POMDP is never directly accessible. The rules we just described must therefore be applied on the belief state. Ultimately, we want to define a function $Rel : \mathfrak{R}^n \rightarrow \mathcal{P}(\mathcal{A})$, which takes as input a point in the belief space and outputs a set of relevant moves. For efficiency reasons, this function can be precomputed offline, by segmenting the state space into distinct regions and assigning a set of relevant moves to each region. The function can then be directly called at runtime by the planning algorithm.

Due to the high dimensionality of the belief space, the above function must be approximated to remain tractable. One way to perform this approximation is to extract, for belief state b , a set S_m of m most likely states, and compute the set of relevant moves for each of them. We then define the global probability estimate of a being a relevant move given b as such:

$$P(\text{RelevantMove}(a) | b, a) \approx \sum_{s \in S_m} P(\text{RelevantMove}(a, s) | s, a) \times b(s) \quad (9)$$

In the limit where $m \rightarrow |S|$, the error margin on the approximation tends to zero.

4 Discussion

4.1 General comments

It is worth noting that the mechanism we just outlined does not intend to *replace* the existing POMDP planning and optimisation algorithms, but rather *complements* them. Each step serves a different purpose: the action space reduction provides an answer to the question “Is this action relevant?”, while the policy optimisation seeks to answer “Is this action useful?”. We believe that such distinction between relevance and usefulness is important and will prove to be beneficial in terms of tractability.

It is also useful to notice that the Markov Logic rules we described provides a “positive” definition of the action space. The rules were applied to produce an exhaustive list of all admissible actions given a state, all actions outside this list being *de facto* labelled as non-admissible. But the rules can also provide a “negative” definition of the action space. That is, instead of generating an exhaustive list of possible actions, the dialogue system can initially consider all actions as admissible, and the rules can then be used to prune this action space by removing irrelevant moves.

The choice of action filter depends mainly on the size of the dialogue domain and the availability of prior domain knowledge. A “positive” filter is a necessity for large dialogue domains, as the action space is likely to grow exponentially with the domain size and become untractable. But the positive definition of the action space is also significantly more expensive for the dialogue developer. There is therefore a trade-off between the costs of tractability issues, and the costs of dialogue domain modelling.

4.2 Related Work

There is a substantial body of existing work in the POMDP literature about the exploitation of the problem structure to tackle the curse of dimensionality (Poupart, 2005; Young et al., 2010), but the vast majority of these approaches retain a propositional structure. A few more theoretical papers also describe first-order MDPs (Wang et al., 2007), and recent work on Markov Logic has extended the MLN formalism to include some decision-theoretic concepts (Nath and Domingos, 2009). To the author’s knowledge, none of these ideas have been applied to dialogue management.

5 Conclusions

This paper described a new approach to exploit relational models of dialogue structure for controlling the action space in POMDPs. This approach is part of an ongoing work to develop a unified framework for adaptive dialogue management in rich, open-ended interactional settings. The dialogue manager is being implemented as part of a larger cognitive architecture for talking robots.

Besides the implementation, future work will focus on refining the theoretical foundations of relational POMDPs for dialogue (including how to specify the transition, observation and reward functions in such a relational framework), as well as investigating the use of reinforcement learning for policy optimisation based on simulated data.

References

- N. Asher and A. Lascarides. 2003. *Logics of Conversation*. Cambridge University Press.
- R. Bellman. 1957. *Dynamic Programming*. Princeton University Press.
- Dan Bohus and Eric Horvitz. 2009. Dialog in the open world: platform and applications. In *ICMI-MLMI '09: Proceedings of the 2009 international conference on Multimodal interfaces*, pages 31–38, New York, NY, USA. ACM.
- G.J.M. Kruijff, M. Brenner, and N.A. Hawes. 2008. Continual planning for cross-modal situated clarification in human-robot interaction. In *Proceedings of the 17th International Symposium on Robot and Human Interactive Communication (RO-MAN 2008)*, Munich, Germany.
- G.-J. M. Kruijff, P. Lison, T. Benjamin, H. Jacobsson, H. Zender, and I. Kruijff-Korbayova. 2010. Situated dialogue processing for human-robot interaction. In H. I. Christensen, A. Sloman, G.-J. M. Kruijff, and J. Wyatt, editors, *Cognitive Systems*. Springer Verlag. (in press).
- O. Lemon and O. Pietquin. 2007. Machine learning for spoken dialogue systems. In *Proceedings of the European Conference on Speech Communication and Technologies (Interspeech'07)*, pages 2685–2688, Anvers (Belgium), August.
- A. Nath and P. Domingos. 2009. A language for relational decision theory. In *Proceedings of the International Workshop on Statistical Relational Learning*.
- J. Pineau, G. Gordon, and S. Thrun. 2006. Anytime point-based approximations for large pomdps. *Artificial Intelligence Research*, 27(1):335–380.
- H. Poon and P. Domingos. 2006. Sound and efficient inference with probabilistic and deterministic dependencies. In *AAAI'06: Proceedings of the 21st national conference on Artificial intelligence*, pages 458–463. AAAI Press.
- P. Poupart. 2005. *Exploiting structure to efficiently solve large scale partially observable markov decision processes*. Ph.D. thesis, University of Toronto, Toronto, Canada.
- M. Richardson and P. Domingos. 2006. Markov logic networks. *Machine Learning*, 62(1-2):107–136.
- Jost Schatzmann, Blaise Thomson, Karl Weilhammer, Hui Ye, and Steve Young. 2007. Agenda-based user simulation for bootstrapping a POMDP dialogue system. In *HLT '07: Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics on Human Language Technologies*, pages 149–152, Rochester, New York, April. Association for Computational Linguistics.
- R. Thomason, M. Stone, and D. DeVault. 2006. Enlightened update: A computational architecture for presupposition and other pragmatic phenomena. In Donna Byron, Craige Roberts, and Scott Schwenter, editors, *Presupposition Accommodation*. Ohio State Pragmatics Initiative.
- B. Thomson and S. Young. 2009. Bayesian update of dialogue state: A pomdp framework for spoken dialogue systems. *Computer Speech & Language*, August.
- Ch. Wang, S. Joshi, and R. Khardon. 2007. First order decision diagrams for relational mdps. In *IJCAI'07: Proceedings of the 20th international joint conference on Artificial intelligence*, pages 1095–1100, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.
- J. Williams and S. Young. 2007. Partially observable markov decision processes for spoken dialog systems. *Computer Speech and Language*, 21(2):231–422.
- S. Young, M. Gašić, S. Keizer, F. Mairesse, J. Schatzmann, B. Thomson, and K. Yu. 2010. The hidden information state model: A practical framework for pomdp-based spoken dialogue management. *Computer Speech & Language*, 24(2):150–174.

WSD as a Distributed Constraint Optimization Problem

Siva Reddy
IIIT Hyderabad
India

gvsreddy@students.iiit.ac.in

Abhilash Inumella
IIIT Hyderabad
India

abhilashi@students.iiit.ac.in

Abstract

This work models Word Sense Disambiguation (WSD) problem as a Distributed Constraint Optimization Problem (DCOP). To model WSD as a DCOP, we view information from various knowledge sources as constraints. DCOP algorithms have the remarkable property to jointly maximize over a wide range of utility functions associated with these constraints. We show how utility functions can be designed for various knowledge sources. For the purpose of evaluation, we modelled all words WSD as a simple DCOP problem. The results are competitive with state-of-art knowledge based systems.

1 Introduction

Words in a language may carry more than one sense. The correct sense of a word can be identified based on the context in which it occurs. In the sentence, *He took all his money from the bank*, *bank* refers to a *financial institution* sense instead of other possibilities like *the edge of river* sense. Given a word and its possible senses, as defined by a dictionary, the problem of Word Sense Disambiguation (WSD) can be defined as the task of assigning the most appropriate sense to the word within a given context.

WSD is one of the oldest problems in computational linguistics which dates back to early 1950's. A range of knowledge sources have been found to be useful for WSD. (Agirre and Stevenson, 2006; Agirre and Martínez, 2001; McRoy, 1992; Hirst, 1987) highlight the importance of various knowledge sources like part of speech, morphology, collocations, lexical knowledge base (sense taxonomy, gloss), sub-categorization, semantic word associations, selectional preferences,

semantic roles, domain, topical word associations, frequency of senses, collocations, domain knowledge. etc. Methods for WSD exploit information from one or more of these knowledge sources.

Supervised approaches like (Yarowsky and Florian, 2002; Lee and Ng, 2002; Martínez et al., 2002; Stevenson and Wilks, 2001) used collective information from various knowledge sources to perform disambiguation. Information from various knowledge sources is encoded in the form of a feature vector and models were built by training on sense-tagged corpora. These approaches pose WSD as a classification problem. They crucially rely on hand-tagged sense corpora which is hard to obtain. Systems that do not need hand-tagging have also been proposed. Agirre and Martínez (Agirre and Martínez, 2001) evaluated the contribution of each knowledge source separately. However, this does not combine information from more than one knowledge source.

In any case, little effort has been made in formalizing the way in which information from various knowledge sources can be collectively used within a single framework: a framework that allows interaction of evidence from various knowledge sources to arrive at a global optimal solution.

Here we present a way for modelling information from various knowledge sources in a multi agent setting called distributed constraint optimization problem (DCOP). In DCOP, agents have constraints on their values and each constraint has a utility associated with it. The agents communicate with each other and choose values such that a global optimum solution (maximum utility) is attained. We aim to solve WSD by modelling it as a DCOP.

To the best of our knowledge, ours is the first attempt to model WSD as a DCOP. In DCOP framework, information from various knowledge sources can be used combinedly to perform WSD.

In section 2, we give a brief introduction of

DCOP. Section 3 describes modelling WSD as a DCOP. Utility functions for various knowledge sources are described in section 4. In section 5, we conduct a simple experiment by modelling all-words WSD problem as a DCOP and perform disambiguation on Senseval-2 (Cotton et al., 2001) and Senseval-3 (Mihalcea and Edmonds, 2004) data-set of all-words task. Next follow the sections on related work, discussion, future work and conclusion.

2 Distributed Constraint Optimization Problem (DCOP)

A DCOP (Modi, 2003; Modi et al., 2004) consists of n variables $V = x_1, x_2, \dots, x_n$ each assigned to an agent, where the values of the variables are taken from finite, discrete domains D_1, D_2, \dots, D_n respectively. Only the agent has knowledge and control over values assigned to variables associated to it. The goal for the agents is to choose values for variables such that a given global objective function is maximized. The objective function is described as the summation over a set of utility functions.

DCOP can be formalized as a tuple (A, V, D, C, F) where

- $A = \{a_1, a_2, \dots, a_n\}$ is a set of n agents,
- $V = \{x_1, x_2, \dots, x_n\}$ is a set of n variables, each one associated to an agent,
- $D = \{D_1, D_2, \dots, D_n\}$ is a set of finite and discrete domains each one associated to the corresponding variable,
- $C = \{f_k : D_i \times D_j \times \dots \times D_m \rightarrow \mathfrak{R}\}$ is a set of constraints described by various utility functions f_k . The utility function f_k is defined over a subset of variables V . The domain of f_k represent the constraints C_{f_k} and $f_k(c)$ represents the utility associated with the constraint c , where $c \in C_{f_k}$.
- $F = \sum_k z_k \cdot f_k$ is the objective function to be maximized where z_k is the weight of the corresponding utility function f_k

An agent is allowed to communicate only with its neighbours. Agents communicate with each other to agree upon a solution which maximizes the objective function.

3 WSD as a DCOP

Given a sequence of words $W = \{w_1, w_2, \dots, w_n\}$ with corresponding admissible senses $D_{w_i} = \{s_{w_i}^1, s_{w_i}^2, \dots\}$, we model WSD as DCOP as follows.

3.1 Agents

Each word w_i is treated as an agent. The agent (word) has knowledge and control of its values (senses).

3.2 Variables

Sense of a word varies and it is the one to be determined. We define the sense of a word as its variable. Each agent w_i is associated with the variable s_{w_i} . The value assigned to this variable indicates the sense assigned by the algorithm.

3.3 Domains

Senses of a word are finite in number. The set of senses D_{w_i} , is the domain of the variable s_{w_i} .

3.4 Constraints

A constraint specifies a particular configuration of the agents involved in its definition and has a utility associated with it. For e.g. If c_{ij} is a constraint defined on agents w_i and w_j , then c_{ij} refers to a particular instantiation of w_i and w_j , say $w_i = s_{w_i}^p$ and $w_j = s_{w_j}^q$.

A utility function $f_k : C_{f_k} \rightarrow \mathfrak{R}$ denote a set of constraints $C_{f_k} = \{D_{w_i} \times D_{w_j} \dots D_{w_m}\}$, defined on the agents $w_i, w_j \dots w_m$ and also the utilities associated with the constraints. We model information from each knowledge source as a utility function. In section 4, we describe in detail about this modelling.

3.5 Objective function

As already stated, various knowledge sources are identified to be useful for WSD. It is desirable to use information from these sources collectively, to perform disambiguation. DCOP provides such framework where an objective function is defined over all the knowledge sources (f_k) as below

$$F = \sum_k z_k \cdot f_k$$

where F denotes the total utility associated with a solution and z_k is the weight given to a knowledge source i.e. information from various sources

can be weighted. (Note: It is desirable to normalize utility functions of different knowledge sources in order to compare them.)

Every agent (word) choose its value (sense) in a such a way that the objective function (global solution) is maximized. This way an agent is assigned a best value which is the target sense in our case.

4 Modelling information from various knowledge sources

In this section, we discuss the modelling of information from various knowledge sources.

4.1 Part-of-speech (POS)

Consider the word *play*. It has 47 senses out of which only 17 senses correspond to *noun* category. Based on the POS information of a word w_i , its domain D_{w_i} is restricted accordingly.

4.2 Morphology

Noun *orange* has at least two senses, one corresponding to *a color* and other to *a fruit*. But plural form of this word *oranges* can only be used in the *fruit* sense. Depending upon the morphological information of a word w_i , its domain D_{w_i} can be restricted.

4.3 Domain information

In the sports domain, *cricket* likely refers to *a game* than *an insect*. Such information can be captured using a unary utility function defined for every word. If the sense distributions of a word w_i are known, a function $f : D_{w_i} \rightarrow \mathfrak{R}$ is defined which return higher utility for the senses favoured by the domain than to the other senses.

4.4 Sense Relatedness

Sense relatedness between senses of two words w_i, w_j is captured by a function $f : D_{w_i} \times D_{w_j} \rightarrow \mathfrak{R}$ where f returns sense relatedness (utility) between senses based on sense taxonomy and gloss overlaps.

4.5 Discourse

Discourse constraints can be modelled using a n-ary function. For instance, to the extent one sense per discourse (Gale et al., 1992) holds true, higher utility can be returned to the solutions which favour same sense to all the occurrences of a word in a given discourse. This information can be modeled as follows: If $w_i, w_j, \dots w_m$ are

the occurrences of a same word, a function $f : D_i \times D_j \times \dots D_m \rightarrow \mathfrak{R}$ is defined which returns higher utility when $s_{w_i} = s_{w_j} = \dots s_{w_m}$ and for the rest of the combinations it returns lower utility.

4.6 Collocations

Collocations of a word are known to provide strong evidence for identifying correct sense of the word. For example: if in a given context *bank* co-occur with *money*, it is likely that *bank* refers to *financial institution* sense rather than *the edge of a river* sense. The word *cancer* has at least two senses, one corresponding to the astrological sign and the other a disease. But its derived form *cancerous* can only be used in disease sense. When the words *cancer* and *cancerous* co-occur in a discourse, it is likely that the word *cancer* refers to *disease sense*.

Most supervised systems work through collocations to identify correct sense of a word. If a word w_i co-occurs with its collocate v , collocational information from v can be modeled by using the following function

$$coll_inform_v_{w_i} : D_{w_i} \rightarrow \mathfrak{R}$$

where $coll_inform_v_{w_i}$ returns high utility to collocationally preferred senses of w_i than other senses.

Collocations can also be modeled by assigning more than one variable to the agents or by adding a dummy agent which gives collocational information but in view of simplicity we do not go into those details.

Topical word associations, semantic word associations, selectional preferences can also be modeled similar to collocations. Complex information involving more than two entities can be modelled by using n-ary utility functions.

5 Experiment: DCOP based All Words WSD

We carried out a simple experiment to test the effectiveness of DCOP algorithm. We conducted our experiment in an all words setting and used only WordNet (Fellbaum, 1998) based relatedness measures as knowledge source so that results can be compared with earlier state-of-art knowledge-based WSD systems like (Agirre and Soroa, 2009; Sinha and Mihalcea, 2007) which used similar knowledge sources as ours.

Our method performs disambiguation on sentence by sentence basis. A utility function based on semantic relatedness is defined for every pair of words falling in a particular window size. Restricting utility functions to a window size reduces the number of constraints. An objective function is defined as sum of these restricted utility functions over the entire sentence and thus allowing information flow across all the words. Hence, a DCOP algorithm which aims to maximize this objective function leads to a globally optimal solution.

In our experiments, we used the best similarity measure settings of (Sinha and Mihalcea, 2007) which is a sum of normalized similarity measures jcn, lch and lesk. We used Distributed Pseudotree Optimization Procedure (DPOP) algorithm (Petcu and Faltings, 2005), which solves DCOP using linear number of messages among agents. The implementation provided with the open source toolkit FRODO¹ (Léauté et al., 2009) is used.

5.1 Data

To compare our results, we ran our experiments on SENSEVAL-2 and SENSEVAL -3 English all-words data sets.

5.2 Results

Table 1 shows results of our experiments. All these results are carried out using a window size of four. Ideally, precision and recall values are expected to be equal in our setting. But in certain cases, the tool we used, FRODO, failed to find a solution with the available memory resources.

Results show that our system performs consistently better than (Sinha and Mihalcea, 2007) which uses exactly same knowledge sources as used by us (with an exception of adverbs in Senseval-2). This shows that DCOP algorithm perform better than page-rank algorithm used in their graph based setting. Thus, for knowledge-based WSD, DCOP framework is a potential alternative to graph based models.

Table 1 also shows the system (Agirre and Soroa, 2009), which obtained best results for knowledge based WSD. A direct comparison between this and our system is not quantitative since they used additional knowledge such as extended WordNet relations (Mihalcea and

Moldovan, 2001) and sense disambiguated gloss present in WordNet3.0.

Senseval-2 All Words data set					
	noun	verb	adj	adv	all
P_dcop	67.85	37.37	62.72	56.87	58.63
R_dcop	66.44	35.47	61.28	56.65	57.09
F_dcop	67.14	36.39	61.99	56.76	57.85
P_Sinha07	67.73	36.05	62.21	60.47	58.83
R_Sinha07	65.63	32.20	61.42	60.23	56.37
F_Sinha07	66.24	34.07	61.81	60.35	57.57
Agirre09	70.40	38.90	58.30	70.1	58.6
MFS	71.2	39.0	61.1	75.4	60.1
Senseval-3 All Words data set					
P_dcop	62.31	43.48	57.14	100	54.68
R_dcop	60.97	42.81	55.17	100	53.51
F_dcop	61.63	43.14	56.14	100	54.09
P_Sinha07	61.22	45.18	54.79	100	54.86
R_Sinha07	60.45	40.57	54.14	100	52.40
F_Sinha07	60.83	42.75	54.46	100	53.60
Agirre09	64.1	46.9	62.6	92.9	57.4
MFS	69.3	53.6	63.7	92.9	62.3

Table 1: Evaluation results on Senseval-2 and Senseval-3 data-set of all words task.

5.3 Performance analysis

We conducted our experiment on a computer with two 2.94 GHz process and 2 GB memory. Our algorithm just took 5 minutes 31 seconds on Senseval-2 data set, and 5 minutes 19 seconds on Senseval-3 data set. This is a singable reduction compared to execution time of page rank algorithms employed in both Sinha07 and Agirre09. In Agirre09, it falls in the range 30 to 180 minutes on much powerful system with 16 GB memory having four 2.66 GHz processors. On our system, time taken by the page rank algorithm in (Sinha and Mihalcea, 2007) is 11 minutes when executed on Senseval-2 data set.

Since DCOP algorithms are truly distributed in nature the execution times can be further reduced by running them parallely on multiple processors.

6 Related work

Earlier approaches to WSD which encoded information from variety of knowledge sources can be classified as follows:

- Supervised approaches: Most of the supervised systems (Yarowsky and Florian, 2002;

¹<http://liawww.epfl.ch/frodo/>

Lee and Ng, 2002; Martínez et al., 2002; Stevenson and Wilks, 2001) rely on the sense tagged data. These are mainly discriminative or aggregative models which essentially pose WSD a classification problem. Discriminative models aim to identify the most informative feature and aggregative models make their decisions by combining all features. They disambiguate word by word and do not collectively disambiguate whole context and thereby do not capture all the relationships (e.g sense relatedness) among all the words. Further, they lack the ability to directly represent constraints like one sense per discourse.

- Graph based approaches: These approaches crucially rely on lexical knowledge base. Graph-based WSD approaches (Agirre and Soroa, 2009; Sinha and Mihalcea, 2007) perform disambiguation over a graph composed of senses (nodes) and relations between pairs of senses (edges). The edge weights encode information from a lexical knowledge base but lack an efficient way of modelling information from other knowledge sources like collocational information, selectional preferences, domain information, discourse. Also, the edges represent binary utility functions defined over two entities which lacks the ability to encode ternary, and in general, any N-ary utility functions.

7 Discussion

This framework provides a convenient way of integrating information from various knowledge sources by defining their utility functions. Information from different knowledge sources can be weighed based on the setting at hand. For example, in a domain specific WSD setting, sense distributions play a crucial role. The utility function corresponding to the sense distributions can be weighed higher in order to take advantage of domain information. Also, different combination of weights can be tried out for a given setting. Thus for a given WSD setting, this framework allows us to find 1) the impact of each knowledge source individually 2) the best combination of knowledge sources.

Limitations of DCOP algorithms: Solving DCOPs is NP-hard. A variety of search algorithms have therefore been developed to solve DCOPs

(Mailler and Lesser, 2004; Modi et al., 2004; Petcu and Faltings, 2005). As the number of constraints or words increase, the search space increases thereby increasing the time and memory bounds to solve them. Also DCOP algorithms exhibit a trade-off between memory used and number of messages communicated between agents. DPOP (Petcu and Faltings, 2005) use linear number of messages but requires exponential memory whereas ADOPT (Modi et al., 2004) exhibits linear memory complexity but exchange exponential number of messages. So it is crucial to choose a suitable algorithm based on the problem at hand.

8 Future Work

In our experiment, we only used relatedness based utility functions derived from WordNet. Effect of other knowledge sources remains to be evaluated individually and in combination. The best possible combination of weights of knowledge sources is yet to be engineered. Which DCOP algorithm performs better WSD and when has to be explored.

9 Conclusion

We initiated a new line of investigation into WSD by modelling it in a distributed constraint optimization framework. We showed that this framework is powerful enough to encode information from various knowledge sources. Our experimental results show that a simple DCOP based model encoding just word similarity constraints performs comparably with the state-of-the-art knowledge based WSD systems.

Acknowledgement

We would like to thank *Prof. Rajeev Sangal* and *Asrar Ahmed* for their support in coming up with this work.

References

- Eneko Agirre and David Martínez. 2001. Knowledge sources for word sense disambiguation. In *Text, Speech and Dialogue, 4th International Conference, TSD 2001, Zelezna Ruda, Czech Republic, September 11-13, 2001*, Lecture Notes in Computer Science, pages 1–10. Springer.
- Eneko Agirre and Aitor Soroa. 2009. Personalizing pagerank for word sense disambiguation. In *EACL '09: Proceedings of the 12th Conference of the European Chapter of the Association for Computational Linguistics*, pages 33–41, Morristown, NJ, USA. Association for Computational Linguistics.

- Eneko Agirre and Mark Stevenson. 2006. Knowledge sources for wsd. In *Word Sense Disambiguation: Algorithms and Applications*, volume 33 of *Text, Speech and Language Technology*, pages 217–252. Springer, Dordrecht, The Netherlands.
- Scott Cotton, Phil Edmonds, Adam Kilgarriff, and Martha Palmer. 2001. Senseval-2. <http://www.sle.sharp.co.uk/senseval2>.
- Christiane Fellbaum, editor. 1998. *WordNet An Electronic Lexical Database*. The MIT Press, Cambridge, MA ; London, May.
- William A. Gale, Kenneth W. Church, and David Yarowsky. 1992. One sense per discourse. In *HLT '91: Proceedings of the workshop on Speech and Natural Language*, pages 233–237, Morristown, NJ, USA. Association for Computational Linguistics.
- Graeme Hirst. 1987. *Semantic interpretation and the resolution of ambiguity*. Cambridge University Press, New York, NY, USA.
- Thomas Léauté, Brammert Ottens, and Radoslaw Szymanek. 2009. FRODO 2.0: An open-source framework for distributed constraint optimization. In *Proceedings of the IJCAI'09 Distributed Constraint Reasoning Workshop (DCR'09)*, pages 160–164, Pasadena, California, USA, July 13. <http://liawww.epfl.ch/frodo/>.
- Yoong Keok Lee and Hwee Tou Ng. 2002. An empirical evaluation of knowledge sources and learning algorithms for word sense disambiguation. In *EMNLP '02: Proceedings of the ACL-02 conference on Empirical methods in natural language processing*, pages 41–48, Morristown, NJ, USA. Association for Computational Linguistics.
- Roger Mailler and Victor Lesser. 2004. Solving distributed constraint optimization problems using cooperative mediation. In *AAMAS '04: Proceedings of the Third International Joint Conference on Autonomous Agents and Multiagent Systems*, pages 438–445, Washington, DC, USA. IEEE Computer Society.
- David Martínez, Eneko Agirre, and Lluís Màrquez. 2002. Syntactic features for high precision word sense disambiguation. In *COLING*.
- Susan W. McRoy. 1992. Using multiple knowledge sources for word sense discrimination. *COMPUTATIONAL LINGUISTICS*, 18:1–30.
- Rada Mihalcea and Phil Edmonds, editors. 2004. *Proceedings Senseval-3 3rd International Workshop on Evaluating Word Sense Disambiguation Systems*. ACL, Barcelona, Spain.
- Rada Mihalcea and Dan I. Moldovan. 2001. extended wordnet: progress report. In *in Proceedings of NAACL Workshop on WordNet and Other Lexical Resources*, pages 95–100.
- Pragnesh Jay Modi, Wei-Min Shen, Milind Tambe, and Makoto Yokoo. 2004. Adopt: Asynchronous distributed constraint optimization with quality guarantees. *Artificial Intelligence*, 161:149–180.
- Pragnesh Jay Modi. 2003. Distributed constraint optimization for multiagent systems. *PhD Thesis*.
- Adrian Petcu and Boi Faltings. 2005. A scalable method for multiagent constraint optimization. In *IJCAI'05: Proceedings of the 19th international joint conference on Artificial intelligence*, pages 266–271, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.
- Ravi Sinha and Rada Mihalcea. 2007. Unsupervised graph-based word sense disambiguation using measures of word semantic similarity. In *ICSC '07: Proceedings of the International Conference on Semantic Computing*, pages 363–369, Washington, DC, USA. IEEE Computer Society.
- Mark Stevenson and Yorick Wilks. 2001. The interaction of knowledge sources in word sense disambiguation. *Comput. Linguist.*, 27(3):321–349.
- David Yarowsky and Radu Florian. 2002. Evaluating sense disambiguation across diverse parameter spaces. *Natural Language Engineering*, 8:2002.

A probabilistic generative model for an intermediate constituency-dependency representation

Federico Sangati

Institute for Logic, Language and Computation

University of Amsterdam, the Netherlands

f.sangati@uva.nl

Abstract

We present a probabilistic model extension to the Tesnière Dependency Structure (TDS) framework formulated in (Sangati and Mazza, 2009). This representation incorporates aspects from both constituency and dependency theory. In addition, it makes use of *junction* structures to handle coordination constructions. We test our model on parsing the English Penn WSJ treebank using a re-ranking framework. This technique allows us to efficiently test our model without needing a specialized parser, and to use the standard evaluation metric on the original Phrase Structure version of the treebank. We obtain encouraging results: we achieve a small improvement over state-of-the-art results when re-ranking a small number of candidate structures, on all the evaluation metrics except for chunking.

1 Introduction

Since its origin, computational linguistics has been dominated by Constituency/Phrase Structure (PS) representation of sentence structure. However, recently, we observe a steady increase in popularity of Dependency Structure (DS) formalisms. Several researchers have compared the two alternatives, in terms of linguistic adequacy (Nivre, 2005; Schneider, 2008), practical applications (Ding and Palmer, 2005), and evaluations (Lin, 1995).

Dependency theory is historically accredited to Lucien Tesnière (1959), although the relation of dependency between words was only one of the various key elements proposed to represent sentence structures. In fact, the original formulation incorporates the notion of *chunk*, as well as a special type of structure to represent *coordination*.

The Tesnière Dependency Structure (TDS) representation we propose in (Sangati and Mazza, 2009), is an attempt to formalize the original work of Tesnière, with the intention to develop a simple but consistent representation which combines constituencies and dependencies. As part of this work, we have implemented an automatic conversion¹ of the English Penn Wall Street Journal (WSJ) treebank into the new annotation scheme.

In the current work, after introducing the key elements of TDS (section 2), we describe a first probabilistic extension to this framework, which aims at modeling the different levels of the representation (section 3). We test our model on parsing the WSJ treebank using a re-ranking framework. This technique allows us to efficiently test our system without needing a specialized parser, and to use the standard evaluation metric on the original PS version of the treebank. In section 3.4 we also introduce new evaluation schemes on specific aspects of the new TDS representation which we will include in the results presented in section 3.4.

2 TDS representation

It is beyond the scope of this paper to provide an exhaustive description of the TDS representation of the WSJ. It is nevertheless important to give the reader a brief summary of its key elements, and compare it with some of the other representations of the WSJ which have been proposed. Figure 1 shows the original PS of a WSJ tree (a), together with 3 other representations: (b) TDS, (c) DS², and (d) CCG (Hockenmaier and Steedman, 2007).

¹staff.science.uva.nl/~fsangati/TDS

²The DS representation is taken from the conversion procedure used in the CoNLL 2007 Shared Task on dependency parsing (Nivre et al., 2007). Although more elaborate representations have been proposed (de Marneffe and Manning, 2008; Cinková et al., 2009) we have chosen this DS representation because it is one of the most commonly used within the CL community, given that it relies on a fully automatic conversion procedure.

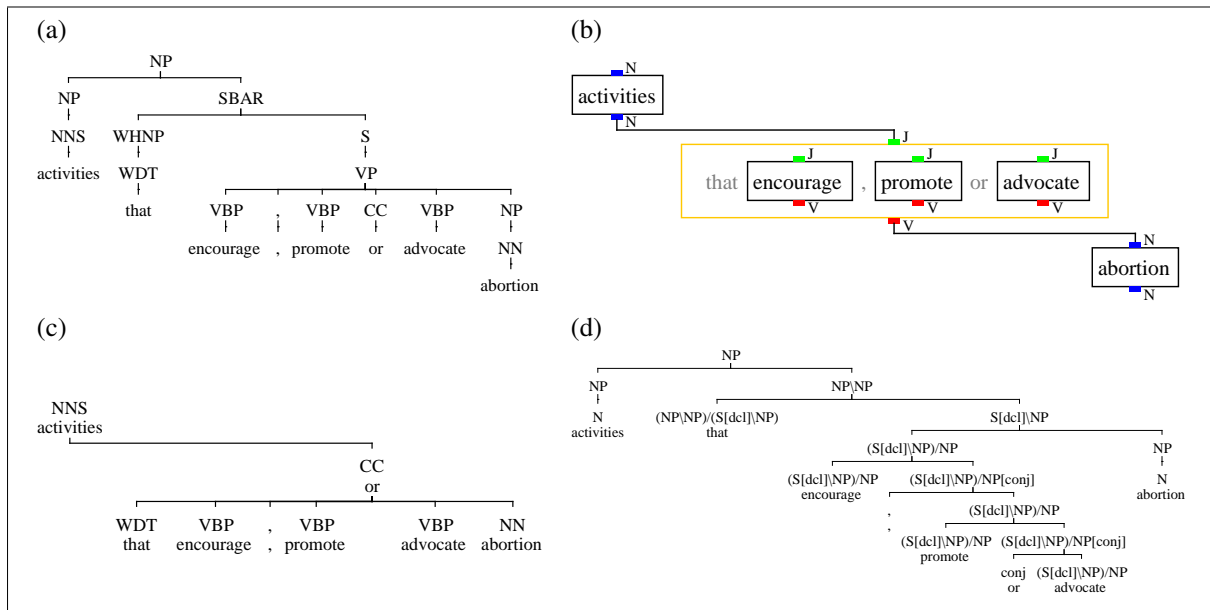


Figure 1: Four different structure representations, derived from a sentence of the WSJ treebank (section 00, #977). (a) PS (original), (b) CCG, (c) DS, (d) TDS.

Words and Blocks In TDS, words are divided in *functional* words (determiners, prepositions, etc.) and *content* words (verbs, nouns, etc.). *Blocks* are the basic elements (chunks) of a structure, which can be combined either via the dependency relation or the junction operation. Blocks can be of two types: *standard* and *junction* blocks. Both types may contain any sequence of functional words. Standard blocks (depicted as black boxes) represent the elementary chunks of the original PS, and include exactly one content word.

Coordination Junction blocks (depicted as yellow boxes) are used to represent coordinated structures. They contain two or more blocks (conjuncts) possibly coordinated by means of functional words (conjunctions). In Figure 1(d) the yellow junction block contains three separate standard blocks. This representation allows to capture the fact that these conjuncts occupy the same role: they all share the relativizer ‘that’, they all depend on the noun ‘activities’, and they all govern the noun ‘abortion’. In Figure 1(a,c), we can notice that both PS and DS do not adequately represent coordination structures: the PS annotation is rather flat, avoiding to group the three verbs in a unique unit, while in the DS the last noun ‘abortion’ is at the same level of the verbs it should be a dependent of. On the other hand, the CCG structure of Figure 1(d), properly represents the coordination. It does so by grouping the first three verbs in a unique constituent which is in turn bi-

narized in a right-branching structure. One of the strongest advantages of the CCG formalism, is that every structure can be automatically mapped to a logical-form representation. This is one reason why it needs to handle coordinations properly. Nevertheless, we conjecture that this representation of coordination might introduce some difficulties for parsing: it is very hard to capture the relation between ‘advocate’ and ‘abortion’ since they are several levels away in the structure.

Categories and Transference There are 4 different block categories, which are indicated with little colored bricks (as well as one-letter abbreviation) on top and at the bottom of the corresponding blocks: verbs (red, V), nouns (blue, N), adverbs (yellow, A), and adjectives (green, J). Every block displays at the bottom the *original category* determined by the content word (or the original category of the conjuncts if it is a junction structure), and at the top, the *derived category* which relates to the grammatical role of the whole block in relation to the governing block. In several cases we can observe a shift in the categories of a block, from the original to the derived category. This phenomenon is called *transference* and often occurs by means of functional words in the block. In Figure 1(b) we can observe the transference of the junction block, which has the original category of a verb, but takes the role of an adjective (through the relativizer ‘that’) in modifying the noun ‘activities’.

$$P(S) = P_{BGM}(S) \cdot P_{BEM}(S) \cdot P_{WFM}(S) \quad (1)$$

$$P_{BGM}(S) = \prod_{B \in dependentBlocks(S)} P(B|parent(B), direction(B), leftSibling(B)) \quad (2)$$

$$P_{BEM}(S) = \prod_{B \in blocks(S)} P(elements(B)|derivedCat(B)) \quad (3)$$

$$P_{WFM}(S) = \prod_{B \in standardBlocks(S)} P(cw(B)|cw(parent(B)), cats(B), fw(B), context(B)) \quad (4)$$

Table 1: Equation (1) gives the likelihood of a structure S as the product of the likelihoods of generating three aspects of the structure, according to the three models (BGM, BEM, WFM) specified in equations (2-4) and explained in the main text.

3 A probabilistic Model for TDS

This section describes the probabilistic generative model which was implemented in order to disambiguate TDS structures. We have chosen the same strategy we have described in (Sangati et al., 2009). The idea consists of utilizing a state of the art parser to compute a list of k -best candidates of a test sentence, and evaluate the new model by using it as a reranker. How well does it select the most probable structure among the given candidates? Since no parser currently exists for the TDS representation, we utilize a state of the art parser for PS trees (Charniak, 1999), and transform each candidate to TDS. This strategy can be considered a first step to efficiently test and compare different models before implementing a full-fledged parser.

3.1 Model description

In order to compute the probability of a given TDS structure, we make use of three separate probabilistic generative models, each responsible for a specific aspect of the structure being generated. The probability of a TDS structure is obtained by multiplying its probabilities in the three models, as reported in the first equation of Table 2.

The first model (equation 2) is the **Block Generation Model (BGM)**. It describes the event of generating a block B as a dependent of its parent block (governor). The dependent block B is identified with its categories (both original and derived), and its functional words, while the parent block is characterized by the original category only. Moreover, in the conditioning context we specify the direction of the dependent with respect to the par-

ent³, and its adjacent left sister (*null* if not present) specified with the same level of details of B . The model applies only to dependent blocks⁴.

The second model (equation 3) is the **Block Expansion Model (BEM)**. It computes the probability of a generic block B of known derived category, to expand to the list of elements it is composed of. The list includes the category of the content word, in case the expansion leads to a standard block. In case of a junction structure, it contains the conjunctions and the conjunct blocks (each identified with its categories and its functional words) in the order they appear. Moreover, all functional words in the block are added to the list⁵. The model applies to all blocks.

The third model (equation 4) is the **Word Filling Model (WFM)**, which applies to each standard block B of the structure. It models the event of filling B with a content word (cw), given the content word of the governing block, the categories ($cats$) and functional words (fw) of B , and further information about the context⁶ in which B occurs. This model becomes particularly interest-

³A dependent block can have three different positions with respect to the parent block: left, right, inner. The first two are self-explanatory. The *inner* case occurs when the dependent block starts after the beginning of the parent block but ends before it (e.g. *a nice dog*).

⁴A block is a dependent block if it is not a conjunct. In other words, it must be connected with a line to its governor.

⁵The attentive reader might notice that the functional words are generated twice (in BGM and BEM). This decision, although not fully justified from a statistical viewpoint, seems to drive the model towards a better disambiguation.

⁶ $context(B)$ comprises information about the grandparent block (original category), the adjacent left sibling block (derived category), the direction of the content word with respect to its governor (in this case only left and right), and the absolute distance between the two words.

ing when a standard block is a dependent of a junction block (such as ‘abortion’ in Figure 1(d)). In this case, the model needs to capture the dependency relation between the content word of the dependent block and each of the content words belonging to the junction block⁷.

3.2 Smoothing

In all the three models we have adopted a smoothing techniques based on back-off level estimation as proposed by Collins (1999). The different back-off estimates, which are listed in decreasing levels of details, are interpolated with confidence weights⁸ derived from the training corpus.

The first two models are implemented with two levels of back-off, in which the last is a constant value (10^{-6}) to make the overall probability small but not zero, for unknown events.

The third model is implemented with three levels of back-off: the last is set to the same constant value (10^{-6}), the first encodes the dependency event using both pos-tags and lexical information of the governor and the dependent word, while the second specifies only pos-tags.

3.3 Experiment Setup

We have tested our model on the WSJ section of Penn Treebank (Marcus et al., 1993), using sections 02-21 as training and section 22 for testing. We employ the Max-Ent parser, implemented by Charniak (1999), to generate a list of k -best PS candidates for the test sentences, which are then converted into TDS representation.

Instead of using Charniak’s parser in its original settings, we train it on a version of the corpus in which we add a special suffix to constituents which have circumstantial role⁹. This decision is based on the observation that the TDS formalism well captures the argument structure of verbs, and

⁷In order to derive the probability of this multi-event we compute the average between the probabilities of the single events which compose it.

⁸Each back-off level obtains a confidence weight which decreases with the increase of the *diversity of the context* ($\theta(C_i)$), which is the number of separate events occurring with the same context (C_i). More formally if $f(C_i)$ is the frequency of the conditioning context of the current event, the weight is obtained as $f(C_i)/(f(C_i) \cdot \mu \cdot \theta(C_i))$; see also (Bikel, 2004). In our model we have chosen μ to be 5 for the first model, and 50 for the second and the third.

⁹Those which have certain function tags (e.g. ADV, LOC, TMP). The full list is reported in (Sangati and Mazza, 2009). It was surprising to notice that the performance of this slightly modified parser (in terms of F-score) is only slightly lower than how it performs out-of-the-box (0.13%).

we believe that this additional information might benefit our model.

We then applied our probabilistic model to re-rank the list of available k -best TDS, and evaluate the selected candidates using several metrics which will be introduced next.

3.4 Evaluation Metrics for TDS

The re-ranking framework described above, allows us to keep track of the original PS of each TDS candidate. This provides an implicit advantage for evaluating our system, viz. it allows us to evaluate the re-ranked structures both in terms of the standard evaluation benchmark on the original PS (F-score) as well as on more refined metrics derived from the converted TDS representation. In addition, the specific head assignment that the TDS conversion procedure performs on the original PS, allows us to convert every PS candidate to a standard projective DS, and from this representation we can in turn compute the standard benchmark evaluation for DS, i.e. unlabeled attachment score¹⁰ (UAS) (Lin, 1995; Nivre et al., 2007).

Concerning the TDS representation, we have formulated 3 evaluation metrics which reflect the accuracy of the chosen structure with respect to the gold structure (the one derived from the manually annotated PS), regarding the different components of the representation:

Block Detection Score (BDS): the accuracy of detecting the correct boundaries of the blocks in the structure¹¹.

Block Attachment Score (BAS): the accuracy of detecting the correct governing block of each block in the structure¹².

Junction Detection Score (JDS): the accuracy of detecting the correct list of content-words composing each junction block in the structure¹³.

¹⁰UAS measures the percentage of words (excluding punctuation) having the correct governing word.

¹¹It is calculated as the harmonic mean between recall and precision between the test and gold set of blocks, where each block is identified with two numerical values representing the start and the end position (punctuation words are discarded).

¹²It is computed as the percentage of words (both functional and content words, excluding punctuation) having the correct governing block. The governing block of a word, is defined as the governor of the block it belongs to. If the block is a conjunct, its governing block is computed recursively as the governing block of the junction block it belongs to.

¹³It is calculated as the harmonic mean between recall and precision between the test and gold set of junction blocks expansions, where each expansion is identified with the list of content words belonging to the junction block. A recursive junction structure expands to a list of lists of content-words.

	F-Score	UAS	BDS	BAS	JDS
Charniak ($k = 1$)	89.41	92.24	94.82	89.29	75.82
Oracle Best F-Score ($k = 1000$)	97.47	96.98	97.03	95.79	82.26
Oracle Worst F-Score ($k = 1000$)	57.04	77.04	84.71	70.10	43.01
Oracle Best JDS ($k = 1000$)	90.54	93.77	96.20	90.57	93.55
PCFG-reranker ($k = 5$)	89.03	92.12	94.86	88.94	75.88
PCFG-reranker ($k = 1000$)	83.52	87.04	92.07	82.32	69.17
TDS-reranker ($k = 5$)	89.65	92.33	94.77	89.35	76.23
TDS-reranker ($k = 10$)	89.10	92.11	94.58	88.94	75.47
TDS-reranker ($k = 100$)	86.64	90.24	93.11	86.34	69.60
TDS-reranker ($k = 500$)	84.94	88.62	91.97	84.43	65.30
TDS-reranker ($k = 1000$)	84.31	87.89	91.42	83.69	63.65

Table 2: Results of Charniak’s parser, the TDS-reranker, and the PCFG-reranker according to several evaluation metrics, when the number k of best-candidates increases.

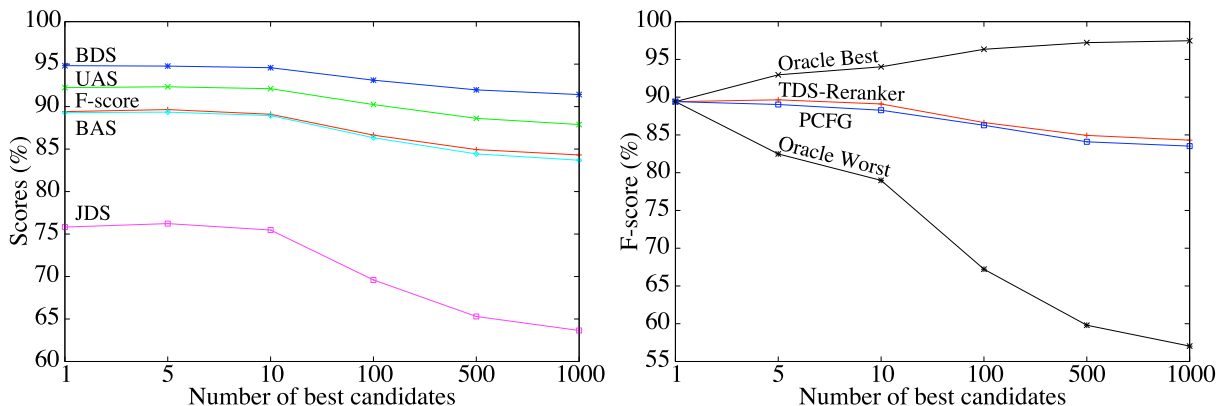


Figure 2: **Left:** results of the TDS-reranking model according to several evaluation metrics as in Table 2. **Right:** comparison between the F-scores of the TDS-reranker and a vanilla PCFG-reranker (together with the lower and the upper bound), with the increase of the number of best candidates.

3.5 Results

Table 2 reports the results we obtain when re-ranking with our model an increasing number of k -best candidates provided by Charniak’s parser (the same results are shown in the left graph of Figure 2). We also report the results relative to a PCFG-reranker obtained by computing the probability of the k -best candidates using a standard vanilla-PCFG model derived from the same training corpus. Moreover, we evaluate, by means of an oracle, the upper and lower bound of the F-Score and JDS metric, by selecting the structures which maximizes/minimizes the results.

Our re-ranking model performs rather well for a limited number of candidate structures, and outperforms Charniak’s model when $k = 5$. In this case we observe a small boost in performance for the detection of junction structures, as well as for

all other evaluation metrics, except for the BDS.

The right graph in Figure 2 compares the F-score performance of the TDS-reranker against the PCFG-reranker. Our system consistently outperforms the PCFG model on this metric, as for UAS, and BAS. Concerning the other metrics, as the number of k -best candidates increases, the PCFG model outperforms the TDS-reranker both according to the BDS and the JDS.

Unfortunately, the performance of the re-ranking model worsens progressively with the increase of k . We find that this is primarily due to the lack of robustness of the model in detecting the block boundaries. This suggests that the system might benefit from a separate preprocessing step which could chunk the input sentence with higher accuracy (Sang et al., 2000). In addition the same module could detect local (intra-clausal) coordinations, as illustrated by (Marinčič et al., 2009).

4 Conclusions

In this paper, we have presented a probabilistic generative model for parsing TDS syntactic representation of English sentences. We have given evidence for the usefulness of this formalism: we consider it a valid alternative to commonly used PS and DS representations, since it incorporates the most relevant features of both notations; in addition, it makes use of junction structures to represent coordination, a linguistic phenomena highly abundant in natural language production, but often neglected when it comes to evaluating parsing resources. We have therefore proposed a special evaluation metrics for junction detection, with the hope that other researchers might benefit from it in the future. Remarkably, Charniak’s parser performs extremely well in all the evaluation metrics besides the one related to coordination.

Our parsing results are encouraging: the overall system, although only when the candidates are highly reliable, can improve on Charniak’s parser on all the evaluation metrics with the exception of chunking score (BDS). The weakness on performing chunking is the major factor responsible for the lack of robustness of our system. We are considering to use a dedicated pre-processing module to perform this step with higher accuracy.

Acknowledgments The author gratefully acknowledge funding by the Netherlands Organization for Scientific Research (NWO): this work is funded through a Vici-grant “Integrating Cognition” (277.70.006) to Rens Bod. We also thank 3 anonymous reviewers for very useful comments.

References

- Daniel M. Bikel. 2004. Intricacies of Collins’ Parsing Model. *Comput. Linguist.*, 30(4):479–511.
- Eugene Charniak. 1999. A Maximum-Entropy-Inspired Parser. Technical report, Providence, RI, USA.
- Silvie Cinková, Josef Toman, Jan Hajič, Kristýna Čermáková, Václav Klimeš, Lucie Mladová, Jana Šindlerová, Kristýna Tomšů, and Zdeněk Žabokrtský. 2009. Tectogrammatical Annotation of the Wall Street Journal. *The Prague Bulletin of Mathematical Linguistics*, (92).
- Michael J. Collins. 1999. *Head-Driven Statistical Models for Natural Language Parsing*. Ph.D. thesis, University of Pennsylvania.
- Marie-Catherine de Marneffe and Christopher D. Manning. 2008. The Stanford Typed Dependencies Representation. In *Coling 2008: Proceedings of the workshop on Cross-Framework and Cross-Domain Parser Evaluation*, pages 1–8, Manchester, UK.
- Yuan Ding and Martha Palmer. 2005. Machine Translation Using Probabilistic Synchronous Dependency Insertion Grammars. In *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL’05)*, pages 541–548.
- Julia Hockenmaier and Mark Steedman. 2007. CCG-bank: A Corpus of CCG Derivations and Dependency Structures Extracted from the Penn Treebank. *Computational Linguistics*, 33(3):355–396.
- Dekang Lin. 1995. A Dependency-based Method for Evaluating Broad-Coverage Parsers. In *In Proceedings of IJCAI-95*, pages 1420–1425.
- Mitchell P. Marcus, Mary Ann Marcinkiewicz, and Beatrice Santorini. 1993. Building a Large Annotated Corpus of English: The Penn Treebank. *Computational Linguistics*, 19(2):313–330.
- Domen Marinčič, Matjaž Gams, and Tomaž Šef. 2009. Intraclausal Coordination and Clause Detection as a Preprocessing Step to Dependency Parsing. In *TSD’09: Proceedings of the 12th International Conference on Text, Speech and Dialogue*, pages 147–153, Berlin, Heidelberg. Springer-Verlag.
- Joakim Nivre, Johan Hall, Sandra Kübler, Ryan McDonald, Jens Nilsson, Sebastian Riedel, and Deniz Yuret. 2007. The CoNLL 2007 Shared Task on Dependency Parsing. In *Proceedings of the CoNLL Shared Task Session of EMNLP-CoNLL 2007*, pages 915–932, Prague, Czech Republic.
- Joakim Nivre. 2005. Dependency Grammar and Dependency Parsing. Technical report, Växjö University: School of Mathematics and Systems Engineering.
- Erik F. Tjong Kim Sang, Sabine Buchholz, and Kim Sang. 2000. Introduction to the CoNLL-2000 Shared Task: Chunking. In *Proceedings of CoNLL-2000 and LLL-2000, Lisbon, Portugal*.
- Federico Sangati and Chiara Mazza. 2009. An English Dependency Treebank à la Tesnière. In *The 8th International Workshop on Treebanks and Linguistic Theories*, pages 173–184, Milan, Italy.
- Federico Sangati, Willem Zuidema, and Rens Bod. 2009. A generative re-ranking model for dependency parsing. In *Proceedings of the 11th International Conference on Parsing Technologies (IWPT’09)*, pages 238–241, Paris, France, October.
- Gerold Schneider. 2008. *Hybrid long-distance functional dependency parsing*. Ph.D. thesis, University of Zurich.
- Lucien Tesnière. 1959. *Éléments de syntaxe structurale*. Editions Klincksieck, Paris.

Sentiment Translation through Lexicon Induction

Christian Scheible

Institute for Natural Language Processing
University of Stuttgart

scheibcn@ims.uni-stuttgart.de

Abstract

The translation of sentiment information is a task from which sentiment analysis systems can benefit. We present a novel, graph-based approach using Sim-Rank, a well-established vertex similarity algorithm to transfer sentiment information between a source language and a target language graph. We evaluate this method in comparison with SO-PMI.

1 Introduction

Sentiment analysis is an important topic in computational linguistics that is of theoretical interest but also implies many real-world applications. Usually, two aspects are of importance in sentiment analysis. The first is the detection of subjectivity, i.e. whether a text or an expression is meant to express sentiment at all; the second is the determination of sentiment orientation, i.e. what sentiment is to be expressed in a structure that is considered subjective.

Work on sentiment analysis most often covers resources or analysis methods in a single language, usually English. However, the transfer of sentiment analysis between languages can be advantageous by making use of resources for a source language to improve the analysis of the target language.

This paper presents an approach to the transfer of sentiment information between languages. It is built around an algorithm that has been successfully applied for the acquisition of bilingual lexicons. One of the main benefits of the method is its ability of handling sparse data well.

Our experiments are carried out using English as a source language and German as a target language.

2 Related Work

The translation of sentiment information has been the topic of multiple publications.

Mihalcea et al. (2007) propose two methods for translating sentiment lexicons. The first method simply uses bilingual dictionaries to translate an English sentiment lexicon. A sentence-based classifier built with this list achieved high precision but low recall on a small Romanian test set. The second method is based on parallel corpora. The source language in the corpus is annotated with sentiment information, and the information is then projected to the target language. Problems arise due to mistranslations, e.g., because irony is not recognized.

Banea et al. (2008) use machine translation for multilingual sentiment analysis. Given a corpus annotated with sentiment information in one language, machine translation is used to produce an annotated corpus in the target language, by preserving the annotations. The original annotations can be produced either manually or automatically.

Wan (2009) constructs a multilingual classifier using co-training. In co-training, one classifier produces additional training data for a second classifier. In this case, an English classifier assists in training a Chinese classifier.

The induction of a sentiment lexicon is the subject of early work by (Hatzivassiloglou and McKeown, 1997). They construct graphs from coordination data from large corpora based on the intuition that adjectives with the same sentiment orientation are likely to be coordinated. For example, *fresh and delicious* is more likely than *rotten and delicious*. They then apply a graph clustering algorithm to find groups of adjectives with the same orientation. Finally, they assign the same label to all adjectives that belong to the same cluster. The authors note that some words cannot be assigned a unique label since their sentiment depends on con-

text.

Turney (2002) suggests a corpus-based extraction method based on his pointwise mutual information (PMI) synonymy measure. He assumes that the sentiment orientation of a phrase can be determined by comparing its pointwise mutual information with a positive (*excellent*) and a negative phrase (*poor*). An introduction to SO-PMI is given in Section 5.1

3 Bilingual Lexicon Induction

Typical approaches to the induction of bilingual lexicons involve gathering new information from a small set of known identities between the languages which is called a *seed lexicon* and incorporating intralingual sources of information (e.g. cooccurrence counts). Two examples of such methods are a graph-based approach by Dorow et al. (2009) and a vector-space based approach by Rapp (1999). In this paper, we will employ the graph-based method.

SimRank was first introduced by Jeh and Widom (2002). It is an iterative algorithm that measures the similarity between all vertices in a graph. In SimRank, two nodes are similar if their neighbors are similar. This defines a recursive process that ends when the two nodes compared are identical. As proposed by Dorow et al. (2009), we will apply it to a graph \mathcal{G} in which vertices represent words and edges represent relations between words. SimRank will then yield similarity values between vertices that indicate the degree of relatedness between them with regard to the property encoded through the edges. For two nodes i and j in \mathcal{G} , similarity according to SimRank is defined as

$$\text{sim}(i, j) = \frac{c}{|N(i)||N(j)|} \sum_{k \in N(i), l \in N(j)} \text{sim}(k, l),$$

where $N(x)$ is the neighborhood of x and c is a weight factor that determines the influence of neighbors that are farther away. The initial condition for the recursion is $\text{sim}(i, i) = 1$.

Dorow et al. (2009) further propose the application of the SimRank algorithm for the calculation of similarities between a source graph \mathcal{S} and a target graph \mathcal{T} . Initially, some relations between the two graphs need to be known. When operating on word graphs, these can be taken from a bilingual lexicon. This provides us with a framework for the induction of a bilingual lexicon which can be

constructed based on the obtained similarity values between the vertices of the two graphs.

One problem of SimRank observed in experiments by Laws et al. (2010) was that while words with high similarity were semantically related, they often were not exact translations of each other but instead often fell into the categories of hyponymy, hypernymy, holonymy, or meronymy. However, this makes the similarity values applicable for the translation of sentiment since it is a property that does not depend on exact synonymy.

4 Sentiment Transfer

Although unsupervised methods for the design of sentiment analysis systems exist, any approach can benefit from using resources that have been established in other languages. The main problem that we aim to deal with in this paper is the transfer of such information between languages. The SimRank lexicon induction method is suitable for this purpose since it can produce useful similarity values even with a small seed lexicon.

First, we build a graph for each language. The vertices of these graphs will represent adjectives while the edges are coordination relations between these adjectives. An example for such a graph is given in Figure 1.

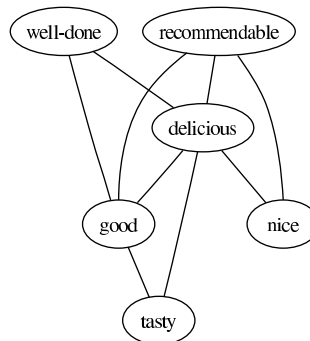


Figure 1: Sample graph showing English coordination relations.

The use of coordination information has been shown to be beneficial for example in early work by Hatzivassiloglou and McKeown (1997).

Seed links between those graphs will be taken from a universal dictionary. Figure 2 shows an example graph. Here, intralingual coordination relations are represented as black lines, seed relations as solid grey lines, and relations that are induced through SimRank as dashed grey lines.

After computing similarities in this graph, we

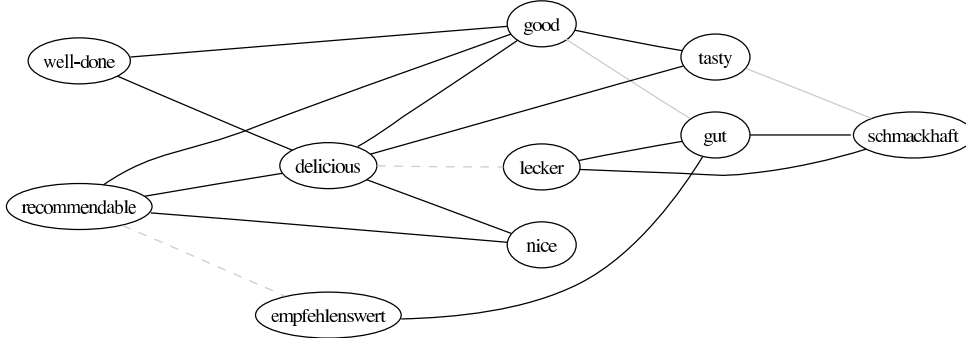


Figure 2: Sample graph showing English and German coordination relations. Solid black lines represent coordinations, solid grey lines represent seed relations, and dashed grey lines show induced relations.

need to obtain sentiment values. We will define the sentiment score (sent) as

$$\text{sent}(n_t) = \sum_{n_s \in \mathcal{S}} \text{sim}_{\text{norm}}(n_s, n_t) \text{sent}(n_s),$$

where n_t is a node in the target graph \mathcal{T} , and \mathcal{S} the source graph. This way, the sentiment score of each node is an average over all nodes in \mathcal{S} weighted by their normalized similarity, sim_{norm} .

We define the normalized similarity as

$$\text{sim}_{\text{norm}}(n_s, n_t) = \frac{\text{sim}(n_s, n_t)}{\sum_{n_s \in \mathcal{S}} \text{sim}(n_s, n_t)}.$$

Normalization guarantees that all sentiment scores lie within a specified range. Scores are not a direct indicator for orientation since the similarities still include a lot of noise. Therefore, we interpret the scores by assigning each word to a category by finding score thresholds between the categories.

5 Experiments

5.1 Baseline Method (SO-PMI)

We will compare our method to the well-established SO-PMI algorithm by Turney (2002) to show an improvement over an unsupervised method. The algorithm works with cooccurrence counts on large corpora. To determine the semantic orientation of a word w , the hits near positive (P words) and negative (N words) seed words is used. The SO-PMI equation is given as

$$\begin{aligned} \text{SO-PMI}(\text{word}) = & \log_2 \left(\frac{\prod_{pword \in Pwords} \text{hits}(\text{word NEAR } pword)}{\prod_{nword \in Nwords} \text{hits}(\text{word NEAR } nword)} \right) \\ & \times \frac{\prod_{nword \in Nwords} \text{hits}(nword)}{\prod_{pword \in Pwords} \text{hits}(pword)} \end{aligned}$$

5.2 Data Acquisition

We used the English and German Wikipedia branches as our corpora. We extracted coordinations from the corpus using a simple CQP pattern search (Christ et al., 1999). For our experiments, we looked only at coordinations with *and*. For the English corpus, we used the pattern `[pos = "JJ"] ([pos = ", "] [pos = "JJ"])* ([pos = ", "]? "and" [pos = "JJ"])+`, and for the German corpus, the pattern `[pos = "ADJ.*"] ([pos = ", "] [pos = "ADJ.*"])* ("und" [pos = "ADJ"])+` was used. This yielded 477,291 pairs of coordinated English adjectives and 44,245 German pairs. We used the dict.cc dictionary¹ as a seed dictionary. It contained a total of 30,551 adjectives.

After building a graph out of this data as described in Section 4, we apply the SimRank algorithm using 7 iterations.

Data for the SO-PMI method had to be collected from queries to search engines since the information available in the Wikipedia corpus was too sparse. Since Google does not provide a stable NEAR operator, we used coordinations instead. For each of the test words w and the SO-PMI seed words s we made two queries `"w und s"` and `"s und w"` to Google. The quotes and + were added to ensure that no spelling correction or synonym replacements took place. Since the original experiments were designed for an English corpus, a set of German seed words had to be constructed. We chose *gut*, *nett*, *richtig*, *schön*, *ordentlich*, *angenehm*, *aufrechtig*, *gewissenhaft*, and *hervorragend* as positive seeds, and *schlecht*, *teuer*, *falsch*, *böse*, *feindlich*, *verhasst*, *widerlich*, *fehlerhaft*, and

¹<http://www.dict.cc/>

word	value
strongpos	1.0
weakpos	0.5
neutral	0.0
weakneg	-0.5
strongneg	-1.0

Table 1: Assigned values for positivity labels

mangelhaft as negative seeds.

We constructed a test set by randomly selecting 200 German adjectives that occurred in a coordination in Wikipedia. We then eliminated adjectives that we deemed uncommon or too difficult to understand or that were mislabeled as adjectives. This resulted in a 150 word test set. To determine the sentiment of these adjectives, we asked 9 human judges, all native German speakers, to annotate them given the classes *neutral*, *slightly negative*, *very negative*, *slightly positive*, and *very positive*, reflecting the categories from the training data. In the annotation process, another 7 adjectives had to be discarded because one or more annotators marked them as unknown.

Since human judges tend to interpret scales differently, we examine their agreement using Kendall’s coefficient of concordance (W) including correction for ties (Legendre, 2005) which takes ranks into account. The agreement was calculated as $W = 0.674$ with a significant confidence ($p < .001$), which is usually interpreted as substantial agreement. Manual examination of the data showed that most disagreement between the annotators occurred with adjectives that are tied to political implications, for example *nuklear* (*nuclear*).

5.3 Sentiment Lexicon Induction

For our experiments, we used the polarity lexicon of Wilson et al. (2005). It includes annotations of positivity in the form of the categories *neutral*, weakly positive (*weakpos*), strongly positive (*strongpos*), weakly negative (*weakneg*), and strongly negative (*strongneg*). In order to conduct arithmetic operations on these annotations, mapped them to values from the interval $[-1, 1]$ by using the assignments given in Table 1.

5.4 Results

To compare the two methods to the human raters, we first reproduce the evaluation by Turney (2002)

and examine the correlation coefficients. Both methods will be compared to an average over the human rater values. These values are calculated on values asserted based on Table 1. The correlation coefficients between the automatic systems and the human ratings, SO-PMI yields $r = 0.551$, and SimRank yields $r = 0.587$ which are not significantly different. This shows that SO and SR have about the same performance on this broad measure.

Since many adjectives do not express sentiment at all, the correct categorization of neutral adjectives is as important as the scalar rating. Thus, we divide the adjectives into three categories – positive, neutral, and negative. Due to disagreements between the human judges there exists no clear threshold between these categories. In order to try different thresholds, we assume that sentiment is symmetrically distributed with mean 0 on the human scores. For $x \in \{\frac{i}{20} | 0 \leq i \leq 19\}$, we then assign word w with human rating $score(w)$ to negative if $score(w) \leq -x$, to neutral if $-x < score(w) < x$ and to positive otherwise. This gives us a three-category gold standard for each x that is then the basis for computing evaluation measures. Each category contains a certain percentile of the list of adjectives. By mapping these percentiles to the rank-ordered scores for SO-PMI and SimRank, we can create three-category partitions for them. For example if for $x = 0.35$ 21% of the adjectives are negative, then the 21% of adjectives with the lowest SO-PMI scores are deemed to have been rated negative by SO-PMI.

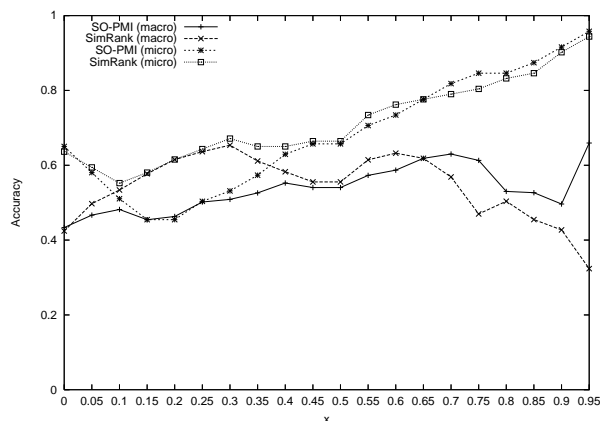


Figure 3: Macro- and micro-averaged Accuracy

First, we will look at the macro- and micro-averaged accuracies for both methods (cf. Figure 3). Overall, SimRank performs better for x

between 0.05 and 0.4 which is a plausible interval for the neutral threshold on the human ratings. The results diverge for very low and high values of x , however these values can be considered unrealistic since they implicate neutral areas that are too small or too large. When comparing the accuracies for each of the classes (cf. Figure 4), we observe that in the aforementioned interval, SimRank has higher accuracy values than SO-PMI for all of them.

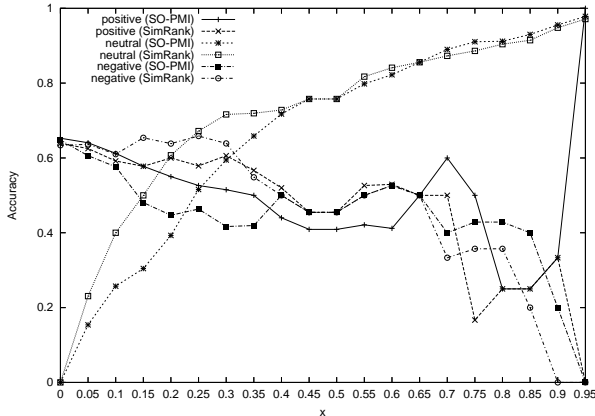


Figure 4: Accuracy for individual classes

Table 2 lists some interesting example words including their human ratings and SO-PMI and SimRank scores which illustrate advantages and possible shortcomings of the two methods. The medians of SO-PMI and SimRank scores are -15.58 and -0.05 , respectively. The mean values are -9.57 for SO-PMI and 0.08 for SimRank, the standard deviations are 13.75 and 0.22 . SimRank values range between -0.67 and 0.41 , SO-PMI ranges between -46.21 and 46.59 . We will assume that the medians mark the center of the set of neutral adjectives.

Ausdrucksvoll receives a positive score from SO-PMI which matches the human rating, however not from SimRank, which assigns a score close to 0 and would likely be considered neutral. This error can be explained by examining the similarity distribution for *ausdrucksvoll* which reveals that there are no nodes that are similar to this node, which was most likely caused by its low degree. *Auferstanden* (resurrected) is perceived as a positive adjective by the human judges, however it is misclassified by SimRank as negative due to its occurrence with words like *gestorben* (deceased) and *gekreuzigt* (crucified) which have negative as-

word (translation)	SR	SO	judges
ausdrucksvoll (expressive)	0.069	22.93	0.39
grafisch (graphic)	-0.050	-4.75	0.00
kriminell (criminal)	-0.389	-15.98	-0.94
auferstanden (resurrected)	-0.338	-10.97	0.34

Table 2: Example adjectives including translation, and their scores

sociations. This suggests that coordinations are sometimes misleading and should not be used as the only data source. *Grafisch* (graphics-related) is an example for a neutral word misclassified by SO-PMI due to its occurrence in positive contexts on the web. Since SimRank is not restricted to relations between an adjective and a seed word, all adjective-adjective coordinations are used for the estimation of a sentiment score. *Kriminell* is also misclassified by SO-PMI for the same reason.

6 Conclusion and Outlook

We presented a novel approach to the translation of sentiment information that outperforms SO-PMI, an established method. In particular, we could show that SimRank outperforms SO-PMI for values of the threshold x in an interval that most likely leads to the correct separation of positive, neutral, and negative adjectives. We intend to compare our system to other available work in the future. In addition to our findings, we created an initial gold standard set of sentiment-annotated German adjectives that will be publicly available.

The two methods are very different in nature; while SO-PMI is suitable for languages in which very large corpora exist, this might not be the case for knowledge-sparse languages. For some German words (e.g. *schwerstkrank* (seriously ill)), SO-PMI lacked sufficient results on the web whereas SimRank correctly assigned negative sentiment. SimRank can leverage knowledge from neighbor words to circumvent this problem. In turn, this information can turn out to be misleading (cf. *auferstanden*). An advantage of our method is that it uses existing resources from another language and can thus be applied without much knowledge about the target language. Our future work will include a further examination of the merits of its application for knowledge-sparse languages.

The underlying graph structure provides a foundation for many conceivable extensions. In this paper, we presented a fairly simple experiment restricted to adjectives only. However, the method

is suitable to include arbitrary parts of speech as well as phrases, as used by Turney (2002). Another conceivable application would be the direct combination of the SimRank-based model with a statistical model.

Currently, our input sentiment list exists only of prior sentiment values, however work by Wilson et al. (2009) has advanced the notion of contextual polarity lists. The automatic translation of this information could be beneficial for sentiment analysis in other languages.

Another important problem in sentiment analysis is the treatment of ambiguity. The sentiment expressed by a word or phrase is context-dependent and is for example related to word sense (Akkaya et al., 2009). Based on regularities in graph structure and similarity, ambiguity resolution might become possible.

References

- C. Akkaya, J. Wiebe, and R. Mihalcea. 2009. Subjectivity Word Sense Disambiguation. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing*, pages 190–199.
- Carmen Banea, Rada Mihalcea, Janyce Wiebe, and Samer Hassan. 2008. Multilingual subjectivity analysis using machine translation. In *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*, pages 127–135, Honolulu, Hawaii, October. Association for Computational Linguistics.
- O. Christ, B.M. Schulze, A. Hofmann, and E. Koenig. 1999. The IMS Corpus Workbench: Corpus Query Processor (CQP): User’s Manual. *University of Stuttgart, March*, 8:1999.
- Beate Dorow, Florian Laws, Lukas Michelbacher, Christian Scheible, and Jason Utt. 2009. A graph-theoretic algorithm for automatic extension of translation lexicons. In *Proceedings of the Workshop on Geometrical Models of Natural Language Semantics*, pages 91–95, Athens, Greece, March. Association for Computational Linguistics.
- Vasileios Hatzivassiloglou and Kathleen R. McKeown. 1997. Predicting the semantic orientation of adjectives. In *Proceedings of the 35th Annual Meeting of the Association for Computational Linguistics*, pages 174–181, Madrid, Spain, July. Association for Computational Linguistics.
- Glen Jeh and Jennifer Widom. 2002. Simrank: a measure of structural-context similarity. In *KDD ’02: Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 538–543, New York, NY, USA. ACM.
- F. Laws, L. Michelbacher, B. Dorow, U. Heid, and H. Schütze. 2010. Building a Cross-lingual Relatedness Thesaurus Using a Graph Similarity Measure. *Submitted on Nov 7, 2009, to the International Conference on Language Resources and Evaluation (LREC)*.
- P. Legendre. 2005. Species associations: the Kendall coefficient of concordance revisited. *Journal of Agricultural Biological and Environment Statistics*, 10(2):226–245.
- Rada Mihalcea, Carmen Banea, and Janyce Wiebe. 2007. Learning multilingual subjective language via cross-lingual projections. In *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, pages 976–983, Prague, Czech Republic, June. Association for Computational Linguistics.
- Reinhard Rapp. 1999. Automatic identification of word translations from unrelated english and german corpora. In *Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics*, pages 519–526, College Park, Maryland, USA, June. Association for Computational Linguistics.
- Peter Turney. 2002. Thumbs up or thumbs down? semantic orientation applied to unsupervised classification of reviews. In *Proceedings of 40th Annual Meeting of the Association for Computational Linguistics*, pages 417–424, Philadelphia, Pennsylvania, USA, July. Association for Computational Linguistics.
- Xiaojun Wan. 2009. Co-training for cross-lingual sentiment classification. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, pages 235–243, Suntec, Singapore, August. Association for Computational Linguistics.
- Theresa Wilson, Janyce Wiebe, and Paul Hoffmann. 2005. Recognizing contextual polarity in phrase-level sentiment analysis. In *Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing*, pages 347–354, Vancouver, British Columbia, Canada, October. Association for Computational Linguistics.
- Theresa Wilson, Janyce Wiebe, and Paul Hoffmann. 2009. Recognizing Contextual Polarity: an Exploration of Features for Phrase-level Sentiment Analysis. *Computational Linguistics*, 35(3):399–433.

Unsupervised Search for The Optimal Segmentation for Statistical Machine Translation

Coşkun Mermer^{1,3} and Ahmet Afşın Akın^{2,3}

¹Boğaziçi University, Bebek, Istanbul, Turkey

²Istanbul Technical University, Sarıyer, Istanbul, Turkey

³TÜBİTAK-UEKAE, Gebze, Kocaeli, Turkey

{coskun, ahmetaa}@uekae.tubitak.gov.tr

Abstract

We tackle the previously unaddressed problem of unsupervised determination of the optimal morphological segmentation for statistical machine translation (SMT) and propose a segmentation metric that takes into account both sides of the SMT training corpus. We formulate the objective function as the posterior probability of the training corpus according to a generative segmentation-translation model. We describe how the IBM Model-1 translation likelihood can be computed incrementally between adjacent segmentation states for efficient computation. Submerging the proposed segmentation method in a SMT task from morphologically-rich Turkish to English does not exhibit the expected improvement in translation BLEU scores and confirms the robustness of phrase-based SMT to translation unit combinatorics. A positive outcome of this work is the described modification to the sequential search algorithm of Morfessor (Creutz and Lagus, 2007) that enables arbitrary-fold parallelization of the computation, which unexpectedly improves the translation performance as measured by BLEU.

1 Introduction

In statistical machine translation (SMT), words are normally considered as the building blocks of translation models. However, especially for morphologically complex languages such as Finnish, Turkish, Czech, Arabic etc., it has been shown that using sub-lexical units obtained after morphological preprocessing can improve the machine translation performance over a word-based system (Habash and Sadat, 2006; Oflazer and Durgar El-Kahlout, 2007; Bisazza and Federico, 2009). However, the effect of segmentation on transla-

tion performance is indirect and difficult to isolate (Lopez and Resnik, 2006).

The challenge in designing a sub-lexical SMT system is the decision of what segmentation to use. Linguistic morphological analysis is intuitive, but it is language-dependent and could be highly ambiguous. Furthermore, it is not necessarily optimal in that (i) manually engineered segmentation schemes can outperform a straightforward linguistic morphological segmentation, e.g., (Habash and Sadat, 2006), and (ii) it may result in even worse performance than a word-based system, e.g., (Durgar El-Kahlout and Oflazer, 2006).

A SMT system designer has to decide what segmentation is optimal for the translation task at hand. Existing solutions to this problem are predominantly heuristic, language-dependent, and as such are not easily portable to other languages. Another point to consider is that the optimal degree of segmentation might decrease as the amount of training data increases (Lee, 2004; Habash and Sadat, 2006). This brings into question: For the particular language pair and training corpus at hand, what is the optimal (level of) sub-word segmentation? Therefore, it is desirable to learn the optimal segmentation in an unsupervised manner.

In this work, we extend the method of Creutz and Lagus (2007) so as to maximize the translation posterior in unsupervised segmentation. The learning process is tailored to the particular SMT task via the same parallel corpus that is used in training the statistical translation models.

2 Related Work

Most works in SMT-oriented segmentation are supervised in that they consist of manual experimentation to choose the best among a set of segmentation schemes, and are language(pair)-dependent. For Arabic, Sadat and Habash (2006) present several morphological preprocessing schemes that entail varying degrees of decomposition and com-

pare the resulting translation performances in an Arabic-to-English task. Shen et al. (2007) use a subset of the morphology and apply only a few simple rules in segmenting words. Durgar El-Kahlout and Oflazer (2006) tackle this problem when translating from English to Turkish, an agglutinative language. They use a morphological analyzer and disambiguation to arrive at morphemes as tokens. However, training the translation models with morphemes actually degrades the translation performance. They outperform the word-based baseline only after some selective morpheme grouping. Bisazza and Federico (2009) adopt an approach similar to the Arabic segmentation studies above, this time in a Turkish-to-English translation setting.

Unsupervised segmentation by itself has garnered considerable attention in the computational linguistics literature (Poon et al., 2009; Snyder and Barzilay, 2008; Dasgupta and Ng, 2007; Creutz and Lagus, 2007; Brent, 1999). However, few works report their performance in a translation task. Virpioja et al. (2007) used Morfessor (Creutz and Lagus, 2007) to segment both sides of the parallel training corpora in translation between Danish, Finnish, and Swedish, but without a consistent improvement in results.

Morfessor, which gives state of the art results in many tests (Kurimo et al., 2009), uses only monolingual information in its objective function. It is conceivable that we can achieve a better segmentation for translation by considering not one but both sides of the parallel corpus. A possible choice is the post-segmentation alignment accuracy. However, Elming et al. (2009) show that optimizing segmentation with respect to alignment error rate (AER) does not improve and even degrades machine translation performance. Snyder and Barzilay (2008) use bilingual information but the segmentation is learned independently from translation modeling.

In Chang et al. (2008), the granularity of the Chinese word segmentation is optimized by training SMT systems for several values of a granularity bias parameter and it is found that the value that maximizes translation performance (as measured by BLEU) is different than the value that maximizes segmentation accuracy (as measured by precision and recall).

One motivation in morphological preprocessing before translation modeling is “morphology

matching” as in Lee (2004) and in the scheme “EN” of Habash and Sadat (2006). In Lee (2004), the goal is to match the lexical granularities of the two languages by starting with a fine-grained segmentation of the Arabic side of the corpus and then merging or deleting Arabic morphemes using alignments with a part-of-speech tagged English corpus. But this method is not completely unsupervised since it requires external linguistic resources in initializing the segmentation with the output of a morphological analyzer and disambiguator. Talbot and Osborne (2006) tackle a special case of morphology matching by identifying redundant distinctions in the morphology of one language compared to another.

3 Method

Maximizing translation performance directly would require SMT training and decoding for each segmentation hypothesis considered, which is computationally infeasible. So we make some conditional independence assumptions using a generative model and decompose the posterior probability $P(M_f|e, f)$. In this notation e and f denote the two sides of a parallel corpus and M_f denotes the segmentation model hypothesized for f . Our approach is an extension of Morfessor (Creutz and Lagus, 2007) so as to include the translation model probability in its cost calculation. Specifically, the segmentation model takes into account the likelihood of both sides of the parallel corpus while searching for the optimal segmentation. The joint likelihood is decomposed into a prior, a monolingual likelihood, and a translation likelihood, as shown in Eq. 1.

$$P(e, f, M_f) = P(M_f)P(f|M_f)P(e|f, M_f) \quad (1)$$

Assuming conditional independence between e and M_f given f , the maximum *a posteriori* (MAP) objective can be written as:

$$\hat{M}_f = \arg \max_{M_f} P(M_f)P(f|M_f)P(e|f) \quad (2)$$

The role of the bilingual component $P(e|f)$ in Eq. 2 can be motivated with a simple example as follows. Consider an occurrence of two phrase pairs in a Turkish-English parallel corpus and the two hypothesized sets of segmentations for the Turkish phrases as in Table 1. Without access to the English side of the corpus, a monolingual segmenter can quite possibly score Seg. #1

	Phrase #1	Phrase #2
Turkish phrase:	anahtar	anahtarım
English phrase:	key	my key
Seg. #1:	anahtar	anahtarı +m
Seg. #2:	anahtar	anahtar +ım

Table 1: Example segmentation hypotheses

higher than Seg. #2 (e.g., due to the high frequency of the observed morph “+m”). On the other hand, a bilingual segmenter is expected to assign a higher alignment probability $P(e|f)$ to Seg. #2 than Seg. #1, because of the aligned words `key|anahtar`, therefore ranking Seg. #2 higher.

The two monolingual components of Eq. 2 are computed as in Creutz and Lagus (2007). To summarize briefly, the prior $P(M_f)$ is assumed to only depend on the frequencies and lengths of the individual morphs, which are also assumed to be independent. The monolingual likelihood $P(f|M_f)$ is computed as the product of morph probabilities estimated from their frequencies in the corpus.

To compute the bilingual (translation) likelihood $P(e|f)$, we use IBM Model 1 (Brown et al., 1993). Let an aligned sentence pair be represented by (s_e, s_f) , which consists of word sequences $s_e = e_1, \dots, e_l$ and $s_f = f_1, \dots, f_m$. Using a purely notational switch of the corpus labels from here on to be consistent with the SMT literature, where the derivations are in the form of $P(f|e)$, the desired translation probability is given by the expression:

$$P(f|e) = \frac{P(m|e)}{(l+1)^m} \prod_{j=1}^m \sum_{i=0}^l t(f_j|e_i), \quad (3)$$

The sentence length probability distribution $P(m|e)$ is assumed to be Poisson with the expected sentence length equal to m .

3.1 Incremental computation of Model-1 likelihood

During search, the translation likelihood $P(e|f)$ needs to be calculated according to Eq. 3 for every hypothesized segmentation.

To compute Eq. 3, we need to have at hand the individual morph translation probabilities $t(f_j|e_i)$. These can be estimated using the EM algorithm given by (Brown, 1993), which is guaranteed to converge to a global maximum of the likelihood for Model 1. However, running the EM algorithm to optimization for each considered segmentation

model can be computationally expensive, and can result in overtraining. Therefore, in this work we used the likelihood computed after the first EM iteration, which also has the nice property that $P(f|e)$ can be computed incrementally from one segmentation hypothesis to the next.

The incremental updates are derived from the equations for the count collection and probability estimation steps of the EM algorithm as follows. In the count collection step, in the first iteration, we need to compute the fractional counts $c(f_j|e_i)$ (Brown et al., 1993):

$$c(f_j|e_i) = \frac{1}{l+1} (\#f_j)(\#e_i), \quad (4)$$

where $(\#f_j)$ and $(\#e_i)$ denote the number of occurrences of f_j in s_f and e_i in s_e , respectively.

Let f_k denote the word hypothesized to be segmented. Let the resulting two sub-words be f_p and f_q , any of which may or may not previously exist in the vocabulary. Then, according to Eq. (4), as a result of the segmentation no update is needed for $c(f_j|e_i)$ for $j = 1 \dots N$, $j \neq p, q$, $i = 1 \dots M$ (note that f_k no longer exists); and the necessary updates $\Delta c(f_j|e_i)$ for $c(f_j|e_i)$, where $j = p, q$; $i = 1 \dots M$ are given by:

$$\Delta c(f_j|e_i) = \frac{1}{l+1} (\#f_k)(\#e_i). \quad (5)$$

Note that Eq. (5) is nothing but the previous count value for the segmented word, $c(f_k|e_i)$. So, all needed in the count collection step is to copy the set of values $c(f_k|e_i)$ to $c(f_p|e_i)$ and $c(f_q|e_i)$, adding if they already exist.

Then in the probability estimation step, the normalization is performed including the newly added fractional counts.

3.2 Parallelization of search

In an iteration of the algorithm, all words are processed in random order, computing for each word the posterior probability of the generative model after each possible binary segmentation (splitting) of the word. If the highest-scoring split increases the posterior probability compared to not splitting, that split is accepted (for all occurrences of the word) and the resulting sub-words are explored recursively for further segmentations. The process is repeated until an iteration no more results in a significant increase in the posterior probability.

The search algorithm of Morfessor is a greedy algorithm where the costs of the next search points

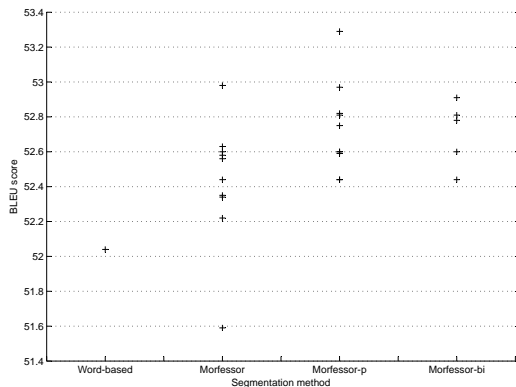


Figure 1: BLEU scores obtained with different segmentation methods. Multiple data points for a system correspond to different random orders in processing the data (Creutz and Lagus, 2007).

are affected by the decision in the current step. This leads to a sequential search and does not lend itself to parallelization.

We propose a slightly modified search procedure, where the segmentation decisions are stored but not applied until the end of an iteration. In this way, the cost calculations (which is the most time-consuming component) can all be performed independently and in parallel. Since the model is not updated at every decision, the search path can differ from that in the sequential greedy search and hence result in different segmentations.

4 Results

We performed *in vivo* testing of the segmentation algorithm on the Turkish side of a Turkish-to-English task. We compared the segmentations produced by Morfessor, Morfessor modified for parallel search (Morfessor-p), and Morfessor with bilingual cost (Morfessor-bi) against the word-based performance. We used the ATR Basic Travel Expression Corpus (BTEC) (Kikui et al., 2006), which contains travel conversation sentences similar to those in phrase-books for tourists traveling abroad. The training corpus contained 19,972 sentences with average sentence length 5.6 and 7.7 words for Turkish and English, respectively. The test corpus consisted of 1,512 sentences with 16 reference translations. We used GIZA++ (Och and Ney, 2003) for post-segmentation token alignments and the Moses toolkit (Koehn et al., 2007) with default parameters for phrase-based translation model generation and decoding. Target language models were

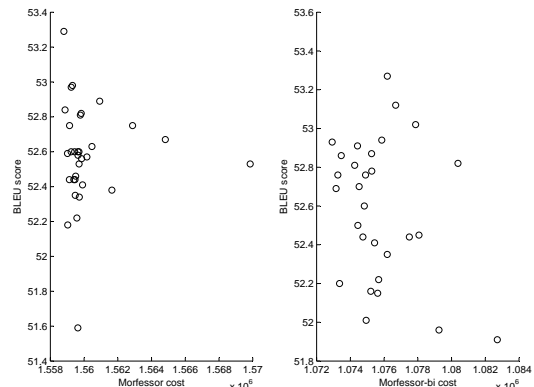


Figure 2: Cost-BLEU plots of Morfessor and Morfessor-bi. Correlation coefficients are -0.005 and -0.279 , respectively.

trained on the English side of the training corpus using the SRILM toolkit (Stolcke, 2002). The BLEU metric (Papineni et al., 2002) was used for translation evaluation.

Figure 1 compares the translation performance obtained using the described segmentation methods. All segmentation methods generally improve the translation performance (Morfessor and Morfessor-p) compared to the word-based models. However, Morfessor-bi, which utilizes both sides of the parallel corpus in segmenting, does not convincingly outperform the monolingual methods.

In order to investigate whether the proposed bilingual segmentation cost correlates any better than the monolingual segmentation cost of Morfessor, we show several cost-BLEU pairs obtained from the final and intermediate segmentations of Morfessor and Morfessor-bi in Fig. 2. The correlation coefficients show that the proposed bilingual metric is somewhat predictive of the translation performance as measured by BLEU, while the monolingual Morfessor cost metric has almost no correlation. Yet, the strong noise in the BLEU scores (vertical variation in Fig. 2) diminishes the effect of this correlation, which explains the inconsistency of the results in Fig. 1. Indeed, in our experiments even though the total cost kept decreasing at each iteration of the search algorithm, the BLEU scores obtained by those intermediate segmentations fluctuated without any consistent improvement.

Table 2 displays sample segmentations produced by both the monolingual and bilingual segmentation algorithms. We can observe that utilizing the English side of the corpus enabled

Count	Morfessor	Morfessor-bi	English Gloss
7	anahtar	anahtar	(the) key
6	anahtar + ımı	anahtar + ımı	my key (ACC.)
5	anahtarla	anahtar + la	with (the) key
4	anahtarı	anahtar + ı	¹ (the) key (ACC.); ² his/her key
3	anahtarı + m	anahtar + ım	my key
3	anahtarı + n	anahtar + ın	¹ your key; ² of (the) key
1	anahtarı + nız	anahtar + ınız	your (pl.) key
1	anahtarı + nı	anahtar + ını	¹ your key (ACC.); ² his/her key (ACC.)
1	anahtar + ınız	anahtar + ınız	your (pl.) key (ACC.)
1	oyun + lar	oyunlar	(the) games
2	oyun + ları	oyunlar + ı	¹ (the) games (ACC.); ² his/her games; ³ their game(s)
1	oyun + ların	oyunlar + ı + n	¹ of (the) games; ² your games
1	oyun + larınızı	oyunlar + ı + n + ızı	your (pl.) games (ACC.)

Table 2: Sample segmentations produced by Morfessor and Morfessor-bi

Morfessor-bi: (i) to consistently identify the root word “anahtar” (top portion), and (ii) to match the English plural word form “games” with the Turkish plural word form “oyunlar” (bottom portion). Monolingual Morfessor is unaware of the target segmentation, and hence it is up to the subsequent translation model training to learn that “oyun” is sometimes translated as “game” and sometimes as “games” in the segmented training corpus.

5 Conclusion

We have presented a method for determining optimal sub-word translation units automatically from a parallel corpus. We have also showed a method of incrementally computing the first iteration parameters of IBM Model-1 between segmentation hypotheses. Being language-independent, the proposed algorithm can be added as a one-time preprocessing step prior to training in a SMT system without requiring any additional data/linguistic resources. The initial experiments presented here show that the translation units learned by the proposed algorithm improves on the word-based baseline in both translation directions.

One avenue for future work is to relax some of the several independence assumptions made in the generative model. For example, independence of consecutive morphs could be relaxed by an HMM model for transitions between morphs (Creutz and Lagus, 2007). Other future work includes optimizing the segmentation of both sides of the corpus and experimenting with other language pairs.

It is also possible that the probability distributions are not discriminative enough to outweigh

the model prior tendencies since the translation probabilities are estimated only crudely (single iteration of Model-1 EM algorithm). A possible candidate solution would be to weigh the translation likelihood more in calculating the overall cost. In fact, this idea could be generalized into a log-linear modeling (e.g., (Poon et al., 2009)) of the various components of the joint corpus likelihood and possibly other features.

Finally, integration of sub-word segmentation with the phrasal lexicon learning process in SMT is desirable (e.g., translation-driven segmentation in Wu (1997)). Hierarchical models (Chiang, 2007) could cover this gap and provide a means to seamlessly integrate sub-word segmentation with statistical machine translation.

Acknowledgements

The authors would like to thank Murat Saraçlar for valuable discussions and guidance in this work, and the anonymous reviewers for very useful comments and suggestions. Murat Saraçlar is supported by the TÜBA-GEBİP award.

References

- Arianna Bisazza and Marcello Federico. 2009. Morphological Pre-Processing for Turkish to English Statistical Machine Translation. In *Proc. of the International Workshop on Spoken Language Translation*, pages 129–135, Tokyo, Japan.
- M.R. Brent. 1999. An efficient, probabilistically sound algorithm for segmentation and word discovery. *Machine Learning*, 34(1):71–105.

- P.F. Brown, V.J. Della Pietra, S.A. Della Pietra, and R.L. Mercer. 1993. The mathematics of statistical machine translation: Parameter estimation. *Computational Linguistics*, 19(2):263–311.
- Pi-Chuan Chang, Michel Galley, and Christopher D. Manning. 2008. Optimizing Chinese word segmentation for machine translation performance. In *Proceedings of the Third Workshop on Statistical Machine Translation*, pages 224–232, Columbus, Ohio.
- David Chiang. 2007. Hierarchical phrase-based translation. *Computational Linguistics*, 33(2):201–228.
- M. Creutz and K. Lagus. 2007. Unsupervised models for morpheme segmentation and morphology learning. *ACM Transactions on Speech and Language Processing*, 4(1):1–34.
- Sajib Dasgupta and Vincent Ng. 2007. High-performance, language-independent morphological segmentation. In *Proceedings of HLT-NAACL*, pages 155–163, Rochester, New York.
- İlknur Durgar El-Kahlout and Kemal Oflazer. 2006. Initial explorations in English to Turkish statistical machine translation. In *Proceedings of the Workshop on Statistical Machine Translation*, pages 7–14, New York City, New York, USA.
- Jakob Elming, Nizar Habash, and Josep M. Crego. 2009. Combination of statistical word alignments based on multiple preprocessing schemes. In Cyrill Goutte, Nicola Cancedda, Marc Dymetman, and George Foster, editors, *Learning Machine Translation*, chapter 5, pages 93–110. MIT Press.
- Nizar Habash and Fatiha Sadat. 2006. Arabic preprocessing schemes for statistical machine translation. In *Proc. of the HLT-NAACL, Companion Volume: Short Papers*, pages 49–52, New York City, USA.
- G. Kikui, S. Yamamoto, T. Takezawa, and E. Sumita. 2006. Comparative study on corpora for speech translation. *IEEE Transactions on Audio, Speech and Language Processing*, 14(5):1674–1682.
- Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondrej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics, Companion Volume: Proceedings of the Demo and Poster Sessions*, pages 177–180, Prague, Czech Republic.
- M. Kurimo, S. Virpioja, V.T. Turunen, G.W. Blackwood, and W. Byrne. 2009. Overview and Results of Morpho Challenge 2009. In *Working notes of the CLEF workshop*.
- Young-Suk Lee. 2004. Morphological analysis for statistical machine translation. In *Proceedings of HLT-NAACL, Companion Volume: Short Papers*, pages 57–60, Boston, Massachusetts, USA.
- Adam Lopez and Philip Resnik. 2006. Word-based alignment, phrase-based translation: What’s the link? In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas (AMTA-06)*, pages 90–99.
- Franz Josef Och and Hermann Ney. 2003. A systematic comparison of various statistical alignment models. *Computational Linguistics*, 29(1):19–51.
- Kemal Oflazer and İlknur Durgar El-Kahlout. 2007. Exploring different representational units in English-to-Turkish statistical machine translation. In *Proceedings of the Second Workshop on Statistical Machine Translation*, pages 25–32, Prague, Czech Republic.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In *Proceedings of 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA.
- Hoifung Poon, Colin Cherry, and Kristina Toutanova. 2009. Unsupervised morphological segmentation with log-linear models. In *Proceedings of HLT-NAACL*, pages 209–217, Boulder, Colorado.
- Fatiha Sadat and Nizar Habash. 2006. Combination of Arabic preprocessing schemes for statistical machine translation. In *Proc. of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics*, pages 1–8, Sydney, Australia.
- Wade Shen, Brian Delaney, and Tim Anderson. 2007. The MIT-LL/AFRL IWSLT-2007 MT system. In *Proc. of the International Workshop on Spoken Language Translation*, Trento, Italy.
- Benjamin Snyder and Regina Barzilay. 2008. Unsupervised multilingual learning for morphological segmentation. In *Proceedings of the 46th Annual Meeting of the Association for Computational Linguistics: HLT*, pages 737–745, Columbus, Ohio.
- A. Stolcke. 2002. SRILM-an extensible language modeling toolkit. In *Seventh International Conference on Spoken Language Processing*, volume 3.
- David Talbot and Miles Osborne. 2006. Modelling lexical redundancy for machine translation. In *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics*, pages 969–976, Sydney, Australia.
- S. Virpioja, J.J. Väyrynen, M. Creutz, and M. Sade-niemi. 2007. Morphology-aware statistical machine translation based on morphs induced in an unsupervised manner. In *Machine Translation Summit XI*, pages 491–498, Copenhagen, Denmark.
- D. Wu. 1997. Stochastic inversion transduction grammars and bilingual parsing of parallel corpora. *Computational Linguistics*, 23(3):377–403.

How spoken language corpora can refine current speech motor training methodologies

Daniil Umanski, Niels O. Schiller

Leiden Institute for Brain and Cognition
Leiden University, The Netherlands

daniil.umanski@gmail.com

N.O.Schiller@hum.leidenuniv.nl

Federico Sangati

Institute for Logic,
Language and Computation

University of Amsterdam, the Netherlands

f.sangati@uva.nl

Abstract

The growing availability of spoken language corpora presents new opportunities for enriching the methodologies of speech and language therapy. In this paper, we present a novel approach for constructing speech motor exercises, based on linguistic knowledge extracted from spoken language corpora. In our study with the Dutch Spoken Corpus, syllabic inventories were obtained by means of automatic syllabification of the spoken language data. Our experimental syllabification method exhibited a reliable performance, and allowed for the acquisition of syllabic tokens from the corpus. Consequently, the syllabic tokens were integrated in a tool for clinicians, a result which holds the potential of contributing to the current state of speech motor training methodologies.

1 Introduction

Spoken language corpora are often accessed by linguists, who need to manipulate specifically defined speech stimuli in their experiments. However, this valuable resource of linguistic information has not yet been systematically applied for the benefit of speech therapy methodologies. This is not surprising, considering the fact that spoken language corpora have only appeared relatively recently, and are still not easily accessible outside the NLP community. Existing applications for selecting linguistic stimuli, although undoubtedly useful, are not based on spoken language data, and are generally not designed for utilization by speech therapists per se (Aichert et al., 2005). As a first attempt to bridge this gap, a mechanism is proposed for utilizing the relevant linguistic information to the service of clinicians. In coordination with speech pathologists, the domain of

speech motor training was identified as an appropriate area of application. The traditional speech motor programs are based on a rather static inventory of speech items, and clinicians do not have access to a modular way of selecting speech targets for training.

Therefore, in this project, we deal with developing an interactive interface to assist speech therapists with constructing individualized speech motor practice programs for their patients. The principal innovation of the proposed system in regard to existing stimuli selection applications is twofold: first, the syllabic inventories are derived from spoken word forms, and second, the selection interface is integrated within a broader platform for conducting speech motor practice.

2 Principles of speech motor practice

2.1 Speech Motor Disorders

Speech motor disorders (SMD) arise from neurological impairments in the motor systems involved in speech production. SMD include acquired and developmental forms of dysarthria and apraxia of speech. Dysarthria refers to the group of disorders associated with weakness, slowness and inability to coordinate the muscles used to produce speech (Duffy, 2005). Apraxia of speech (AOS) is referred to the impaired planning and programming of speech (Ziegler, 2008). Fluency disorders, namely stuttering and cluttering, although not always classified as SMD, have been extensively studied from the speech motor skill perspective (Van Lieshout et al., 2001).

2.2 Speech Motor Training

The goal of speech therapy with SMD patients is establishing and maintaining correct speech motor routines by means of practice. The process of learning and maintaining productive speech motor skills is referred to as speech motor training.

An insightful design of speech motor training exercises is crucial in order to achieve an optimal learning process, in terms of efficiency, retention, and transfer levels (Namasivayam, 2008).

Maas et al. (2008) make the attempt to relate findings from research on non-speech motor learning principles to the case of speech motor training. They outline a number of critical factors in the design of speech motor exercises. These factors include the training program structure, selection of speech items, and the nature of the provided feedback.

It is now generally agreed that speech motor exercises should involve simplified speech tasks. The use of non-sense syllable combinations is a generally accepted method for minimizing the effects of higher-order linguistic processing levels, with the idea of tapping as directly as possible to the motor component of speech production (Smits-Bandstra et al., 2006).

2.3 Selection of speech items

The main considerations in selecting speech items for a specific patient are functional relevance and motor complexity. Functional relevance refers to the specific motor, articulatory or phonetic deficits, and consequently to the treatment goals of the patient. For example, producing correct stress patterns might be a special difficulty for one patient, while producing consonant clusters might be challenging for another. Relative motor complexity of speech segments is much less defined in linguistic terms than, for example, syntactic complexity (Kleinow et al., 2000). Although the part-whole relationship, which works well for syntactic constructions, can be applied to syllabic structures as well (e.g., 'flake' and 'lake'), it may not be the most suitable strategy.

However, in an original recent work, Ziegler presented a non-linear probabilistic model of the phonetic code, which involves units from a sub-segmental level up to the level of metrical feet (Ziegler, 2009). The model is verified on the basis of accuracy data from a large sample of apraxic speakers, and thus provides a quantitative index of a speech segment's motor complexity.

Taken together, it is evident that the task of selecting sets of speech items for an individualized, optimal learning process is far from obvious, and much can be done to assist the clinicians with going through this step.

3 The role of the syllable

The syllable is the primary speech unit used in studies on speech motor control (Namasivayam, 2008). It is also the basic unit used for constructing speech items in current methodologies of speech motor training (Kent, 2000). Since the choice of syllabic tokens is assumed to affect speech motor learning, it would be beneficial to have access to the syllabic inventory of the spoken language. Besides the inventory of spoken syllables, we are interested in the distribution of syllables across the language.

3.1 Syllable frequency effects

The observation that syllables exhibit an exponential distribution in English, Dutch and German has led researchers to infer the existence of a 'mental syllabary' component in the speech production model (Schiller et al., 1996). Since this hypothesis assumes that production of high frequency syllables relies on highly automated motor gestures, it bears direct consequences on the utility of speech motor exercises. In other words, manipulating syllable sets in terms of their relative frequency is expected to have an effect on the learning process of new motor gestures. This argument is supported by a number of empirical findings. In a recent study, Staiger et al. report that syllable frequency and syllable structure play a decisive role with respect to articulatory accuracy in the spontaneous speech production of patients with AOS (Staiger et al., 2008). Similarly, (Laganaro, 2008) confirms a significant effect of syllable frequency on production accuracy in experiments with speakers with AOS and speakers with conduction aphasia.

3.2 Implications on motor learning

In that view, practicing with high-frequency syllables could promote a faster transfer of skills to everyday language, as the most 'required' motor gestures are being strengthened. On the other hand, practicing with low-frequency syllables could potentially promote plasticity (or 'stretching') of the speech motor system, as the learner is required to assemble motor plans from scratch, similar to the process of learning to pronounce words in a foreign language. In the next section, we describe our study with the Spoken Dutch Corpus, and illustrate the performed data extraction strategies.

4 A study with the Spoken Dutch Corpus

The Corpus Gesproken Nederlands (CGN) is a large corpus of spoken Dutch¹. The CGN contains manually verified phonetic transcriptions of 53,583 spoken forms, sampled from a wide variety of communication situations. A spoken form reports the phoneme sequence as it was actually uttered by the speaker as opposed to the canonical form, which represents how the same word would be uttered in principle.

4.1 Motivation for accessing spoken forms

In contrast to written language corpora, such as CELEX (Baayenet al., 1996), or even a corpus like TIMIT (Zue et al., 1996), in which speakers read prepared written material, spontaneous speech corpora offer an access to an informal, unscripted speech on a variety of topics, including speakers from a range of regional dialects, age and educational backgrounds.

Spoken language is a dynamic, adaptive, and generative process. Speakers most often deviate from the canonical pronunciation, producing segment reductions, deletions, insertions and assimilations in spontaneous speech (Mitterer, 2008). The work of Greenberg provides an in-depth account on the pronunciation variation in spoken English. A detailed phonetic transcription of the Switchboard corpus revealed that the spectral properties of many phonetic elements deviate significantly from their canonical form (Greenberg, 1999).

In the light of the apparent discrepancy between the canonical forms and the actual spoken language, it becomes apparent that deriving syllabic inventories from spoken word forms will approximate the reality of spontaneous speech production better than relying on canonical representations. Consequently, it can be argued that clinical applications will benefit from incorporating speech items which optimally converge with the 'live' realization of speech.

4.2 Syllabification of spoken forms

The syllabification information available in the CGN applies only to the canonical forms of words, and no syllabification of spoken word forms exists. The methods of automatic syllabification have been applied and tested exclusively on canonical word forms (Bartlett, 2007). In order to obtain the syllabic inventory of spoken language per se,

¹(see <http://lands.let.kun.nl/cgn/>)

a preliminary study on automatic syllabification of spoken word forms has been carried out. Two methods for dealing with the syllabification task were proposed, the first based on an n-gram model defined over sequences of phonemes, and the second based on statistics over syllable units. Both algorithms accept as input a list of possible segmentations of a given phonetic sequence, and return the one which maximizes the score of the specific function they implement. The list of possible segmentations is obtained by exhaustively generating all possible divisions of the sequence, satisfying the condition of keeping exactly one vowel per segment.

4.3 Syllabification Methods

The first method is a reimplementation of the work of (Schmid et al., 2007). The authors describe the syllabification task as a tagging problem, in which each phonetic symbol of a word is tagged as either a syllable boundary ('B') or as a non-syllable boundary ('N'). Given a set of possible segmentations of a given word, the aim is to select the one, viz. the tag sequence \hat{b}_1^n , which is more probable for the given phoneme sequence p_1^n , as shown in equation (1). This probability in equations (3) is reduced to the joint probability of the two sequences: the denominator of equation (2) is in fact constant for the given list of possible syllabifications, since they all share the same sequence of phonemes. Equation (4) is obtained by introducing a Markovian assumption of order 3 in the way the phonemes and tags are jointly generated

$$\hat{b}_1^n = \arg \max_{b_1^n} P(b_1^n | p_1^n) \quad (1)$$

$$= \arg \max_{b_1^n} P(b_1^n, p_1^n) / P(p_1^n) \quad (2)$$

$$= \arg \max_{b_1^n} P(b_1^n, p_1^n) \quad (3)$$

$$= \arg \max_{b_1^n} \prod_{i=1}^{n+1} P(b_i, p_i | b_{i-3}^{i-1}, p_{i-3}^{i-1}) \quad (4)$$

The second syllabification method relies on statistics over the set of syllables unit and bigram (bisegments) present in the training corpus. Broadly speaking, given a set of possible segmentations of a given phoneme sequence, the algorithm, selects the one which maximizes the presence and frequency of its segments.

Corpus	Phonemes		Syllables	
	Boundaries	Words	Boundaries	Words
CGN_Dutch	98.62	97.15	97.58	94.99
CELEX_Dutch	99.12	97.76	99.09	97.70
CELEX_German	99.77	99.41	99.51	98.73
CELEX_English	98.86	97.96	96.37	93.50

Table 1: Summary of syllabification results on canonical word forms.

4.4 Results

The first step involved the evaluation of the two algorithms on syllabification of canonical word forms. Four corpora comprising three different languages (English, German, and Dutch) were evaluated: the CELEX2 corpora (Baayen et al., 1996) for the three languages, and the Spoken Dutch Corpus (CGN). All the resources included manually verified syllabification transcriptions. A 10-fold cross validation on each of the corpora was performed to evaluate the accuracy of our methods. The evaluation is presented in terms of percentage of correct syllable boundaries², and percentage of correctly syllabified words.

Table 1 summarizes the obtained results. For the CELEX corpora, both methods produce almost equally high scores, which are comparable to the state of the art results reported in (Bartlett, 2007). For the Spoken Dutch Corpus, both methods demonstrate quite high scores, with the phoneme-level method showing an advantage, especially with respect to correctly syllabified words.

4.5 Data extraction

The process of evaluating syllabification of spoken word forms is compromised by the fact that there exists no gold annotation for the pronunciation data in the corpus. Therefore, the next step involved applying both methods on the data set and comparing the two solutions. The results revealed that the two algorithms agree on 94.29% of syllable boundaries and on 90.22% of whole word syllabification. Based on the high scores reported for lexical word forms syllabification, an agreement between both methods most probably implies a correct solution. The 'disagreement' set can be assumed to represent the class of ambiguous cases, which are the most problematic for automatic syllabification. As an example, consider

²Note that recall and precision coincide since the number of boundaries (one less than the number of vowels) is constant for different segmentations of the same word.

the following pair of possible syllabification, on which the two methods disagree: 'bEl-kOm-pjut' vs 'bEl-kOmp-jut'³.

Motivated by the high agreement score, we have applied the phoneme-based method on the spoken word forms in the CGN, and compiled a syllabic inventory. In total, 832,236 syllable tokens were encountered in the corpus, of them 11,054 unique syllables were extracted and listed. The frequencies distribution of the extracted syllabary, as can be seen in Figure 1, exhibits an exponential curve, a result consistent with earlier findings reported in (Schiller et al., 1996). According to our statistics, 4% of unique syllable tokens account for 80% of all extracted tokens, and 10% of unique syllables account for 90% respectively. For each extracted syllable, we have recorded its structure, frequency rank, and the articulatory characteristics of its consonants. Next, we describe the speech items selection tool for clinicians.

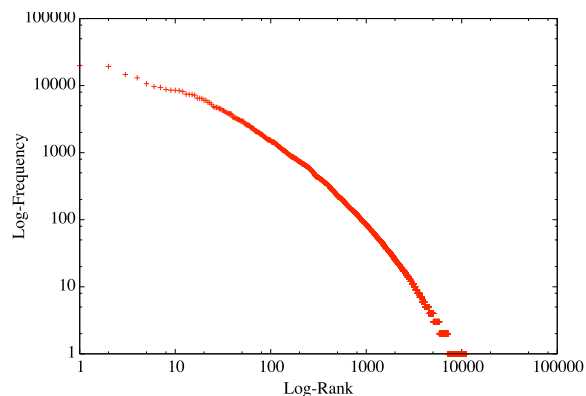


Figure 1: Syllable frequency distribution over the spoken forms in the Dutch Spoken Corpus.

The x-axis represents 625 ranked frequency bins. The y-axis plots the total number of syllable tokens extracted for each frequency bin.

³A manual evaluation of the disagreement set revealed a clear advantage for the phoneme-based method

5 An interface for clinicians

In order to make the collected linguistic information available for clinicians, an interface has been built which enables clinicians to compose individual training programs. A training program consists of several training sessions, which in turn consists of a number of exercises. For each exercise, a number of syllable sets are selected, according to the specific needs of the patient. The main function of the interface, thus, deals with selection of customized syllable sets, and is described next. The rest of the interface deals with the different ways in which the syllable sets can be grouped into exercises, and how exercises are scheduled between treatment sessions.

5.1 User-defined syllable sets

The process starts with selecting the number of syllables in the current set, a number between one and four. Consequently, the selected number of 'syllable boxes' appear on the screen. Each box allows for a separate configuration of one syllable group. As can be seen in Figure 2, a syllable box contains a number of menus, and a text grid at the bottom of the box.



Figure 2: A snapshot of the part of the interface allowing configuration of syllable sets

Here follows the list of the parameters which the user can manipulate, and their possible values:

- Syllable Type⁴
- Syllable Frequency⁵

⁴CV, CVC, CCV, CCVC, etc.

⁵Syllables are divided in three rank groups - high, medium, and low frequency.

- Voiced - Unvoiced consonant⁶
- Manner of articulation⁷
- Place of articulation⁸

Once the user selects a syllable type, he/she can further specify each consonant within that syllable type in terms of voiced/unvoiced segment choice and manner and place of articulation. For the sake of simplicity, syllable frequency ranks have been divided in three rank groups. Alternatively, the user can bypass this criterion by selecting 'any'. As the user selects the parameters which define the desired syllable type, the text grid is continuously filled with the list of syllables satisfying these criteria, and a counter shows the number of syllables currently in the grid.

Once the configuration process is accomplished, the syllables which 'survived' the selection will constitute the speech items of the current exercise, and the user proceeds to select how the syllable sets should be grouped, scheduled and so on.

6 Final remarks

6.1 Future directions

A formal usability study is needed in order to establish the degree of utility and satisfaction with the interface. One question which demands investigation is the degrees of choice that the selection tool should provide. With too many variables and hinges of choice, the configuration process for each patient might become complicated and time consuming. Therefore, a usability study should provide guidelines for an optimal design of the interface, so that its utility for clinicians is maximized.

Furthermore, we plan to integrate the proposed interface within an computer-based interactive platform for speech therapy. A seamless integration of a speech items selection module within biofeedback games for performing exercises with these items seems straight forward, as the selected items can be directly embedded (e.g., as text symbols or more abstract shapes) in the graphical environment where the exercises take place.

⁶when applicable

⁷for a specific consonant. Plosives, Fricatives, Sonorants

⁸for a specific consonant. Bilabial, Labio-Dental, Alveolar, Post-Alveolar, Palatal, Velar, Uvular, Glottal

Acknowledgments

This research is supported with the 'Mosaic' grant from The Netherlands Organisation for Scientific Research (NWO). The authors are grateful for the anonymous reviewers for their constructive feedback.

References

- Aichert, I., Ziegler, W. 2004. *Syllable frequency and syllable structure in apraxia of speech*. *Brain and Language*, 88, 148-159.
- Aichert, I., Marquardt, C., Ziegler, W. 2005. *Frequenzen sublexikalischer Einheiten des Deutschen: CELEX-basierte Datenbanken*. *Neurolinguistik*, 19, 55-81
- Baayen R.H., Piepenbrock R. and Gulikers L. 1996. *CELEX2. Linguistic Data Consortium, Philadelphia*.
- Bartlett, S. 2007. *Discriminative approach to automatic syllabication*. Masters thesis, Department of Computing Science, University of Alberta.
- Duffy, J.R. 2005. *Motor speech disorder: Substrates, Differential Diagnosis, and Management*. (2nd Ed.) 507-524. St. Louis, MO: Elsevier Mosby
- Greenberg, S. 1999. *Speaking in shorthand: a syllable-centric perspective for understanding pronunciation variation*. *Speech Comm.*, 29(2-4):159-176
- Kent, R. 2000. *Research on speech motor control and its disorders, a review and perspectives*. *Speech Comm.*, 29(2-4):159-176 J.
- Kleinow, J., Smith, A. 2000. *Influences of length and syntactic complexity on the speech motor stability of the uent speech of adults who stutter*. *Journal of Speech, Language, and Hearing Research*, 43, 548559.
- Laganaro, M. 2008. *Is there a syllable frequency effect in aphasia or in apraxia of speech or both?* *Aphasiology*, Volume 22, Number 11, November 2008 , pp. 1191-1200(10)
- Maas, E., Robin, D.A., Austermann Hula, S.N., Freedman, S.E., Wulf, G., Ballard, K.J., Schmidt, R.A. 2008. *Principles of Motor Learning in Treatment of Motor Speech Disorders* *American Journal of Speech-Language Pathology*, 17, 277-298.
- Mitterer, H. 2008. *How are words reduced in spontaneous speech?* In A. Botinis (Ed.), *Proceedings of the ISCA Tutorial and Research Workshop on Experimental Linguistics* (pages 165-168). University of Athens.
- Namasivayam, A.K., van Lieshout, P. 2008. *Investigating speech motor practice and learning in people who stutter* *Journal of Fluency Disorders* 33 (2008) 3251
- Schiller, N. O., Meyer, A. S., Baayen, R. H., Levelt, W. J. M. 1996. *A Comparison of Lexeme and Speech Syllables in Dutch*. *Journal of Quantitative Linguistics*, 3, 8-28.
- Schmid H., Möbius B. and Weidenkaff J. 2007. *Tagging Syllable Boundaries With Joint N-Gram Models*. *Proceedings of Interspeech-2007 (Antwerpen)*, pages 2857-2860.
- Smits-Bandstra, S., DeNil, L. F., Saint-Cyr, J. 2006. *Speech and non-speech sequence skill learning in adults who stutter*. *Journal of Fluency Disorders*, 31,116136.
- Staiger, A., Ziegler, W. 2008. *Syllable frequency and syllable structure in the spontaneous speech production of patients with apraxia of speech*. *Aphasiology*, Volume 22, Number 11, November 2008 , pp. 1201-1215(15)
- Tjaden, K. 2000. *Exploration of a treatment technique for prosodic disturbance following stroke training*. *Clinical Linguistics and Phonetics* 2000, Vol. 14, No. 8, Pages 619-641
- Riley, J., Riley, G. 1995. *Speech motor improvement program for children who stutter*. In C.W. Starkweather, H.F.M. Peters (Eds.), *Stuttering* (pp.269-272) New York: Elsevier
- Van Lieshout, P. H. H. M. 2001. *Recent developments in studies of speech motor control in stuttering*. In B. Maassen, W. Hulstijn, R. D. Kent, H. F. M. Peters, P. H. H. M. Van Lieshout (Eds.), *Speech motor control in normal and disordered speech*(pp. 286290). Nijmegen, The Netherlands: Vantilt.
- Ziegler W. 2009. *Modelling the architecture of phonetic plans: Evidence from apraxia of speech*. *Language and Cognitive Processes* 24, 631 - 661
- Ziegler W. 2008. *Apraxia of speech*. In: Goldenberg G, Miller B (Eds.), *Handbook of Clinical Neurology*, Vol. 88 (3rd series), pp. 269 - 285. Elsevier. London
- Zue, V.W. and Seneff, S. 1996. *Transcription and alignment of the TIMIT database*. In *Recent Research Towards Advanced Man-Machine Interface Through Spoken Language*. H. Fujisaki (ed.), Amsterdam: Elsevier, 1996, pp. 515-525.

Mood Patterns and Affective Lexicon Access in Weblogs

Thin Nguyen

Curtin University of Technology
Bentley, WA 6102, Australia

thin.nguyen@postgrad.curtin.edu.au

Abstract

The emergence of social media brings chances, but also challenges, to linguistic analysis. In this paper we investigate a novel problem of discovering patterns based on emotion and the association of moods and affective lexicon usage in blogosphere, a representative for social media. We propose the use of normative emotional scores for English words in combination with a psychological model of emotion measurement and a nonparametric clustering process for inferring meaningful emotion patterns automatically from data. Our results on a dataset consisting of more than 17 million mood-groundtruthed blogposts have shown interesting evidence of the emotion patterns automatically discovered that match well with the core-affect emotion model theorized by psychologists. We then present a method based on information theory to discover the association of moods and affective lexicon usage in the new media.

1 Introduction

Social media provides communication and interaction channels where users can freely participate in, express their opinions, make their own content, and interact with other users. Users in this new media are more comfortable in expressing their feelings, opinions, and ideas. Thus, the resulting user-generated content tends to be more subjective than other written genres, and thus, is more appealing to be investigated in terms of subjectivity and sentiment analysis. Research in sentiment analysis has recently attracted much attention (Pang and Lee, 2008), but modeling emotion

patterns and studying the affective lexicon used in social media have received little attention.

Work in sentiment analysis in social media is often limited to finding the sentiment sign in the dipole pattern (negative/positive) for given text. Extensions to this task include the three-class classification (adding neutral to the polarity) and locating the value of emotion the text carries across a spectrum of valence scores. On the other hand, it is well appreciated by psychologists that sentiment has much richer structures than the aforementioned simplified polarity. For example, emotion – a form of expressive sentiment – was suggested by psychologists to be measured in terms of *valence* and *arousal* (Russell, 2009). Thus, we are motivated to analyze the sentiment in blogosphere in a more fine-grained fashion. In this paper we study the grouping behaviors of the emotion, or emotion patterns, expressed in the blogposts. We are inspired to get insights into the question of whether these structures can be discovered directly from data without the cost of involving human participants as in traditional psychological studies. Next, we aim to study the relationship between the data-driven emotion structures discovered and those proposed by psychologists.

Work on the analysis of effects of sentiment on lexical access is great in a psychology perspective. However, to our knowledge, limited work exists to examine the same tasks in social media context.

The contribution in this paper is twofold. To our understanding, we study a novel problem of emotion-based pattern discovery in blogosphere. We provide an initial solution for the matter using a combination of psychological models, affective norm scores for English words, a novel feature representation scheme, and a nonparametric clustering to automatically group moods into meaningful emotion patterns. We believe that we are the first to consider the matter of data-driven emotion pattern discovery at the scale presented in this

paper. Secondly, we explore a novel problem of detecting the mood – affective lexicon usage correlation in the new media, and propose a novel use of a term-goodness criterion to discover this sentiment – linguistic association.

2 Related Work

Much work in sentiment analysis measures the value of emotion the text convey in a continuum range of valence (Pang and Lee, 2008). Emotion patterns have often been used in sentiment analysis limited to this one-dimensional formulation. On the other hand, in psychology, emotions have often been represented in dimensional and discrete perspectives. In the former, emotion states are conceptualized as combinations of some factors like valence and arousal. In contrast, the latter style argues that each emotion has a unique coincidence of experience, psychology and behavior (Mauss and Robinson, 2009). Our work utilizes the dimensional representation, and in particular, the core-affect model (Russell, 2009), which encodes emotion states along the valence and arousal dimensions. The sentiment scoring for emotion bearing words is available in a lexicon known as Affective Norms for English Words (ANEW) (Bradley and Lang, 1999). Related work making use of ANEW includes (Dodds and Danforth, 2009) for estimating happiness levels in three types of data: song lyrics, blogs, and the State of the Union addresses.

From a psychological perspective, for estimating mood effects in lexicon decisions, (Chastain et al., 1995) investigates the influence of moods on the access of affective words. For learning affect in blogosphere, (Leshed and Kaye, 2006) utilizes Support Vector Machines (SVM) to predict moods for coming blog posts and detect mood synonymy.

3 Moods and Affective Lexicon Access

3.1 Mood Pattern Detection

Livejournal provides a comprehensive set of 132 moods for users to tag their moods when blogging. The provided moods range diversely in the emotion spectrum but typically are observed to fall into soft clusters such as happiness (*cheerful* or *grateful*) or sadness (*discontent* or *uncomfortable*). We call each cluster of these moods an *emotion pattern* and aim to detect them in this paper.

We observe that the blogposts tagged with moods in the same emotion pattern have similar

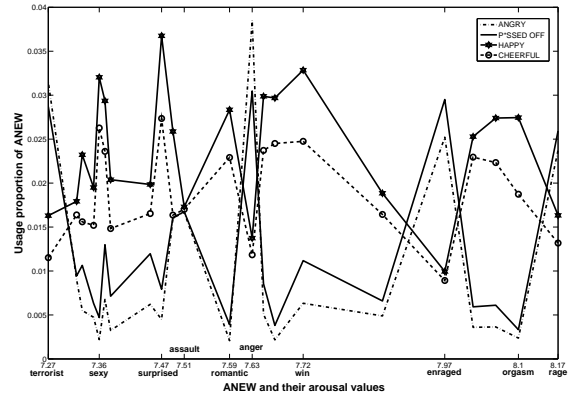


Figure 1: ANEW usage proportion in the posts tagged with *happy/cheerful* and *angry/p*ssed off*

proportions in the usage of ANEW. For example, in Figure 1 – a plot of the usage of ANEW having arousal in the range of 7.2 – 8.2 in the blogposts – we could see that the ANEW usage patterns of *happy/cheerful* and *angry/p*ssed off* are well separated. *Anger*, *enraged*, and *rage* will be most likely found in the *angry/p*ssed off* tagged posts and least likely found in the *happy/cheerful* ones. In contrast, the ANEW as *romantic* or *surprised* are not commonly used in the posts tagged with *angry/p*ssed off* but most popularly used in the *happy/cheerful* ones; suggesting that, the similarity between ANEW usage patterns can be used as a basis to study the structure of mood space.

Let us denote by \mathcal{B} the corpus of all blogposts and by $\mathcal{M} = \{sad, happy, \dots\}$ the predefined set of moods ($|\mathcal{M}| = 132$). Each blogpost $b \in \mathcal{B}$ in the corpus is labeled with a mood $l_b \in \mathcal{M}$. Denote by n the number of ANEW ($n = 1034$). Let $\mathbf{x}^m = [\mathbf{x}_1^m, \dots, \mathbf{x}_i^m, \dots, \mathbf{x}_n^m]$ be the vector representing the usage of ANEW by the mood m . Thus, $\mathbf{x}_i^m = \sum_{b \in \mathcal{B}, l_b = m} c_{ib}$, where c_{ib} is the counting of the ANEW i -th occurrence in the blogpost b tagged with the mood m . The usage vector is normalized so that $\sum_{i=1}^n \mathbf{x}_i^m = 1$ for all $m \in \mathcal{M}$. To discover the grouping of the moods based on the usage vectors we use a nonparametric clustering algorithm known as Affinity Propagation (AP) (Frey and Dueck, 2007). AP is desirable here because it automatically discovers the number of clusters as well as the cluster exemplars. The algorithm only requires the pairwise similarities between moods, which we compute based on the Euclidean distances for simplicity.

To map the emotion patterns detected to their psychological meaning, we proceed to measure

the sentiment scores of those $|\mathcal{M}|$ mood words. In particular, we use ANEW (Bradley and Lang, 1999), which is a set of 1034 sentiment conveying English words. The valence and arousal of moods are assigned by those of the same words in the ANEW lexicon. For those moods which are not in ANEW, their values are assigned by those of the nearest father words in the mood hierarchical tree¹, where those moods conveying the same meaning, to some extent, are in the same level of the tree. Thus, each member of the mood clusters can be placed onto the a 2D representation along the valence and arousal dimensions, making it feasible to compare with the *core-affect* model (Russell, 2009) theorized by psychologists.

3.2 Mood and ANEW Usage Association

To study the statistical strength of an ANEW word with respect to a particular mood, the information gain measure (Mitchell, 1997) is adopted. Given a collection of blog posts \mathcal{B} consisting of those tagged or not tagged with a target class attribute mood m . The entropy of \mathcal{B} relative to this binary classification is

$$\mathcal{H}(\mathcal{B}) = -p_{\oplus} \log_2(p_{\oplus}) - p_{\ominus} \log_2 p_{\ominus}$$

where p_{\oplus} and p_{\ominus} are the proportions of the posts tagged and not tagged with m respectively.

The entropy of \mathcal{B} relative to the binary classification given a binary attribute A (e.g. if the word A present or not) observed is computed as

$$\mathcal{H}(\mathcal{B}|A) = \frac{|\mathcal{B}_{\oplus}|}{|\mathcal{B}|} \mathcal{H}(\mathcal{B}_{\oplus}) + \frac{|\mathcal{B}_{\ominus}|}{|\mathcal{B}|} \mathcal{H}(\mathcal{B}_{\ominus})$$

where \mathcal{B}_{\oplus} is the subset of \mathcal{B} for which attribute A is present in the corpus and \mathcal{B}_{\ominus} is the subset of \mathcal{B} for which attribute A is absent in the corpus.

The information gain of an attribute ANEW A in classifying the collection with respect to the target class attribute mood m , $IG(m, A)$, is the reduction in entropy caused by partitioning the examples according to the attribute A . Thus,

$$IG(m, A) = \mathcal{H}(\mathcal{B}) - \mathcal{H}(\mathcal{B}|A)$$

With respect to a given mood m , those ANEW having high information gain are considered likely to be associated with the mood. This measure, also often considered a term-goodness criterion, outperforms others in feature selection in text categorization (Yang and Pedersen, 1997).

¹<http://www.livejournal.com/moodlist.bml>

4 Experimental Results

4.1 Mood Patterns

We use a large Livejournal blogpost dataset, which contains more than 17 million blogposts tagged with the predefined moods. These journals were posted from May 1, 2001 to April 23, 2005. The ANEW usage vectors of all moods are subjected to a clustering to learn emotion patterns. After running the Affinity Propagation algorithm, 16 patterns of moods are clustered as below (the moods in upper case are the exemplars).

-
1. CHEERFUL, ecstatic, jubilant, giddy, happy, excited, energetic, bouncy, chipper
 2. PENSIVE, determined, contemplative, thoughtful
 3. REJUVENATED, optimistic, relieved, refreshed, hopeful, peaceful
 4. QUIXOTIC, surprised, enthralled, devious, geeky, creative, recumbent, artistic, impressed, amused, complacent, curious, weird
 5. CRAZY, horny, giggly, high, flirty, hyper, drunk, naughty, dorky, ditzy, silly
 6. MELLOW, pleased, satisfied, relaxed, content, anxious, good, full, calm, okay
 7. GRATEFUL, loved, thankful, touched
 8. AGGRAVATED, irritated, bitchy, annoyed, frustrated, cynical
 9. ANGRY, p*ssed off, infuriated, irate, enraged
 10. GLOOMY, jealous, envious, rejected, confused, worried, lonely, guilty, scared, pessimistic, discontent, distressed, indescribable, crushed, depressed, melancholy, numb, morose, sad, sympathetic
 11. PRODUCTIVE, accomplished, working, nervous, busy, rushed
 12. TIRED, sore, lazy, sleepy, awake, groggy, exhausted, lethargic, drained
 13. NAUSEATED, sick
 14. MOODY, disappointed, grumpy, cranky, stressed, uncomfortable, crappy
 15. THIRSTY, nerdy, mischievous, hungry, dirty, hot, cold, bored, blah
 16. EXANIMATE, intimidated, predatory, embarrassed, restless, nostalgic, indifferent, listless, apathetic, blank, shocked
-

Generally, the patterns 1–7 contain moods in high valence (pleasure) and the patterns 8–16 include mood in low valence (displeasure). To examine whether members in these emotion patterns

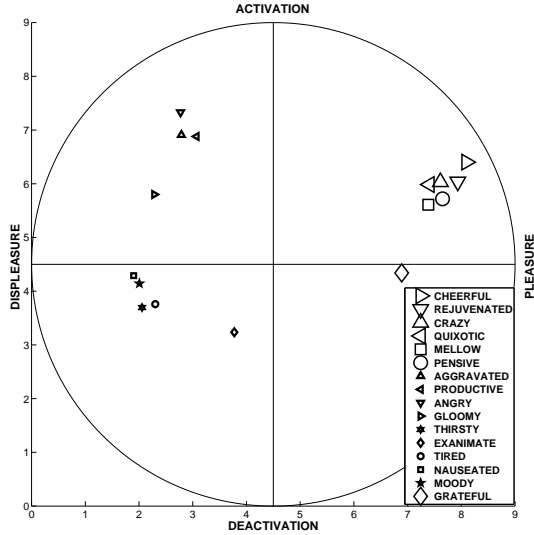


Figure 4: Discovered emotion patterns in the affect circle

follow an affect concept, we place them on the affect circle (Russell, 2009). We learn that nearly all members in the same patterns express a common affect concept. Those moods in the patterns with *cheerful*, *pensive*, and *rejuvenated* as the exemplars are mostly located in the first quarter of the affect circle ($0^{\circ} - 90^{\circ}$), which should contain moods being high in both pleasure and activation measures. Meanwhile, many members of the *angry* and *aggravated* patterns are found in the second quarter ($90^{\circ} - 180^{\circ}$), which roughly means that those moods express the feeling of sadness in the high of activation. The patterns with the exemplars *nauseated* and *tired* contain a majority of moods found in the third quarter ($180^{\circ} - 270^{\circ}$), which could be representatives for the mood fashion of sadness and deactivation. In addition, the *grateful* group could be a representative for moods which are both low in pleasure and in the degree of activation ($270^{\circ} - 360^{\circ}$ of the affect circle). Thus, the clustering process based on the ANEW usage could separate moods having similar affect scores into corresponding segments in the circle proposed in (Russell, 2009).

To visualize mood patterns that have been detected, we plot these emotion modes on the affect circle plane in Figure 4. For each pattern, the valence and arousal are computed by averaging of the values of those moods in the quarter where most of the members in the pattern are.

To further visualize the similarity of moods, the ANEW usage vectors are subject to a classical multidimensional scaling (Borg and Groenen,

Mood	Top ANEW words associated
Cheerful	fun, happy, hate, good, christmas, merry, birthday, cute, sick, love
Happy	happy, hate, fun, good, birthday, sick, love, mind, alone, bored
Angry	angry, hate, fun, mad, love, anger, good, stupid, pretty, movie
P*ssed off	hate, stupid, mad, love, hell, fun, good, god, pretty, movie
Gloomy	sad, depressed, hate, wish, life, alone, lonely, upset, pain, heart
Sad	sad, fun, heart, upset, wish, funeral, hurt, pretty, loved, cancer

(a) Moods and the most associated ANEW words

ANEW	Most likely moods	Least likely moods
Desire	contemplative, thoughtful	enraged, drained
Anger	angry, p*ssed off	nauseated, grateful
Accident	sore, bored	exanimate, indifferent
Terrorist	angry, cynical	rejuvenated, touched
Wine	drunk, p*ssed off	ditzy, okay

(b) ANEW words and the most associated moods

Table 1: Mood and ANEW correlation

2005) (MDS) and a hierarchical clustering. Figure 2 and Figure 3 show views of the distance between moods, based on the Euclidean measure of their corresponding ANEW usage, using MDS and hierarchical clustering respectively.

4.2 Mood and ANEW Association

Based on the IG values between moods and ANEW, we learn the correlation of moods and the affective lexicon. With respect to a given mood, those ANEW having high information gain are most likely to be found in the blogposts tagged with the mood. The ANEW most likely happened in the blogposts tagged with a given mood are shown in Table 1a; the most likely moods for the blog posts containing a given ANEW are shown in Table 1b.

The ANEW used in the blog posts tagged with moods in the same pattern are more similar than those in the posts tagged with moods in different patterns. In Table 1a, the most associated ANEW

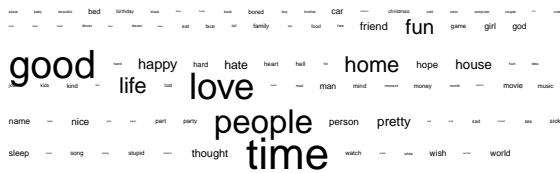


Figure 5: Top 100 ANEW words used in the dataset

in the blogposts tagged with *cheerful* are more similar to those in *happy* ones than those in *angry* or *p*ssed off* ones.

For a given mood, a majority of the ANEW used in the blog posts tagged with the mood is similar in the valence with the mood. The occurrence of some ANEW having valence much different with the tagging mood, e.g. the ANEW *hate* in the posts tagged with *cheerful* or *happy* moods, might be the result of a negation construction used in the text or of other context.

For a given ANEW, the most likely moods tagged to the blog posts containing the word are similar with the word in the affective scores. In addition, the least likely moods are much different with the ANEW in the affect measure. A plot of top ANEWs used in the blogposts is shown in Figure 5.

Other than the ANEW conveying abstract concept, e.g. *desire* or *anger*, those ANEW expressing more concrete existence, e.g. *terrorist* or *accident*, might be a good source for learning opinions from social network towards the things. In the corpus, the posts containing the ANEW *terrorist* are most likely tagged with *angry* or *cynical* moods. Also, the posts containing the ANEW *accident* are most likely tagged with *bored* and *sore* moods.

5 Conclusion and Future Work

We have investigated the problems of emotion-based pattern discovery and mood – affective lexicon usage correlation detection in blogosphere. We presented a method for feature representation based on the affective norms of English scores usage. We then presented an unsupervised approach using Affinity Propagation, a nonparametric clustering algorithm that does not require the number of clusters a priori, for detecting emotion patterns in blogosphere. The results are showing that those automatically discovered patterns match well with the core-affect model for emotion, which is independently formulated in the psychology literature. In addition, we proposed a novel use of a term-

goodness criterion to discover mood–lexicon correlation in blogosphere, giving hints on predicting moods based on the affective lexicon usage and vice versa in the social media. Our results could also have potential uses in sentiment-aware social media applications.

Future work will take into account the temporal dimension to trace changes in mood patterns over time in blogosphere. Another direction is to integrate negation information to learn more cohesive association in affect scores between moods and affective words. In addition, a new affective lexicon could be automatically detected based on learning correlation of the blog text and the moods tagged.

References

- I. Borg and P.J.F. Groenen. 2005. *Modern multidimensional scaling: Theory and applications*. Springer Verlag.
- M.M. Bradley and P.J. Lang. 1999. Affective norms for English words (ANEW): Stimuli, instruction manual and affective ratings. Technical report, University of Florida.
- G. Chastain, P.S. Seibert, and F.R. Ferraro. 1995. Mood and lexical access of positive, negative, and neutral words. *Journal of General Psychology*, 122(2):137–157.
- P.S. Dodds and C.M. Danforth. 2009. Measuring the happiness of large-scale written expression: Songs, blogs, and presidents. *Journal of Happiness Studies*, pages 1–16.
- B.J. Frey and D. Dueck. 2007. Clustering by passing messages between data points. *Science*, 315(5814):972.
- G. Leshed and J.J. Kaye. 2006. Understanding how bloggers feel: recognizing affect in blog posts. In *Proc. of ACM Conf. on Human Factors in Computing Systems (CHI)*.
- I.B. Mauss and M.D. Robinson. 2009. Measures of emotion: A review. *Cognition & emotion*, 23:2(2):209–237.
- T. Mitchell. 1997. *Machine Learning*. McGraw Hill.
- B. Pang and L. Lee. 2008. Opinion mining and sentiment analysis. *Foundations and Trends in Information Retrieval*, 2(1-2):1–135.
- J.A. Russell. 2009. Emotion, core affect, and psychological construction. *Cognition & Emotion*, 23:7(1):1259–1283.
- Y. Yang and J.O. Pedersen. 1997. A comparative study on feature selection in text categorization. In *Proc. of Intl. Conf. on Machine Learning (ICML)*, pages 412–420.

Growing Related Words from Seed via User Behaviors: A Re-ranking Based Approach

Yabin Zheng

Zhiyuan Liu

Lixing Xie

State Key Laboratory on Intelligent Technology and Systems

Tsinghua National Laboratory for Information Science and Technology

Department of Computer Science and Technology, Tsinghua University, Beijing 100084, China

{yabin.zheng, lzy.thu, lavender087}@gmail.com

Abstract

Motivated by Google Sets, we study the problem of growing related words from a single seed word by leveraging user behaviors hiding in user records of Chinese input method. Our proposed method is motivated by the observation that the more frequently two words co-occur in user records, the more related they are. First, we utilize user behaviors to generate candidate words. Then, we utilize search engine to enrich candidate words with adequate semantic features. Finally, we reorder candidate words according to their semantic relatedness to the seed word. Experimental results on a Chinese input method dataset show that our method gains better performance.

1 Introduction

What is the relationship between “自然语言处理” (Natural Language Processing) and “人工智能” (Artificial Intelligence)? We may regard NLP as a research branch of AI. Problems arise when we want to find more words related to the input query/seed word. For example, if seed word “自然语言处理” (Natural Language Processing) is entered into Google Sets (Google, 2010), Google Sets returns an ordered list of related words such as “人工智能” (Artificial Intelligence) and “计算机” (Computer). Generally speaking, it performs a large-scale clustering algorithm that can gather related words.

In this paper, we want to investigate the advantage of user behaviors and re-ranking framework in related words retrieval task using Chinese input method user records. We construct a User-Word bipartite graph to represent the information hiding in user records. The bipartite graph keeps users on one side and words on the other side. The underlying idea is that the more frequently two words co-occur in user records, the more related they are. For example, “机器翻译” (Machine Translation) is quite related to “中

文分词” (Chinese Word Segmentation) because the two words are usually used together by researchers in natural language processing community. As a result, user behaviors offer a new perspective for measuring relatedness between words. On the other hand, we can also recommend related words to users in order to enhance user experiences. Researchers are always willing to accept related terminologies in their research fields.

However, the method is purely statistics based if we only consider co-occurrence aspect. We want to add semantic features. Sahami and Helman (2006) utilize search engine to supply web queries with more semantic context and gains better results for query suggestion task. We borrow their idea in this paper. User behaviors provide statistic information to generate candidate words. Then, we can enrich candidate words with additional semantic features using search engine to retrieve more relevant candidates earlier. Statistical and semantic features can complement each other. Therefore, we can gain better performance if we consider them together.

The contributions of this paper are threefold. First, we introduce user behaviors in related word retrieval task and construct a User-Word bipartite graph from user behaviors. Words are used by users, and it is reasonable to measure relatedness between words by analyzing user behaviors. Second, we take the advantage of semantic features using search engine to reorder candidate words. We aim to return more relevant candidates earlier. Finally, our method is unsupervised and language independent, which means that we do not require any training set or manual labeling efforts.

The rest of the paper is organized as follows. Some related works are discussed in Section 2. Then we introduce our method for related words retrieval in Section 3. Experiment results and discussions are showed in Section 4. Finally, Section 5 concludes the whole paper and gives some future works.

2 Related Work

For related words retrieval task, Google Sets (Google, 2010) provides a remarkably interesting tool for finding words related to an input word. As stated in (Zheng et al., 2009), Google Sets performs poor results for input words in Chinese language. Bayesian Sets (Ghahramani and Heller, 2006) offers an alternative method for related words retrieval under the framework of Bayesian inference. It computes a score for each candidate word by comparing the posterior probability of that word given the input, to the prior probability of that candidate word. Then, it returns a ranked list of candidate words according to their computed scores.

Recently, Zheng et al. (2009) introduce user behaviors in new word detection task via a collaborative filtering manner. They extend their method to related word retrieval task. Moreover, they prove that user behaviors provide a new point for new word detection and related word retrieval tasks. However, their method is purely statistical method without considering semantic features.

We can regard related word retrieval task as problem of measuring the semantic relatedness between pairs of very short texts. Sahami and Helman (2006) introduce a web kernel function for measuring semantic similarities using snippets of search results. This work is followed by Metzler et al., (2007), Yih and Meek, (2007). They combine the web kernel with other metrics of similarity between word vectors, such as Jaccard Coefficient and KL Divergence to enhance the result.

In this paper, we follow the similar idea of using search engine to enrich semantic features of a query word. We regard the returned snippets as the context of a query word. And then we reorder candidate words and expect more relevant candidate words can be retrieved earlier. More details are given in Section 3.

3 Related Words Retrieval

In this section, we will introduce how to find related words from a single seed word via user behaviors and re-ranking framework.

First, we introduce the dataset utilized in this paper. All the resource used in this paper comes from Sogou Chinese pinyin input method (Sogou, 2006). We use Sogou for abbreviation hereafter. Users can install Sogou on their computers and the word lists they have used are kept in their user records. Volunteers are encouraged to upl-

oad their anonymous user records to the server side. In order to preserve user privacy, usernames are hidden using MD5 hash algorithm.

Then we demonstrate how to build a User-Word bipartite graph based on the dataset. The construction can be accomplished while traversing the dataset with linear time cost. We will give more details in Section 3.1.

Second, we adopt conditional probability (Deshpande and Karypis, 2004) to measure the relatedness of two words. Intuitively, two words are supposed to be related if there are a lot of users who have used both of them. In other words, the two words always co-occur in user records. Starting from a single seed word, we can generate a set of candidate words. This is the candidate generation step.

Third, in order to take the advantage of semantic features, we carry out feature extraction techniques to represent generated candidate words with enriched semantic context. In this paper, we generally make use of search engine to conduct the feature extraction step. After this step, input seed word and candidate words are represented as feature vectors in the vector space.

Finally, we can reorder generated candidate words according to their semantic relatedness of the input seed word. We expect to retrieve more relevant candidate words earlier. We will make further explanations about the mentioned steps in the next subsections.

3.1 Bipartite Graph Construction

As stated before, we first construct a User-Word bipartite graph from the dataset. The bipartite graph has two layers, with users on one side and the words on the other side. We traverse the user records, and add a link between user u and word w if w appears in the user record of u . Thus this procedure can be accomplished in linear time.

In order to give better explanations of bipartite graph construction step, we show some user records in Figure 1 and the corresponding bipartite graph in Figure 2.

User ₁	Word ₁ 自然语言(Natural Language) Word ₂ 人工智能(Artificial Intelligence)
User ₂	Word ₃ 机器翻译(Machine Translation) Word ₂ 人工智能(Artificial Intelligence)
User ₃	Word ₄ 信息检索(Information Retrieval) Word ₃ 机器翻译(Machine Translation) Word ₁ 自然语言(Natural Language)

Fig. 1. User Records Sample

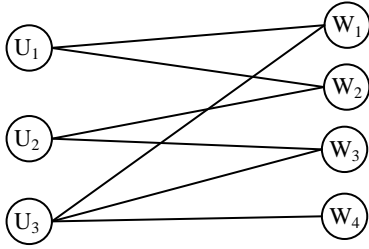


Fig. 2. Corresponding Bipartite Graph

From Figure 1, we can see that $Word_1$ and $Word_2$ appear in $User_1$'s record, which indicates that $User_1$ has used $Word_1$ and $Word_2$. As a result, in Figure 2, node $User_1$ is linked with node $Word_1$ and $Word_2$. The rest can be done in the same manner.

3.2 Candidates Generation

After the construction of bipartite graph, we can measure the relatedness of words from the bipartite graph. Intuitively, if two words always co-occur in user records, they are related to each other. Inspired by (Deshpande and Karypis, 2004), we adopt conditional probability to measure the relatedness of two words.

In particular, the conditional probability of word j occurs given that word i has already appeared is the number of users that used both word i and word j divided by the total number of users that used word i .

$$P(j|i) = \frac{Freq(ij)}{Freq(i)} \quad (1)$$

In formula 1, $Freq(X)$ is the number of users that have used words in the set X . We can clearly see that $P(j|i) \neq P(i|j)$, which means that conditional probability leads to asymmetric relations. The disadvantage is that each word i tends to have a close relationship with stop words that are used quite frequently in user records, such as “的” (of) and “一个” (a).

In order to alleviate this problem, we consider the conditional probabilities $P(j|i)$ and $P(i|j)$ together. Word i and word j is said to be quite related if conditional probabilities $P(j|i)$ and $P(i|j)$ are both relatively high. We borrow the idea proposed in (Li and Sun, 2007). In their paper, a *weighted harmonic averaging* is used to define the relatedness score between word i and word j because either $P(j|i)$ or $P(i|j)$ being too small is a severe detriment.

$$Score(i, j) = \left(\frac{\lambda}{P(i|j)} + \frac{1-\lambda}{P(j|i)} \right)^{-1} \quad (2)$$

In formula 2, parameter $\lambda \in [0, 1]$ is the weight for $P(i|j)$, which denotes how much $P(i|j)$ should be emphasized. We carry out some comparative experiments when parameter λ varies from 0 to 1 stepped by 0.1. We also tried other co-occurrence based measures like mutual information, Euclidean and Jaccard distance, and found that weight harmonic averaging gives relatively better results. Due to space limitation, we are not able to report detailed results.

So far, we have introduced how to calculate the relatedness $Score(i, j)$ between word i and word j . When a user enters an input seed word w , we can compute $Score(w, c)$ between seed word w and each candidate word c , and then sort candidate words in a descending order. Top N candidate words are kept for re-ranking, we aim to reorder top N candidate words and return the more related candidate words earlier. Alternatively, we can also set a threshold for $Score(w, c)$, which keeps the candidate word c with $Score(w, c)$ larger than the threshold. We argue that this threshold is difficult to set because different seed words have different score thresholds.

Note that this candidate generation step is completely statistical method as we only consider the co-occurrence of words. We argue that semantic features can be a complement of statistical method.

3.3 Semantic Feature Representation and Re-ranking

As stated before, we utilize search engine to enrich semantic features of the input seed word and top N candidate words. To be more specific, we issue a word to a search engine (Sogou, 2004) and get top 20 returned snippets. We regard snippets as the context and the semantic representation of this word.

For an input seed word w , we can generate top N candidate words using formula (2). We issue each word to search engine and get returned snippets. Then, each word is represented as a feature vector using bag-of-words model. Following the conventional approach, we calculate the relatedness between the input seed word w and a candidate word c as the cosine similarity between their feature vectors. Intuitively, if we introduce more candidate words, we are more likely to find related words in the candidate sets. However, noisy words are inevitably included. We will show how to tune parameter N in the experiment part.

As a result, candidate words with higher semantic similarities can be returned earlier with enriched semantic features. Re-ranking can be regarded as a complementary step after candidate generation. We can improve the performance of related word retrieval task if we consider user behaviors and re-ranking together.

4 Experiment

In this section, we demonstrate our experiment results. First, we introduce the dataset used in this paper and some statistics of the dataset. Then, we build our ground truth for related word retrieval task using Baidu encyclopedia. Third, we give some example of related word retrieval task. We show that more related words can be returned earlier if we consider semantic features. Finally, we make further analysis of the parameter tuning mentioned before.

4.1 Experiment Settings

We carry out our experiment on Sogou Chinese input method dataset. The dataset contains 10,000 users and 183,870 words, and the number of edges in the constructed bipartite graph is 42,250,718. As we can see, the dataset is quite sparse, because most of the users tend to use only a small number of words.

For related word retrieval task, we need to judge whether a candidate word is related to the input seed word. We can ask domain experts to answer this question. However, it needs a lot of manual efforts. To alleviate this problem, we adopt Baidu encyclopedia (Baidu, 2006) as our ground truth. In Baidu encyclopedia, volunteers give a set of words that are related to the particular seed word. As related words are provided by human, we are confident enough to use them as our ground truth.

We randomly select 2,000 seed words as our validation set. However, whether two words are related is quite subjective. In this paper, Baidu encyclopedia is only used as a relatively accurate standard for evaluation. We just want to investigate whether user behaviors and re-ranking framework is helpful in the related word retrieval task under various evaluation metrics.

We give a simple example of our method in Table 1. The input seed word is “机器学习” (Machine Learning). Generally speaking, all these returned candidate words are relevant to the seed word to certain degree, which indicates the effectiveness of our method.

特征向量(feature vector)	核函数(kernel function)
训练集(training set)	决策树(decision tree)
分类器(classifier)	测试集(test set)
降维(dimension reduction)	特征提取(feature extraction)

Table 1. Words Related to “Machine Learning”

4.2 Evaluation Metrics

In this paper, we use three evaluation metrics to validate the performance of our method:

1. Precision@N (**P@N**). P@N measures how much percent of the topmost results returned are correct. We consider P@5 and P@10.
2. Binary preference measure (**Bpref**) (Buckley and Voorhees, 2004). As we cannot list all the related words of an input seed word, we use Bpref to evaluate our method. For an input seed word with R judged candidate words where r is a related word and n is a nonrelated word. Bpref is defined as follow:

$$Bpref = \frac{1}{R} \sum_r 1 - \frac{|n \text{ ranked higher than } r|}{R} \quad (3)$$

3. Mean reciprocal rank of the first retrieved result (**MRR**). For a sample of input seed words W , $rank_i$ is the rank of the first related candidate word for the input seed word w_i , MRR is the average of the reciprocal ranks of results, which is defined as follow:

$$MRR = \frac{1}{|W|} \sum_i \frac{1}{rank_i} \quad (4)$$

4.3 Candidate Re-ranking

In order to show the effectiveness of semantic features and re-ranking framework, we give an example in Table 2. The input seed word is “爱立信” (Ericsson), and if we only take user behaviors into consideration, top 5 words returned are shown on the left side. After using search engine and semantic representation, we reorder the candidate words as shown on the right side.

Input Seed Word: 爱立信 (Ericsson)	
Top 5 Candidates	After Re-ranking
北电 (Nortel)	索尼爱立信 (Sony Ericsson)
中兴 (ZTE Corporation)	索爱 (Sony Ericsson)
基站 (Base Station)	阿尔卡特 (Alcatel)
阿尔卡特 (Alcatel)	索尼 (Sony)
核心网 (Core Network)	华为 (Huawei)

Table 2. Candidate Re-ranking

As shown in Table 2, we can clearly see that we return the most related candidate words such as “索尼爱立信” (Sony Ericsson) and “索爱” (the abbreviation of Sony Ericsson in Chinese) in the first two places. Moreover, after re-ranking, top candidate words are some famous brands that are quite related to query word “爱立信” (Ericsson). Some words like “核心网” (Core Network) that are not quite related to the query word are removed from the top list. From this observation, we can see that semantic features and re-ranking framework can improve the performance.

4.4 Parameter Tuning

As discussed in Section 3, we have introduced two parameters in this paper. The first is the parameter λ in the candidate generation step, and the other is the parameter N in the re-ranking step. We show how these two parameters affect the performance. In addition, we should emphasize that the ground truth is not a complete answer, so all the results are only useful for comparisons. The absolute value is not very meaningful.

As we have shown in Section 3.2, parameter λ adjusts the weight of conditional probability between two word i, j . The parameter λ is varied from 0 to 1 stepped by 0.1. We record the corresponding values of P@5, P@10, Bpref and MRR. The results are shown in Figure 3.

We can clearly see that all the values increase when λ increases first. And then all the values decrease dramatically when λ is close to 1. This indicates that either $P(j|i)$ or $P(i|j)$ being too small is a severe detriment. The result reaches peak value when $\lambda=0.5$, i.e. we should treat $P(j|i)$ and $P(i|j)$ equally to get the best result. Therefore, we use $\lambda=0.5$ to generate candidate words, those candidates are used for re-ranking.

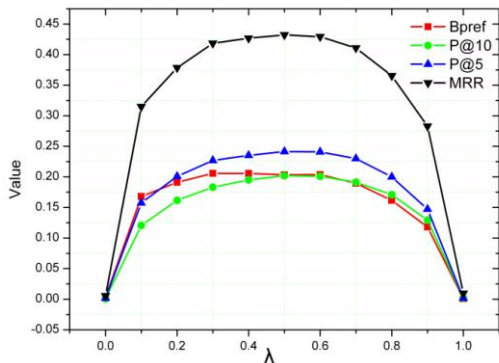


Fig. 3. Parameter λ for Candidate Generation

We also carry out the comparisons with Bayesian Sets, which is shown in Table 3. It is clear

that our method gains better results than Bayesian Sets with different values of parameter λ . Results of Google Sets are omitted here because Zheng et al. (2009) have already showed that Google Sets performs worse than Bayesian Sets with query words in Chinese.

	Bpref	MRR	P@5	P@10
$\lambda = 0.4$	0.2057	0.4267	0.2352	0.195
$\lambda = 0.5$	0.2035	0.4322	0.2414	0.2019
$\lambda = 0.6$	0.2038	0.4292	0.2408	0.2009
Bayesian Sets	0.2033	0.3291	0.1842	0.1512

Table 3. Comparisons with Bayesian Sets

To investigate the effectiveness of re-ranking framework, we also conduct experiments on the parameter N that is used for re-ranking. The experimental results are shown in Figure 4.

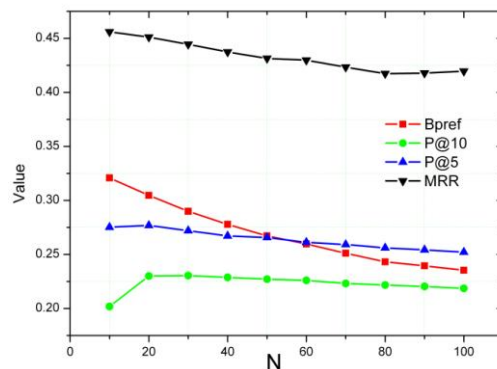


Fig. 4. Top N Candidates for Re-ranking

We can observe that more candidates tend to harm the performance as noisy words are introduced inevitably. For example, Bpref drops to less than 0.25 when $N = 100$. More comparative results are shown in Table 4. We can see that $N = 20$ gives relatively best results, which indicates that we should select Top 20 candidate words for re-ranking.

	Bpref	MRR	P@5	P@10
Non Re-ranking	0.2035	0.4322	0.2414	0.2019
$N = 10$	0.3208	0.456	0.2752	0.2019
$N = 20$	0.3047	0.4511	0.2769	0.2301
$N = 30$	0.2899	0.4444	0.272	0.2305

Table 4. Comparisons with Re-ranking Method

5 Conclusions and Future Work

In this paper, we have proposed a novel method for related word retrieval task. Different from other method, we consider user behaviors, semantic features and re-ranking framework together. We make a reasonable assumption that if two words always co-occur in user records, then

they tend to have a close relationship with each other. Based on this assumption, we first generate a set of candidate words that are related to an input seed word via user behaviors. Second, we utilize search engine to enrich candidates with semantic features. Finally, we can reorder the candidate words to return more related candidates earlier. Experiment results show that our method is effective and gains better results.

However, we also observed some noisy words in the returned results. As our dataset is generated from Chinese input method, users can type whatever they want, which will bring some noise in the dataset. We plan to remove noisy words in the future. Furthermore, we want to take the advantage of learning to rank literature (Liu, 2009) to further improve the performance of related word retrieval task. We may need to extract more features to represent the word pairs and build a labeled training set. Then various machine learning techniques can be used in this task.

Another important issue is how to build a complete and accurate ground truth for related word retrieval task. People may have different opinions about whether two words are related or not, which makes this problem complicate.

Thirdly, our method can only process a single seed word, so we aim to extend our method to process multiple seed words. In addition, we want to build a network of Chinese word association. We can discover how words are organized and connected within this network. And this word association network will be quite useful for foreigners to learn Chinese.

Fourthly, how to deal with ambiguous query word is also left as our future work. For example, query word “apple” can refer to a kind of fruit or an IT company. As a result, we are expected to return two groups of related words instead of mixing them together.

Finally, our dataset provides a new perspective for many interesting research tasks like new word detection, social network analysis, user behavior analysis, and so on. We are trying to release our dataset for research use in the future.

Acknowledgement

We thank Xiance Si and Wufeng Ke for providing the Baidu encyclopedia corpus for evaluation. We also thank the anonymous reviewers for their helpful comments and suggestions. This work is supported by a Tsinghua-Sogou joint research project.

References

- Baidu. 2006. Baidu Encyclopedia. Available at <http://baike.baidu.com>
- Chris Buckley and Ellen M. Voorhees. 2004. Retrieval Evaluation with Incomplete Information. In *Proceedings of the 27th annual international ACM SIGIR conference on Research and development in information retrieval*, pp 25-32
- Mukund Deshpande and George Karypis. 2004. Item-Based Top-N Recommendation Algorithms, *ACM Trans. Information Systems*, 22(1): 143-177
- Zoubin Ghahramani and Katherine A. Heller. 2005. Bayesian Sets. In *Advances in Neural Information Processing Systems*
- Google. Google Sets. Accessed on Feb. 9th, 2010, available at: <http://labs.google.com/sets>
- Jingyang Li and Maosong Sun. 2007. Scalable term selection for text categorization, In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pp. 774-782
- Tie-Yan Liu. 2009. Learning to Rank for Information Retrieval, *Foundation and Trends on Information Retrieval*, Now Publishers
- Donald Metzler, Susan T. Dumais, and Christopher Meek. 2007. Similarity measures for short segments of text. In *Proceeding of the 29th European Conference on Information Retrieval*, pp 16-27
- Mehran Sahami and Timothy D. Heilman. 2006. A web-based kernel function for measuring the similarity of short text snippets. In *Proceedings of the 15th International Conference on World Wide Web*, pp 377-386
- Sogou. 2006. Sogou Chinese Pinyin Input Method. Available at <http://pinyin.sogou.com/>
- Sogou. 2004. Sogou Search Engine. Available at <http://www.sogou.com>
- Wen-Tau Yih and Christopher Meek. 2007. Improving similarity measures for short segments of text. In *Proceedings of AAAI 2007*, pp 1489-1494
- Yabin Zheng, Zhiyuan Liu, Maosong Sun, Liyun Ru, and Yang Zhang. 2009. Incorporating User Behaviors in New Word Detection. In *Proceedings of the Twenty-First International Joint Conference on Artificial Intelligence*, pp 2101-2106

Transition-based parsing with Confidence-Weighted Classification

Martin Haulrich

Dept. of International Language Studies and Computational Linguistics
Copenhagen Business School

mwh.isv@cbs.dk

Abstract

We show that using confidence-weighted classification in transition-based parsing gives results comparable to using SVMs with faster training and parsing time. We also compare with other online learning algorithms and investigate the effect of pruning features when using confidence-weighted classification.

1 Introduction

There has been a lot of work on data-driven dependency parsing. The two dominating approaches have been graph-based parsing, e.g. MST-parsing (McDonald et al., 2005b) and transition-based parsing, e.g. the MaltParser (Nivre et al., 2006a). These two approaches differ radically but have in common that the best results have been obtained using margin-based machine learning approaches. For the MST-parsing MIRA (McDonald et al., 2005a; McDonald and Pereira, 2006) and for transition-based parsing Support-Vector Machines (Hall et al., 2006; Nivre et al., 2006b).

Dredze et al. (2008) introduce a new approach to margin-based online learning called *confidence-weighted classification* (CW) and show that the performance of this approach is comparable to that of Support-Vector Machines. In this work we use confidence-weighted classification with transition-based parsing and show that this leads to results comparable to the state-of-the-art results obtained using SVMs.

We also compare training time and the effect of pruning when using confidence-weighted learning.

2 Transition-based parsing

Transition-based parsing builds on the idea that parsing can be viewed as a sequence of transitions

between states. A transition-based parser (deterministic classifier-based parser) consists of three essential components (Nivre, 2008):

1. A parsing algorithm
2. A feature model
3. A classifier

The focus here is on the classifier but we will briefly describe the parsing algorithm in order to understand the classification task better.

The parsing algorithm consists of two components, a *transition system* and an *oracle*. Nivre (2008) defines a transition system $S = (C, T, c_s, C_t)$ in the following way:

1. C is a set of configurations, each of which contains a buffer β of (remaining) nodes and a set A of dependency arcs,
2. T is a set of transitions, each of which is a partial function $t : C \rightarrow C$,
3. c_s is a initialization function mapping a sentence $x = (w_0, w_1, \dots, w_n)$ to a configuration with $\beta = [1, \dots, n]$,
4. C_t is a set of terminal configurations.

A *transition sequence* for a sentence x in S is a sequence $C_{0,m} = (c_0, c_1, \dots, c_m)$ of configurations, such that

1. $c_0 = c_s(x)$,
2. $c_m \in C_t$,
3. for every i ($1 \leq i \leq m$) $c_i = t(c_{i-1})$ for some $t \in T$

The oracle is used during training to determine a transition sequence that leads to the correct parse. The job of the classifier is to 'imitate' the oracle, i.e. to try to always pick the transitions that

lead to the correct parse. The information given to the classifier is the current configuration. Therefore the training data for the classifier consists of a number of configurations and the transitions the oracle chose with these configurations.

Here we focus on stack-based parsing algorithms. A stack-based configuration for a sentence $x = (w_0, w_1, \dots, w_n)$ is a triple $c = (\sigma, \beta, A)$, where

1. σ is a stack of tokens $i \leq k$ (for some $k \leq n$),
2. β is a buffer of tokens $j > k$,
3. A is a set of dependency arcs such that $G = (0, 1, \dots, n, A)$ is a dependency graph for x . (Nivre, 2008)

In the work presented here we use the NivreEager algorithm which has four transitions:

Shift Push the token at the head of the buffer onto the stack.

Reduce Pop the token on the top of the stack.

Left-Arc_l Add to the analysis an arc with label l from the token at the head of the buffer to the token on the top of the stack, and push the buffer-token onto the stack.

Right-Arc_l Add to the analysis an arc with label l from the token on the top of the stack to the token at the head of the buffer, and pop the stack.

2.1 Classification

Transition-based dependency parsing reduces parsing to consecutive multiclass classification. From each configuration one amongst some predefined number of transitions has to be chosen. This means that any classifier can be plugged into the system. The training instances are created by the oracle so the training is offline. So even though we use online learners in the experiments these are used in a batch setting.

The best results have been achieved using Support-Vector Machines placing the MaltParser very high in both the CoNLL shared tasks on dependency parsing in 2006 and 2007 (Buchholz and Marsi, 2006; Nivre et al., 2007) and it has been shown that SVMs are better for the task than Memory-based learning (Hall et al., 2006). The standard setting in the MaltParser is to use a 2nd-degree polynomial kernel with the SVM.

3 Confidence-weighted classification

Dredze et al. (2008) introduce confidence-weighted linear classifiers which are online-classifiers that maintain a confidence parameter for each weight and uses this to control how to change the weights in each update. A problem with online algorithms is that because they have no memory of previously seen examples they do not know if a given weight has been updated many times or few times. If a weight has been updated many times the current estimation of the weight is probably relatively good and therefore should not be changed too much. On the other hand if it has never been updated before the estimation is probably very bad. CW classification deals with this by having a confidence-parameter for each weight, modeled by a Gaussian distribution, and this parameter is used to make more aggressive updates on weights with lower confidence (Dredze et al., 2008). The classifiers also use Passive-Aggressive updates (Crammer et al., 2006) to try to maximize the margin between positive and negative training instances.

CW classifiers are online-algorithms and are therefore fast to train, and it is not necessary to keep all training examples in memory. Despite this they perform as well or better than SVMs (Dredze et al., 2008). Crammer et al. (2009) extend the approach to multiclass classification and show that also in this setting the classifiers often outperform SVMs. They show that updating only the weights of the best of the wrongly classified classes yields the best results. We also use this approach, called top-1, here.

Crammer et al. (2008) present different update-rules for CW classification and show that the ones based on standard deviation rather than variance yield the best results. Our experiments have confirmed this, so in all experiments the update-rule from equation 10 (Crammer et al., 2008) is used.

4 Experiments

4.1 Software

We use the open-source parser MaltParser¹ for all experiments. We have integrated confidence-weighted, perceptron and MIRA classifiers into the code. The code for the online classifiers has

¹We have used version 1.3.1, available at maltparser.org

been made available by the authors of the CW-papers.

4.2 Data

We have used the 10 smallest data sets from CoNNL-X (Buchholz and Marsi, 2006) in our experiments. Evaluation has been done with the official evaluation script and evaluation data from this task.

4.3 Features

The standard setting for the MaltParser is to use SVMs with polynomial kernels, and because of this it uses a relatively small number of features. In most of our experiments the default feature set of MaltParser consisting of 14 features has been used.

When using a linear-classifier without a kernel we need to extend the feature set in order to achieve good results. We have done this very uncritically by adding all pair wise combinations of all features. This leads to 91 additional features when using the standard 14 features.

5 Results and discussion

We will now discuss various results of our experiments with using CW-classifiers in transition-based parsing.

5.1 Online classifiers

We compare CW-classifiers with other online algorithms for linear classification. We compare with perceptron (Rosenblatt, 1958) and MIRA (Crammer et al., 2006). With both these classifiers we use the same top-1 approach as with the CW-classifiers and also averaging which has been shown to alleviate overfitting (Collins, 2002). Table 2 shows Labeled Attachment Score obtained with the three online classifiers. All classifiers were trained with 10 iterations.

These results confirm those by Crammer et al. (2009) and show that confidence-weighted classifiers are better than both perceptron and MIRA.

5.2 Training and parsing time

The training time of the CW-classifiers depends on the number of iterations used, and this of course affects the accuracy of the parser. Figure 1 shows Labeled Attachment Score as a function of the number of iterations used in training. The horizontal line shows the LAS obtained with SVM.

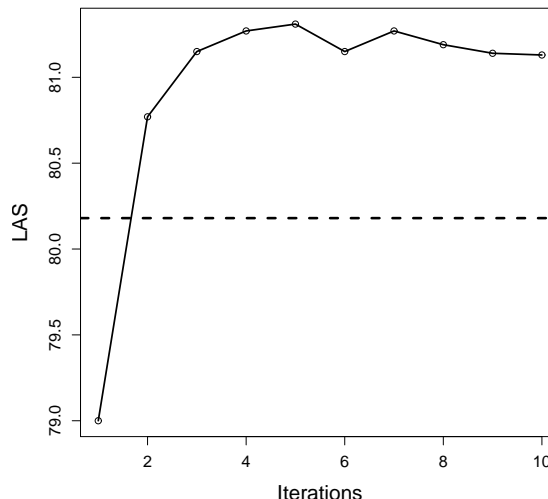


Figure 1: LAS as a function of number of training iterations on Danish data. The dotted horizontal line shows the performance of the parser trained with SVM.

We see that after 4 iterations the CW-classifier has the best performance for the data set (Danish) used in this experiment. In most experiments we have used 10 iterations. Table 1 compares training time (10 iterations) and parsing time of a parser using a CW-classifiers and a parser using SVM on the same data set. We see that training of the CW-classifier is faster, which is to be expected given their online-nature. We also see that parsing is much faster.

	SVM	CW
Training	75 min	8 min
Parsing	29 min	1.5 min

Table 1: Training and parsing time on Danish data.

5.3 Pruning features

Because we explicitly represent pair wise combinations of all of the original features we get an extremely high number of binary features. For some of the larger data sets, the number of features is so big that we cannot hold the weight-vector in memory. For instance the Czech data-set has 16 million binary features, and almost 800 classes - which means that in practice there are 12 billion binary features².

²Which is also why we only have used the 10 smallest data sets from CoNNL-X.

	Perceptron	MIRA	CW, manual fs	CW	SVM
Arabic	58.03	59.19	60.55	† 60.57	59.93
Bulgarian	80.46	81.09	82.57	† 82.76	82.12
Danish	79.42	79.90	81.06	† 81.13	80.18
Dutch	75.75	77.47	77.65	† 78.65	77.76
Japanese	87.74	88.06	88.14	88.19	† 89.47
Portuguese	85.69	85.95	86.11	86.20	86.25
Slovene	64.35	65.38	66.09	† 66.28	65.45
Spanish	74.06	74.86	75.58	75.90	75.46
Swedish	79.79	80.31	81.03	† 81.24	80.56
Turkish	46.48	47.13	46.98	47.09	47.49
All	78.26	79.00	79.68	† 79.86	79.59

Table 2: LAS on development data for three online classifiers, CW-classifiers with manual feature selection and SVM. Statistical significance is measured between CW-classifiers without feature selection and SVMs.

To solve this problem we have tried to use pruning to remove the features occurring fewest times in the training data. If a feature occurs fewer times than a given cutoff limit the feature is not included. This goes against the idea of CW classifiers which are exactly developed so that rare features can be used. Experiments also show that this pruning hurts accuracy. Figure 2 shows the labeled attachment score as a function of the cutoff limit on the Danish data.

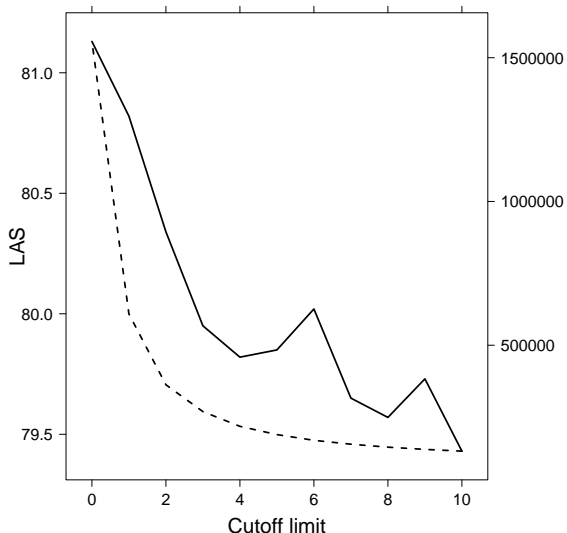


Figure 2: LAS as a function of the cutoff limit when pruning rare features. The dotted line shows the number of features left after pruning.

5.4 Manual feature selection

Instead of pruning the features we tried manually removing some of the pair wise feature combinations. We removed some of the combinations that lead to the most extra features, which is especially the case with combinations of lexical features. In the extended default feature set for instance we removed all combinations of lexical features except the combination of the word form of the token at the top of the stack and of the word form of the token at the head of the buffer.

Table 2 shows that this consistently leads to a small decreases in LAS.

5.5 Results without optimization

Table 2 shows the results for the 10 CoNNL-X data sets used. For comparison we have included the results from using the standard classifier in the MaltParser, i.e. SVM with a polynomial kernel. The hyper-parameters for the SVM have not been optimized, and neither has the number of iterations for the CW-classifiers, which is always 10. We see that in many cases the CW-classifier does significantly³ better than the SVM, but that the opposite is also the case.

5.6 Results with optimization

The results presented above are suboptimal for the SVMs because default parameters have been used for these, and optimizing these can improve ac-

³In all tables statistical significance is marked with †. Significance is calculated using McNemar’s test ($p = 0.05$). These tests were made with MaltEval (Nilsson and Nivre, 2008)

	SVM			CW		
	LAS	UAS	LA	LAS	UAS	LA
Arabic	66.71	77.52	80.34	67.03	77.52	† 81.20
Bulgarian*	87.41	91.72	90.44	87.25	91.56	89.77
Danish	† 84.77	† 89.80	89.16	84.15	88.98	88.74
Dutch*	† 78.59	† 81.35	† 83.69	77.21	80.21	82.63
Japanese	† 91.65	† 93.10	† 94.34	90.41	91.96	93.34
Portuguese*	† 87.60	† 91.22	† 91.54	86.66	90.58	90.34
Slovene	70.30	78.72	80.54	69.84	† 79.62	79.42
Spanish	81.29	84.67	90.06	82.09	† 85.55	90.52
Swedish*	† 84.58	89.50	87.39	83.69	89.11	87.01
Turkish	† 65.68	† 75.82	† 78.49	62.00	73.15	76.12
All	† 79.86	† 85.35	† 86.60	79.04	84.83	85.91

Table 3: Results on the CoNNL-X evaluation data. Manual feature selection has been used for languages marked with an *.

curacy a lot. In this section we will compare results obtained with CW-classifiers with the results for the MaltParser from CoNNL-X. In CoNNL-X both the hyper parameters for the SVMs and the features have been optimized. Here we do not do feature selection but use the features used by the MaltParser in CoNNL-X⁴.

The only hyper parameter for CW classification is the number of iterations. We optimize this by doing 5-fold cross-validation on the training data. Although the manual feature selection has been shown to decrease accuracy this has been used for some languages to reduce the size of the model. The results are presented in table 3.

We see that even though the feature set used are optimized for the SVMs there are not big differences between the parses that use SVMs and the parsers that use CW classification. In general though the parsers with SVMs does better than the parsers with CW classifiers and the difference seems to be biggest on the languages where we did manual feature selection.

6 Conclusion

We have shown that using confidence-weighted classifiers with transition-based dependency parsing yields results comparable with the state-of-the-art results achieved with Support Vector Machines - with faster training and parsing times. Currently we need a very high number of features to achieve these results, and we have shown that pruning this big feature set uncritically hurts performance of

⁴Available at <http://maltparser.org/conll/conllx/>

the confidence-weighted classifiers.

7 Future work

Currently the biggest challenge in the approach outlined here is the very high number of features needed to achieve good results. A possible solution is to use kernels with confidence-weighted classification in the same way they are used with the SVMs.

Another possibility is to extend the feature set in a more critical way than what is done now. For instance the combination of a POS-tag and CPOS-tag for a given word is now included. This feature does not convey any information that the POS-tag-feature itself does not. The same is the case for some word-form and word-lemma features. All in all a lot of non-informative features are added as things are now. We have not yet tried to use automatic features selection to select only the combinations that increase accuracy.

We will also try to do feature selection on a more general level as this can boost accuracy a lot. The results in table 3 are obtained with the features optimized for the SVMs. These are not necessarily the optimal features for the CW-classifiers.

Another comparison we would like to do is with linear SVMs. Unlike the polynomial kernel SVMs used as default in the MaltParser linear SVMs can be trained in linear time (Joachims, 2006). Trying to use the same extended feature set we use with the CW-classifiers with a linear SVM would provide an interesting comparison.

8 Acknowledgements

The author thanks three anonymous reviewers and Anders Søgaard for their helpful comments and the authors of the CW-papers for making their code available.

References

- Sabine Buchholz and Erwin Marsi. 2006. Conll-x shared task on multilingual dependency parsing. In *Proceedings of the Tenth Conference on Computational Natural Language Learning (CoNLL-X)*, pages 149–164, New York City, June. Association for Computational Linguistics.
- Michael Collins. 2002. Discriminative training methods for hidden markov models: theory and experiments with perceptron algorithms. In *EMNLP '02: Proceedings of the ACL-02 conference on Empirical methods in natural language processing*, pages 1–8, Morristown, NJ, USA. Association for Computational Linguistics.
- Koby Crammer, Ofer Dekel, Joseph Keshet, Shai Shalev-Shwartz, and Yoram Singer. 2006. Online passive-aggressive algorithms. *J. Mach. Learn. Res.*, 7:551–585.
- Koby Crammer, Mark Dredze, and Fernando Pereira. 2008. Exact convex confidence-weighted learning. In Daphne Koller, Dale Schuurmans, Yoshua Bengio, and Léon Bottou, editors, *NIPS*, pages 345–352. MIT Press.
- Koby Crammer, Mark Dredze, and Alex Kulesza. 2009. Multi-class confidence weighted algorithms. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing*, pages 496–504, Singapore, August. Association for Computational Linguistics.
- Mark Dredze, Koby Crammer, and Fernando Pereira. 2008. Confidence-weighted linear classification. In *ICML '08: Proceedings of the 25th international conference on Machine learning*, pages 264–271, New York, NY, USA. ACM.
- Johan Hall, Joakim Nivre, and Jens Nilsson. 2006. Discriminative classifiers for deterministic dependency parsing. In *Proceedings of the COLING/ACL 2006 Main Conference Poster Sessions*, pages 316–323, Sydney, Australia, July. Association for Computational Linguistics.
- Thorsten Joachims. 2006. Training linear svms in linear time. In *KDD '06: Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 217–226, New York, NY, USA. ACM.
- Ryan T. McDonald and Fernando C. N. Pereira. 2006. Online learning of approximate dependency parsing algorithms. In *EACL*. The Association for Computer Linguistics.
- Ryan T. McDonald, Koby Crammer, and Fernando C. N. Pereira. 2005a. Online large-margin training of dependency parsers. In *ACL*. The Association for Computer Linguistics.
- Ryan T. McDonald, Fernando Pereira, Kiril Ribarov, and Jan Hajic. 2005b. Non-projective dependency parsing using spanning tree algorithms. In *HLT/EMNLP*. The Association for Computational Linguistics.
- Jens Nilsson and Joakim Nivre. 2008. Malteval: An evaluation and visualization tool for dependency parsing. In *Proceedings of the Sixth International Language Resources and Evaluation*, Marrakech, Morocco, May. LREC.
- Joakim Nivre, Johan Hall, and Jens Nilsson. 2006a. Maltparser: A data-driven parser-generator for dependency parsing. In *Proceedings of the fifth international conference on Language Resources and Evaluation (LREC2006)*, pages 2216–2219, May.
- Joakim Nivre, Johan Hall, Jens Nilsson, Gülşen Eryiğit, and Svetoslav Marinov. 2006b. Labeled pseudo-projective dependency parsing with support vector machines. In *Proceedings of the Tenth Conference on Computational Natural Language Learning (CoNLL-X)*, pages 221–225, New York City, June. Association for Computational Linguistics.
- Joakim Nivre, Johan Hall, Sandra Kübler, Ryan McDonald, Jens Nilsson, Sebastian Riedel, and Deniz Yuret. 2007. The CoNLL 2007 shared task on dependency parsing. In *Proceedings of the CoNLL Shared Task Session of EMNLP-CoNLL 2007*, pages 915–932.
- Joakim Nivre. 2008. Algorithms for deterministic incremental dependency parsing. *Computational Linguistics*, 34(4):513–553.
- Frank Rosenblatt. 1958. The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, 65(6):386–408.

Expanding Verb Coverage in Cyc With VerbNet

Clifton J. McFate

Northwestern University

Evanston, IL, USA.

c-mcfate@northwestern.edu

Abstract

A robust dictionary of semantic frames is an essential element of natural language understanding systems that use ontologies. However, creating lexical resources that accurately capture semantic representations en masse is a persistent problem. Where the sheer amount of content makes hand creation inefficient, computerized approaches often suffer from over generality and difficulty with sense disambiguation. This paper describes a semi-automatic method to create verb semantic frames in the Cyc ontology by converting the information contained in VerbNet into a Cyc usable format. This method captures the differences in meaning between types of verbs, and uses existing connections between WordNet, VerbNet, and Cyc to specify distinctions between individual verbs when available. This method provides 27,909 frames to OpenCyc which currently has none and can be used to extend ResearchCyc as well. We show that these frames lead to a 20% increase in sample sentences parsed over the Research Cyc verb lexicon.

1 Introduction

The Cyc¹ knowledge base represents general purpose knowledge across a vast array of domains. Low level event and individual facts are contained in larger definitional hierarchical representations and contextualized through microtheories (Matuszek *et al.*, 2006). Higher order predicates built into Cyc's formal language, CycL, allow efficient inferencing about context and meta-language reasoning above and beyond first-order logic rules (Ramachandran *et al.*, 2005).

Because of the expressiveness and size of the ontology, Cyc has been used in NL applications

including word sense disambiguation and rule acquisition by reading (Curtis, Cabral, & Baxter, 2006; Curtis *et al.*, 2009). Such applications use NL-to-CycL parsers which use Cyc semantic frames to convert natural language into Cyc representations. These frames represent sentence content through a set of propositional logic assertions that first reify the sentence in terms of a real world event and then define the semantic relationships between the elements of the sentence, as described later. Because these parsers require semantic frames to represent sentence content, existing parsers are limited due to Cyc's limited coverage (Curtis *et al.*, 2009). The goal is to increase this coverage by automatically translating the class frames in VerbNet into individual verb templates.

2 Previous Work

The Cyc knowledge base is continuously expanding and much work has been done on automatic fact acquisition as well as merging ontologies. However, the semantic frames remain mostly hand-made in ResearchCyc² and non-existent in the open-license OpenCyc³. Translating VerbNet frames into Cyc will expand the natural language capabilities of both.

There has been previous research on mapping existing Cyc templates to VerbNet, but thus far these approaches have not created new templates to address Cyc's lapses in coverage. One such attempt, King and Crouch's (2005) unified lexicon, compiled many lexical resources into a unified representation. While this research created a valuable resource, it did not extend the existing Cyc coverage. Of the 45,704 entries in the UL only 3,544 have Cyc entries (King & Crouch, 2005).

Correspondences between a few VerbNet frames and ResearchCyc templates have also been mapped out through the VxC VerbNet Cyc

¹ <http://www.opencyc.org/cyc>

² <http://research.cyc.com>

³ <http://opencyc.org>

Mapper (Trumbo 2006). These mappings became a standard that we later used to evaluate the quality of our created frames.

A notable exception to the hand-made paradigm is Curtis *et al*'s (2009) TextLearner which uses rules and existing semantic frames to handle novel sentence structures. Given an existing template that fits some of the syntactic constraints of the sentence, TextLearner will attempt to create a new frame by suggesting a predicate that fits the missing part. Often these are general underspecified predicates, but TextLearner is able to use common sense reasoning and existing facts to find better matches (Curtis *et al*, 2009).

While TextLearner improves its performance with time, it is not an attempt to create new frames on a large scale. Creating generalized frames based on verb classes will increase the depth of the Cyc Lexicon quickly. Furthermore, automatic processes like those in TextLearner could be used to make individual verb semantic frames more specific.

3 VerbNet

VerbNet is an extension of Levin's (1993) verb classes that uses the class structure to apply general syntactic frames to member verbs that have those syntactic uses and similar semantic meanings (Kipper *et al*, 2000). The current version has been expanded to include class distinctions not included in Levin's original proposal (Kipper *et al*, 2006).

VerbNet is an appealing lexical resource for this task because it represents semantic meaning as the union of both syntactic structure and semantic predicates. VerbNet uses Lexicalized Tree Adjoining Grammar to generate the syntactic frames. The syntactic roles in the frame are appended with general thematic roles that fill arguments of semantic predicates. Each event is broken down into a tripartite structure as described by Moens & Steedman (1988) and uses a time modifier for each predicate to indicate when specific predicates occur in the event. This allows for a dynamic representation of change over an event. (Kipper *et al*, 2000).

This approach is transferable to Cyc's semantic templates in which syntactic slots fill predicate arguments in the context of a specific syntactic frame. Both also have extensive connections to WordNet2.0, an electronic edition of Miller's (1985) WordNet (Fellbaum, 1998).

4 Method

The general method for creating semantic templates in Cyc requires creating Verb Class Frames and then using Cyc predicates and heuristic rules to create individual frames for each member verb.

4.1 OpenCyc

The existing semantic templates are accessible through the ResearchCyc KB. However, for the purposes of this study the OpenCyc KB was used. The OpenCyc KB is an open source version of ResearchCyc that contains much of the definitional information and higher order predicates, but has had much of the lower level specific facts and the entire word lexicon removed (Matuszek *et al*, 2006). However, the assertions generated by this method are fully usable in ResearchCyc. OpenCyc was used so as to minimize the effect of existing semantic frames on new frame creation. Since OpenCyc and VerbNet are open-licensed, our translation provides an open-license extension to OpenCyc to support its use in natural language research.

4.2 Knowledge Representation

The primary difficulty with integrating VerbNet frames into Cyc was overcoming differences in knowledge representation. Cyc semantic templates reify events as an instance of a collection of events. The arguments correspond to syntactic roles. The following is a semantic template for a ditransitive use of the word *give* from ResearchCyc.

```
(verbSemTrans Give-TheWord 0
 (PPCompFrameFn
  DitransitivePPFrameType To-TheWord)
 (and
  (isa ACTION GivingSomething)

  (objectGiven ACTION OBJECT)
  (giver ACTION SUBJECT)
  (givee ACTION OBLIQUE-OBJECT)))
```

However, VerbNet uses semantic predicates that describe relationships between two thematic roles. The following is a frame for the VerbNet class *Give* as presented in the Unified Verb Index⁴.

```
NP V NP PP.recipient
example
```

⁴ <http://verbs.colorado.edu/verb-index/>

"They lent a bicycle to me."

syntax

Agent V Theme {to} Recipient

semantics

-has_possession(start(E), Agent, Theme)
-has_possession(end(E), Recipient, Theme)
-transfer(during(E), Theme)
-cause(Agent, E)

The predicate `has_possession` occurs twice, at the beginning and end of the event. In one case the Agent has possession and in the second the Recipient does. Both refer to the Theme which is being transferred.

In Cyc the `hasPossession` relationship to Agent and Recipient is represented with the predicates `giver` and `givee`. The subject and oblique-object of the sentence fill those arguments, and the actual change of possession is represented by the collection of events `GivingSomething`. The VerbNet Theme is the object in `objectGiven`. Thus an individual VerbNet semantic predicate often has a many-to-one mapping with Cyc predicates.

4.3 Predicates

To account for representation differences, a single Cyc predicate was mapped to a unique combination of Verbnet predicate and thematic role (ie. `Has_Possession Agent at start(E) => givee`). 56 of these mappings were done by hand. Though far from exhaustive, these hand mappings represent many frequently used predicates in VerbNet. The hand mapping was done by looking at the uses of the predicate across different classes.

Because the mappings were not exhaustive, a safety net automatically catches predicates that haven't been mapped. The VerbNet predicates `Cause` and `InReactionTo` corresponded to the Cyc predicates `performedBy`, `doneBy`, and `causes-Underspecified`. These predicates were selected whenever the VerbNet predicates occurred with a theme role that was the subject of the sentence. The more specific `performedBy` was selected in cases where the frame's temporal structure suggested a result. The predicate `doneBy` was selected in other cases. The `causes-Underspecified` predicate was used in frames whose time modifiers suggested that they were continuous states. The predicates `patientGeneric` and `patientGeneric-Direct` were used when a

predicate was not found for a required object or oblique object.

Some Cyc templates don't have predicates that reference the event. For example, the verb *touch* can be efficiently represented with the relation (`objectsInContact :SUBJECT :OBJECT`). Situations like this were hand assigned.

4.4 Collections

In Cyc, concepts are represented by collections. Inheritance between collections is specified by the `genls` relationship, which can be viewed as subset. Most verb frames have an associated collection of events of which each use is an instance. The associated collection of the class frame templates was automatically selected using the common link that both resources share with WordNet (Fellbaum, 1998). To do this, the WordNet synsets of the member verbs for a class were matched with their Cyc-WordNet `synonymousExternalConcept` assertion. The Cyc representation became a denoted collection. The most general collection out of the list of viable collections was chosen as the general class frame collection. The number of `genls` links to a collection was used as a proxy for generality. In the case of a tie the first was chosen.

While the most general collection was used for the class semantic frame, at the level of individual verb frames the specific synset denoted collection was substituted for the more general one when applicable. Verbs with multiple meanings across classes were given a unique index number for each sense. However, within a given class each word only received one denotation. The general class level collection was used in cases where no Cyc-WordNet-VerbNet link existed. If no verb had a synset in Cyc, the general collection `Situation` was used.

4.5 Subcategorization Frames

Each syntactic frame is a subcategorization frame or a subset of one. In this case, the naming conventions were different between VerbNet and Cyc. Frames with prepositions kept Cyc's notation for prepositional phrases. However, since VerbNet had a much broader coverage the VerbNet subcat names were kept.

4.6 Assertions

The process above was used to create general class frames, for example,

```
(verbClassSemTrans give-13.1  
(TransitiveNPFrame)
```

```
(and
(isa :ACTION
MakingSomethingAvailable)
(patient-GenericDirect :ACTION
:OBJECT)
(performedBy :ACTION :SUBJECT)
(fromPossessor :ACTION :SUBJECT)
(objectOfPossessionTransfer :ACTION
:OBJECT)))
```

These frames use more generic collections and apply to a VerbNet class rather than a specific verb.

Specific verb semantic templates were created by inferring that each member verb of a VerbNet class participated in every template in a class. Again, collections were taken from existing WordNet connections if possible. The output was assertions in the Cyc semantic template format:

```
(verbSemTrans Loan-TheWord 0
(PPCompFrameFn NP-PP (WordFn to))
(and
(isa :ACTION Lending)
(patient-GenericDirect :ACTION
:OBJECT)
(performedBy :ACTION :SUBJECT)
(fromPossessor :ACTION :SUBJECT)
(toPossessor :ACTION :OBLIQUE-
OBJECT)
(objectOfPossessionTransfer :ACTION
:OBJECT)))
```

This method for giving class templates to each verb in a class was written as a Horn clause for the FIRE reasoning engine. FIRE is a reasoning engine that incorporates both logical inference based on axioms and analogy-based reasoning over a Cyc-derived knowledge base (Forbus, Mostek, & Ferguson, 2002). FIRE could then be queried for implied verb templates which became the final list of verb templates.

4.7 Subclasses

VerbNet has an extensive classification system involving subclasses. Subclasses contain verbs that take all of the syntactic formats of the main class plus additional frames that verbs in the main class cannot.

Verbs in a subclass inherit frames from their superordinate classes. FIRE was used again to create the verb semantic templates.

Each subclass template's collection was selected using the same process as the main class. If no subclass member had a Cyc denotation, then the main class collection was used.

5 Results

The end result of this process was the creation of 27,909 verb semantic template assertions for 5,050 different verbs. This substantially increases the number of frames for ResearchCyc and creates frames for OpenCyc.

To test the accuracy of the results and their contribution to the knowledge base we ran two tests. The first was to compare our frames with the 139 hand-checked VxC matches by hand. Of the 139 frames from VxC, 81 were qualified as "good" matches, and 58 as "maybe" (Trumbo, 2006). Since these frames already existed in Cyc and were hand matched we used them as the current gold standard for what a VerbNet frame translated into Cyc should look like.

Matches between frames were evaluated along several criteria. First was whether the frame had as good a syntactic parse as the manual version. This was defined as having predicates that addressed all syntactic roles in the sentence or, if not enough, as many as the VxC match. Secondly we asked if the collection was similar to the manual version. Frames with collections that were too specific, unrelated, or just *Situation* were discarded. Because frame-specific predicates were not created on a large scale, a frame was not rejected for using general predicates.

It is important to note a difference in matching methodology between the VxC matches and our frames. First, the VxC mappings included frames in Cyc that only partially matched more syntactically robust VerbNet frames. Our frames were only included if they matched the intended VerbNet syntactic frame. Because of this some of our frames beat the VxC gold standard for syntactic completeness. The VxC frames also included multiple similar senses for an individual verb. Our verbs had one denotation per class or subclass. Thus in some cases our frames failed not from over generalizing but because they were only meant to represent one meaning per class. Since the strength of our approach lies in generating a near exhaustive list of syntactic frames and not multiple word senses, these kinds of failures are not necessarily representative of the success of the frames as a whole.

A total of 55 frames (39.5%) were correct with seventeen (30.9%) of the correct frames having a more complete syntactic parse than the manually mapped frame. 48 frames (34.5%) were rejected only for having too general or specific a collection; however ten (20.8%) of the collection

rejected frames had a more complete parse than their manual counterparts. Thus 103 frames (74.1%) were as syntactically correct or better than the existing Cyc frame mapped to that VerbNet frame. Nine (6.47%) frames failed syntactically, with four (44.4%) of the syntax failures also having the wrong collection. Thirteen frames (9.3%) were not matched.

Fifteen frames (10.8%) from the *Hold* class, were separated out for a formatting error that resulted in a duplicate, though not syntactically incorrect, predicate. The predicate repeated was (`objectsInContact :ACTION :OBJECT`). 12 of 15 frames (80%) had accurate collections.

The second test compared the results of a natural language understanding system using either ResearchCyc alone or a version of ResearchCyc with our frames substituted for theirs. The test corpus was 50 randomly selected example sentences from the VerbNet frame examples. We used the EA NLU parser, which uses a bottom-up chart parser and compositional semantics to convert the semantic content of a sentence in CycL (Tomai & Forbus 2009). Possible frames are returned in choice sets. A parse was judged correct if it returned a verb frame for the central verb of the example sentence that either wholly or in combination with preposition frames addressed the syntactic constituents of the sentence with an acceptable collection and acceptable predicates. Again general predicates were acceptable.

ResearchCyc got sixteen out of 50 frames correct (32%). Eleven frames (22%) did not return a template but did return a denotation to a Cyc collection. Twelve verbs (24%) returned nothing, while eleven (22%) returned frames that were either not the correct syntactic frame or were a different sense of the verb.

EA NLU running the VerbNet generated frames got 26 out of 50 (52%) frames correct. Twelve frames (24%) returned nothing. Eight frames, (16%) failed because of a too specific or too general collection. Four generated frames (8%) were either not the correct syntactic frame or were for a different sense of the verb. This was an overall 20% improvement in accuracy.

Five (10%) parses using the VerbNet generated correct frames that were labeled as noisy. Noisy frames had duplicate predicates or more general predicates in addition to the specific ones. The *Hold* frames separated out in the VxC test are an example of noisy frames. None of these frames were syntactically incorrect or contradictory. The redundant predicates arise

because the predicate safety net had to be greedy. This was in the interest of capturing more complex frames that may have multiple relations for the same thematic role in a sentence.

This evaluation is based on parser recall and frame semantic accuracy only. As would be expected, adding more frames to the knowledge base did result in more parser retrievals and possible interpretations. The implications for this on word sense disambiguation is evaluated further in the discussion. To improve predicate specificity, the next phase of research with these frames will be to implement predicate strengthening methods that move down the hierarchy to find more specific predicates to replace the generalized ones. Thus in the future precision both in terms of frame retrieval and predicate specificity will be a vital metric for evaluating success.

6 Discussion

As has been demonstrated in this approach and in previous research like Curtis *et al's* (2009) TextLearner, Cyc provides powerful reasoning capabilities that can be used to successfully infer more specific information from general existing facts. We hope that future research is able to use this feature to provide more specific individual frames. Because Cyc is consistently changing and growing, an approach that uses Cyc relationships will be able to improve as the knowledge base improves its coverage.

While many of the frames are general, they provide a solid foundation for further research. As they are now, the added 27,909 frames increase the language capabilities of OpenCyc which previously had none. For ResearchCyc the contribution is less clear-cut. The 27,909 VerbNet frames have approximately 7.93 times the coverage of the existing 3,517 ResearchCyc frames⁵ and they improved ResearchCyc parser performance by 20%. However, with 35% of frames in the VxC comparison and 16% in the parse test failing because of collections, and 10.8% of the VxC comparison set and 10% of correct parses classified as noisy, these frames are not as precise as the existing frames. The goal of these frames is not necessarily to replace the existing frames, but rather to extend coverage and provide a platform for further development whether by hand or through automatic methods. Precision can be improved upon in future

⁵ D. Lenat briefing, March 15, 2006

research and is facilitated by the expressiveness of Cyc. Predicate strengthening, using existing relationships to infer more specific predicates, is the next step in creating robust frames.

Additionally, there is a tradeoff between the number of frames covered and efficiency of disambiguation. More frame choices make it harder for parsers to choose the correct frame, but it will hopefully improve their handling of more complex sentence structures.

One possible solution to competition and overgenerality is to add verbs incrementally by class. The class based approach makes it easy to separate verbs by types, such as verbs that relate to mechanical processes or emotion verbs. One could use classes of frames to strengthen specific areas of parsing while choosing not to take verbs from a class covering a domain that the parser already performs strongly in. This approach can reduce interference with existing domains that have been hand built and extended beyond the standard Cyc KB for individual research.

Furthermore, semi-automatic approaches like this generate information more quickly than one could do by hand. Thus an approach to computational verb semantic representation that is rooted in classes can take advantage of modern reasoning sources like Cyc to efficiently create semantic knowledge.

Acknowledgments

This research was supported by the Air Force Office of Scientific Research and Northwestern University. A special thanks to Kenneth Forbus and the members of QRG for their continued invaluable guidance.

References

Crouch, Dick, and Tracy Holloway King. 2005. Unifying Lexical Resources. In *Proceedings of the Interdisciplinary Workshop on the Identification and Representation of Verb Features and Verb Classes*, Saarbruecken, Germany

Curtis, John, David Baxter, Peter Wagner, John Cabral, Dave Schneider, and Michael Witbrock. 2009. Methods of Rule Acquisition in the TextLearner System. In *Proceedings of the 2009 AAAI Spring Symposium on Learning by Reading and Learning to Read*, pages 22-28, Palo Alto, CA. AAAI Press.

Curtis, John, John Cabral, and David Baxter. 2006. On the Application of the Cyc Ontology to Word Sense Disambiguation. In *Proceedings of the Nineteenth International FLAIRS Conference*, pages 652-657, Melbourne Beach, FL.

Fellbaum, Christiane. Ed. 1998. *WordNet: An Electronic Database*. MIT Press, Cambridge, MA.

Forbus, Kenneth, Thomas Mostek, and Ron Ferguson. 2002. An Analogy Ontology for Integrating Analogical Processing and First-principle Reasoning. In *Proceedings of the Thirteenth Conference on Innovative Applications of Artificial Intelligence*. Menlo Park, CA. AAAI Press.

Kipper, Karin, Hoa Trang Dang, and Martha Palmer. 2000. Class-Based Construction of a Verb Lexicon. In *AAAI-2000 Seventeenth National Conference on Artificial Intelligence*, Austin, TX.

Kipper, Karin, Anna Korhonen, Neville Ryant, and Martha Palmer. 2006. Extending VerbNet with Novel Verb Classes. In *Fifth International Conference on Language Resources and Evaluation (LREC 2006)*. Genoa, Italy.

Levin, Beth. 1993. *English Verb Classes and Alternation: A Preliminary Investigation*. The University of Chicago Press, Chicago.

Matuszek, Cynthia, John Cabral, Michael Witbrock, and John DeOliveira. 2006. An Introduction to the Syntax and Content of Cyc. In *Proceedings of the 2006 AAAI Spring Symposium on Formalizing and Compiling Background Knowledge and Its Applications to Knowledge Representation and Question Answering*, Stanford, CA.

Moens, Marc, and Mark Steedman. 1988. Temporal Ontology and Temporal Reference. *Computational Linguistics*. 14(2):15-28.

Miller, G. 1985. WORDNET: A Dictionary Browser. In *Proceedings of the First International Conference on Information in Data*.

Ramachandran, Deepak, Pace Reagan, and Keith Goolsbey. 2005. First-Orderized Research Cyc: Expressivity and Efficiency in a Common-Sense Ontology. In *Papers from the AAAI Workshop on Contexts and Ontologies: Theory, Practice and Applications*. Pittsburgh, PA.

Tomai, Emmet, and Kenneth Forbus. 2009. EA NLU: Practical Language Understanding for Cognitive Modeling. In *Proceedings of the 22nd International Florida Artificial Intelligence Research Society Conference*, Sanibel Island, FL.

Trumbo, Derek. 2006. VxC: A VerbNet-Cyc Mapper. <http://verbs.colorado.edu/verb-index/vxc/>

A Framework for Figurative Language Detection Based on Sense Differentiation

Daria Bogdanova

University of Saint Petersburg

Saint Petersburg

dasha.bogdanova@gmail.com

Abstract

Various text mining algorithms require the process of feature selection. High-level semantically rich features, such as figurative language uses, speech errors etc., are very promising for such problems as e.g. writing style detection, but automatic extraction of such features is a big challenge. In this paper, we propose a framework for figurative language use detection. This framework is based on the idea of sense differentiation. We describe two algorithms illustrating the mentioned idea. We show then how these algorithms work by applying them to Russian language data.

1 Introduction

Various text mining algorithms require the process of feature selection. For example, authorship attribution algorithms need to determine features to quantify the writing style. Previous work on authorship attribution among computer scientists is mostly based on low-level features such as word frequencies, sentence length counts, n-grams etc. A significant advantage of such features is that they can be easily extracted from any corpus. But the study by Batov and Sorokin (1975) shows that such features do not always provide accurate measures for authorship attribution. The linguistic approach to the problem involves such high-level characteristics as the use of figurative language, irony, sound devices and so on. Such characteristics are very promising for the mentioned above tasks, but the extraction of these features is extremely hard to automate. As a result, very few attempts have been made to exploit high-level features for stylometric purposes (Stamatatos, 2009). Therefore, our long-term objective is the extraction of high-level semantically rich features.

Since the mentioned topic is very broad, we focus our attention only on some particular prob-

lems and approaches. In this paper, we examine one of such problems, the problem of automatic figurative language use detection. We propose a framework for figurative language detection based on the idea of sense differentiation. Then, we describe two algorithms illustrating the mentioned idea. One of them is intended to decide whether a usage is literal by comparing the texts related to the target expression and the set of texts related to the context itself. The other is aimed at grouping instances into literal and non-literal uses and is based on DBSCAN clustering (Ester et al, 1996). We illustrate then how these algorithms work by applying them to Russian language data. Finally, we propose some ideas on modifications which can significantly improve the accuracy of the algorithms.

2 Related Work

Sporleder and Li (April 2009) proposed an unsupervised method for recognition of literal and non-literal use of idiomatic expressions. Given an idiom the method detects the presence or absence of cohesive links between the words the idiom consists of and the surrounding text. When such links exist, the occurrence is considered as a literal usage and as a non-literal when there are no such links. For most idioms the experiments showed an accuracy above 50% (it varies between 11% and 98% for different idioms). The authors then proposed an improvement of the algorithm (Li and Sporleder, August 2009) by adding the Support Vector Machine classifier as a second stage. They use the mentioned above unsupervised algorithm to label the training data for the supervised classifier. The average accuracy of the improved algorithm is about 90%. Our approach is also based on the idea of the relatedness between the expression and the surrounding context. Unlike the mentioned study, we do not focus our attention only on idioms. So far we have mostly dealt with ex-

pressions, which are not necessarily idiomatic by themselves, but become metaphors in a particular context (e.g. "she is the *sunshine*", "life is a *journey*") and expressions that are invented by an author (e.g. "*my heart's in the Highlands*"). Moreover, the improved algorithm (Li and Sporleder, August 2009) is supervised, and our approach is unsupervised.

The study by Katz and Giesbrecht (2006) is also supervised, unlike ours. It also considers multi-word expressions that have idiomatic meanings. They propose an algorithm, which computes the vectors for literal and non-literal usages and then use the nearest neighbor classifier to label an unseen occurrence of the given idiom.

The approach proposed by Birke and Sarkar (2006) is nearly unsupervised. They constructed two seed sets: one consists of literal usages of different expressions and the other consists of non-literal usages. They calculate the distance between an occurrence in question and these two sets and assign to the occurrence the label of the closest set. This work, as well as ours, refers to the ideas from Word Sense Disambiguation area. Unlike our approach, the authors focus their attention only on the detection of figuratively used verbs and, whereas we only refer to the concepts and ideas of WSD, they adapt a particular existing one-word disambiguation method.

As we have already said, we deal with different types of figurative language (metaphors, metonymies etc.). However, there are some works aimed at extracting particular types of figurative language. For example, Nissim and Markert (2003) proposed a machine learning algorithm for metonymy resolution. They state the problem of metonymy resolution as a classification task between literal use of a word and a number of predefined metonymy types.

3 Sense Differentiation

We could treat a figurative meaning of a word as an additional, not common meaning of this word. Actually, some metaphors are quite common (e.g. *eye of a needle*, *leg of a table*, etc.) and are called catachretic metaphors. They appear in a language to remedy the gap in vocabulary (Black, 1954). These metaphors do not indicate an author's writing style: an author uses such metaphor for an object because the language has no other name for

that object. Therefore the algorithms we are developing do not work with this type of metaphors.

Our approach to figurative language detection is based on the following idea: the fact that the sense of a word significantly differs from the sense of the surrounding text usually indicates that the word is used figuratively. Two questions arise immediately:

1. How do we represent the sense of both the word and the surrounding context?
2. How do we find out that these senses differ significantly?

To answer the first question, we refer to the ideas popular in the Word Sense Disambiguation community: sense is a group of contextually similar occurrences of a word (Schütze, 1996). Hence, we represent the senses of both a word and its context as sets of documents related to the word and the context respectively. These sets can be obtained e.g. by searching Wikipedia, Google or another web search engine. For a word the query can be the word itself. As for a text, this query can be formulated as the whole text or as a set of some words contained in this text. It seems to us that querying the lexical chains (Halliday and Hasan, 1976) extracted from the text should provide better results than querying the whole text.

As soon as we have a sense representation for such objects as a word and a text, we should find a way to measure the difference between these sense representations and find out what difference is strong enough for the considered occurrence to be classified as a non-literal usage. One way to do this is representing sets of documents as sets of vectors and measuring the distance between the centers of the obtained vector sets. Another way is to apply clustering techniques to the sets and to measure the accuracy of the produced clustering. The higher the accuracy is, the more different the sets are.

Besides, this can be done by calculating text-to-text semantic similarity using for example the measure proposed by Mihalcea et al (2006). This is rather difficult in case of the Russian language because at the moment there is no WordNet-like taxonomies for Russian.

In the next section, we propose two algorithms based on the mentioned above idea. We state the algorithms generally and try to find out experi-

mentally what combination of the described techniques provides the best results.

4 Finding the Distance to the Typical Context Set

The algorithm is intended to determine whether a word (or an expression) in a given context is used literally or not.

As it was mentioned above, we decided to represent senses of both an expression and a context as sets of documents. Our hypothesis is that these document sets differ significantly if and only if an expression is used figuratively. Thus, the algorithm decides whether the occurrence is literal by comparing two sets of documents: the *typical context set*, which represents a sense of the expression, and the *related context set*, which represents a sense of the context. A naive way to construct the *typical context set* is searching some searching engine (e.g. Google) for the expression. Given a context with a target expression, the *related context set* can be constructed as follows:

1. Remove the target expression from the context;
2. Extract the longest lexical chains from the resulting context;
3. For every chain put to the set the first N articles retrieved by searching a searching engine for the chain;

After constructing the sets the algorithm should estimate the similarity between these two sets. This, for example, can be done by applying any clustering algorithm to the data and measuring the accuracy. Evidently, the higher the accuracy of the obtained clustering is, the more separated the sets are. It means that, when the usage is literal, the accuracy should be lower because we try to make two clusters out of data that should appear as the only cluster.

We hypothesize that in case of non-literal usages these two sets should be significantly separated.

Our experiments include two stages. During the first one we test our idea and estimate the parameters of the algorithms. During the second stage we test the more precise algorithm obtained during the first stage.

For the first stage, we found literal and non-literal occurrences of the following Russian words and expressions:

вьюга (snowstorm), дыхание (breath), кинжальный (dagger), плясать (dance), стебель гибкий (flexible (flower) stalk), утонуть (be drowned), хрустальный (crystal), шотландская волынка (bagpipes), мед (honey), лекарство (medicine).

For every expression, the *typical context set* consists of the first 10 articles retrieved by searching Google for the expression. In order to construct the second set we removed the target expression from the context and manually extracted lexical chains from the texts, although, the process of lexical chains extraction can be done automatically. However the algorithms on lexical chains extraction usually use WordNet to calculate the relatedness, but as it was already mentioned WordNet for the Russian language does not exist yet. Another way to calculate semantic relatedness is using Wikipedia (Mihalcea, 2007; Turdakov and Veliikhov, 2008), but it takes much effort. The second set for each occurrence consists of the first 10 articles retrieved by searching Google for the extracted chains. Then we applied k-means clustering algorithm ($k = 2$) to these sets. To evaluate the clustering we used measures from the clustering literature. We denote our sets by $G = g_1, g_2$ and the clusters obtained by k-means as $C = c_1, c_2$. We define a mapping f from the elements of G to the elements of C , such that each set g_i is mapped to a cluster $c_j = f(g_i)$ that has the highest percentage of common elements with g_i . Precision and recall for a cluster $g_i, i = 1, 2$ are defined as follows:

$$Pr_i = \frac{|f(g_i) \cap g_i|}{|f(g_i)|} \text{ and } Re_i = \frac{|f(g_i) \cap g_i|}{|g_i|}$$

Precision, Pr , and recall, Re , of the clustering are defined as the weighted averages of the precision and recall values over the sets:

$$Pr = \frac{1}{2}(Pr_1 + Pr_2) \text{ and } Re = \frac{1}{2}(Re_1 + Re_2)$$

F_1 -measure is defined as the harmonic mean of precision and recall, i.e.,

$$F_1 = \frac{2 \times Pr \times Re}{Pr + Re}.$$

Table 1 shows the results of the clustering. For 9 expressions out of 10, the clustering accuracy is higher in case of a metaphorical usage than in case of a literal one. Moreover, for 9 out of 10

	Figurative usage			Literal usage		
	<i>Pr</i>	<i>Re</i>	<i>F</i>	<i>Pr</i>	<i>Re</i>	<i>F</i>
вьюга	0,85	0,85	0,85	0,50	0,50	0,50
дыхание	0,83	0,75	0,79	0,65	0,60	0,63
кинжальный	0,85	0,85	0,85	0,70	0,65	0,67
плясать	0,95	0,95	0,95	0,66	0,65	0,66
стебель гибкий	0,85	0,85	0,85	0,88	0,85	0,86
утонул	0,85	0,85	0,85	0,81	0,70	0,75
хрустальный	0,95	0,95	0,95	0,83	0,75	0,78
шотландская волынка	0,88	0,85	0,86	0,70	0,70	0,70
мед	0,90	0,90	0,90	0,88	0,85	0,87
лекарство	0,90	0,90	0,90	0,81	0,70	0,75

Table 1: Results provided by k-means clustering

metaphorical usages, F-measure is 0,85 or higher. And for 7 out of 10 literal usages, F-measure is 0,75 or less.

The first stage of the experiments illustrates the idea of sense differentiation. Based on the obtained results, we have concluded, that F-measure value equal to 0,85 or higher indicates a figurative usage, and the value equal to 0,75 or less indicates a literal usage.

At the second stage, we applied the algorithm to several Russian language expressions used literally or figuratively. The accuracy of the k-means clustering is shown in Table 2.

Figurative usages			
живой костер из снега и вина	0,76	0,55	0,64
лев	1,00	1,00	1,00
иней	0,90	0,90	0,90
ключ	0,95	0,93	0,94
лютый зверь	0,88	0,85	0,87
рогатый	0,92	0,90	0,91
терлась о локоть	0,88	0,85	0,86
иглою снежного огня	0,95	0,95	0,95
клавишей стая	0,76	0,55	0,64
горели глаза	0,95	0,95	0,95
цветок	0,80	0,80	0,80
загорелся	0,91	0,90	0,90
Literal usages			
ловил рыбу	0,71	0,70	0,70
играл в футбол	0,74	0,70	0,71
детство	0,66	0,65	0,66
кухня	0,88	0,85	0,87
снег	0,95	0,95	0,95
весна	0,50	0,50	0,50
пить кофе	0,85	0,85	0,85
танцы	0,90	0,90	0,90
платье	0,65	0,65	0,65
человек	0,81	0,70	0,75
ветер	0,85	0,85	0,85
дождь	0,91	0,90	0,90

Table 2: Testing the algorithm. Accuracy of the k-means clustering

For 75% of metaphorical usages F-measure is 0,85 or more as was expected and for 50% of literal usages F-measure is 0,75 or less.

5 Figurative Language Uses as Outliers

The described above approach is to decide whether a word in a context is used literally or not. Unlike the first one, the second approach we propose, deals with a set of occurrences of a word as to label every occurrence as 'literal' or 'non-literal'. We formulate this task as a clustering problem and apply DBSCAN (Ester et al, 1996) clustering algorithm to the data. Miller and Charles (1991) hypothesized that words with similar meanings are often used in similar contexts. As it was mentioned, we can treat a meaning of a metaphoric usage of an expression as an additional, not common for the expression. That's why we expect metaphorical usages to be outliers, while clustering together with common (i.e. literal) usages. Theoretically, the algorithm should also distinguish between all literal senses so that the contexts of the same meaning appear in the same cluster and the contexts of different meanings - in different clusters. Therefore, ideally, the algorithm should solve word sense discrimination and non-literal usages detection tasks simultaneously.

For each Russian word shown in Table 3, we extracted from the Russian National Corpora (<http://ruscorpora.ru/>) several literal and non-literal occurrences. Some of these words have more than one meaning in Russian, e.g. *ключ* can be translated as a *key* or *water spring* and the word *коса* as a *plait*, *scythe* or *spit*.

word	literal	non-literal
бабочка (butterfly, bow-tie)	12	2
иней (frost)	14	2
ключ (key, spring(water))	14	2
коса (plait, scythe, spit)	21	2
лев (lion, Bulgarian lev)	17	5
лук (onion, bow)	17	1
мука (flour, pain)	21	2
пыль (dust)	14	4

Table 3: Data used in the first experiment

All the documents are stemmed and all stop-words are removed with the SnowBall Stemmer (<http://snowball.tartarus.org/>) for the Russian language.

As it was mentioned above, this algorithm is aimed at providing word sense discrimination and non-literal usages detection simultaneously. So far we have paid attention only to the non-literal usages detection aspects. DBSCAN algorithm is a density-based clustering algorithm designed to

discover clusters of arbitrary shape. This algorithm requires two parameters: ε (eps) and the minimum number of points in a cluster (minPts).

We set minPts to 3 and run the algorithm for different eps between 1.45 and 1.55.

As was mentioned, so far we have considered only figurative language detection issues: The algorithm marks an instance as a figurative usage iff the instance is labeled as an outlier. Thus, we measure the accuracy of the algorithm as follows:

$$precision = \frac{|figurative\ uses \cap outliers|}{|outliers|},$$

$$recall = \frac{|figurative\ uses \cap outliers|}{|figurative\ uses|}.$$

Figures 1 and 2 show the dependency between the eps parameter and the algorithm's accuracy for different words.

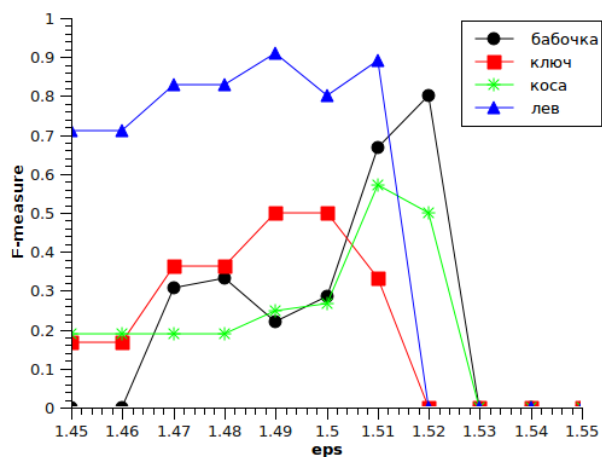


Figure 1: Dependency between eps and F-measure

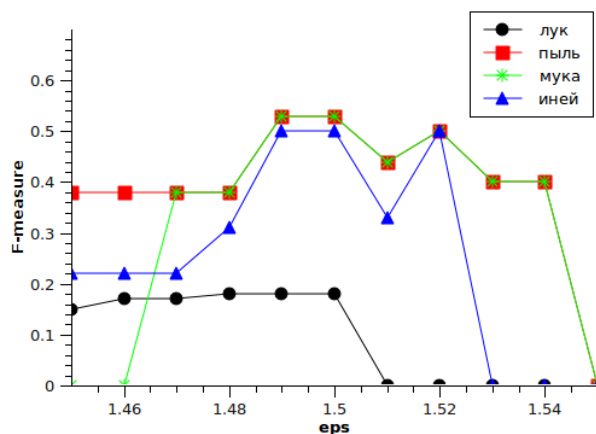


Figure 2: Dependency between eps and F-measure

Table 4 shows "the best" eps for each word and the corresponding accuracies of metaphor detection

word	eps	precision	recall
бабочка	1.520	0.66	1.00
иней	1.520	0.50	0.50
ключ	1.500	0.33	1.00
коса	1.510	0.40	1.00
лев	1.490	1.00	0.83
лук	1.505	0.17	1.00
мука	1.525	0.67	0.50
пыль	1.505	0.50	0.60

Table 4: The best eps parameters and corresponding accuracies of the algorithm

6 Future Work

So far we have worked only with tf-idf and word frequency model for both algorithms. The next step in our study is utilizing different text representation models, e.g. second order context vectors. We are also going to develop an efficient parameter estimation procedure for the algorithm based on DBSCAN clustering.

As for the other algorithm, we are going to distinguish between different figurative language expressions:

- one word expressions
 - monosemous word
 - polysemous word
- multiword expressions

We expect the basic algorithm to provide different accuracy in case of different types of expressions. Dealing with multiword expressions and monosemous words should be easier than with polysemous words: i.e., for monosemous word we expect the second set to appear as one cluster, whereas this set for a polysemous word is expected to have the number of clusters equal to the number of senses it has.

Another direction of the future work is developing an algorithm for figurative language uses extraction. The algorithm has to find figuratively used expressions in a text.

7 Conclusion

In this paper, we have proposed a framework for figurative language detection based on the idea of sense differentiation. We have illustrated how this

idea works by presenting two clustering-based algorithms. The first algorithm deals with only one context. It is based on comparing two context sets: one is related to the expression and the other is semantically related to the given context. The second algorithm groups the given contexts in literal and non-literal usages. This algorithm should also distinguish between different senses of a word, but we have not yet paid enough attention to this aspect. By applying these algorithms to small data sets we have illustrated how the idea of sense differentiation works. These algorithms show quite good results and are worth further work.

Acknowledgments

This work was partially supported by Russian Foundation for Basic Research RFBR, grant 10-07-00156.

References

- Vitaly I. Batov and Yury A. Sorokin. 1975. Text attribution based on objective characteristics. *Seriya yazyka i literatury*, 34, 1.
- Julia Birke and Anoop Sarkar. 2006. A Clustering Approach for the Nearly Unsupervised Recognition of Nonliteral Language. *Proceedings of EACL-06*
- Max Black. 1954. Metaphor. *Proceedings of the Aristotelian Society*, 55, pp. 273-294.
- Martin Ester, Hans-Peter Kriegel, Jörg Sander and Xiaowei Xu. 1996. A density-based algorithm for discovering clusters in large spatial databases with noise. *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD-96)*, AAAI Press, pp. 226-231
- Michael Halliday and Ruqaiya Hasan. 1976. Cohesion in English. *Longman, London*
- Graham Katz and Eugenie Giesbrecht. 2006. Automatic identification of non-compositional multiword expressions using latent semantic analysis. *Proceedings of the ACL/COLING-06 Workshop on Multiword Expressions: Identifying and Exploiting Underlying Properties*
- Linlin Li and Caroline Sporleder. August 2009. Classifier combination for contextual idiom detection without labeled data. *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing*, pp. 315-323.
- Rada Mihalcea, Courtney Corley and Carlo Strapparava. 2006. Corpus-based and Knowledge-based Measures of Text Semantic Similarity. *Proceedings of AAAI-06*
- Rada Mihalcea. 2007. Using Wikipedia for Automatic Word Sense Disambiguation. *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL 2007)*
- George A. Miller and Walter G. Charles. 1991. Contextual correlates of semantic similarity. *Language and Cognitive Processes*, 6(1):128.
- Malvina Nissim and Katja Markert. 2003. Syntactic features and word similarity for supervised metonymy resolution. *Proceedings of the 41st Annual Meeting of the Association for Computational Linguistics (ACL-03)* (Sapporo, Japan, 2003). 56-63.
- Hinrich Schütze. 1998. Automatic word sense discrimination. *Computational Linguistics*, 24(1), pp. 97-123
- Caroline Sporleder and Linlin Li. April 2009. Unsupervised recognition of literal and non-literal use of idiomatic expressions. *Proceedings of EACL-09*
- Efstathios Stamatatos. 2009. A survey of modern authorship attribution methods. *Journal of the American Society for information Science and Technology*, 60(3): 538-556.
- Denis Turdakov and Pavel Velikhov. 2008. Semantic Relatedness Metric for Wikipedia Concepts Based on Link Analysis and its Application to Word Sense Disambiguation *SYRCoDIS 2008*

Automatic Selectional Preference Acquisition for Latin verbs

Barbara McGillivray

University of Pisa

Italy

b.mcgillivray@ling.unipi.it

Abstract

We present a system that automatically induces Selectional Preferences (SPs) for Latin verbs from two treebanks by using Latin WordNet. Our method overcomes some of the problems connected with data sparseness and the small size of the input corpora. We also suggest a way to evaluate the acquired SPs on unseen events extracted from other Latin corpora.

1 Introduction

Automatic acquisition of semantic information from corpora is a challenge for research on low-resourced languages, especially when semantically annotated corpora are not available. Latin is definitely a high-resourced language for what concerns the number of available texts and traditional lexical resources such as dictionaries. Nevertheless, it is a low-resourced language from a computational point of view (McGillivray et al., 2009).

As far as NLP tools for Latin are concerned, parsing experiments with machine learning techniques are ongoing (Bamman and Crane, 2008; Passarotti and Ruffolo, forthcoming), although more work is still needed in this direction, especially given the small size of the training data. As a matter of fact, only three syntactically annotated Latin corpora are available (and still in progress): the Latin Dependency Treebank (LDT, 53,000 tokens) for classical Latin (Bamman and Crane, 2006), the *Index Thomisticus* Treebank (IT-TB, 54,000 tokens) for Thomas Aquinas's works (Passarotti, 2007), and the PROIEL treebank (approximately 100,000 tokens) for the Bible (Haug and Jøndal, 2008). In addition, a Latin version of WordNet – Latin WordNet (LWN; Minozzi, (2009) – is being compiled, consisting of around 10,000 lemmas inserted in the multilingual structure of MultiWordNet (Bentivogli et al., 2004).

The number and the size of these resources are small when compared with the corpora and the lexicons for modern languages, e. g. English.

Concerning semantic processing, no semantically annotated Latin corpus is available yet; building such a corpus manually would take considerable time and energy. Hence, research in computational semantics for Latin would benefit from exploiting the existing resources and tools through automatic lexical acquisition methods.

In this paper we deal with automatic acquisition of verbal selectional preferences (SPs) for Latin, i. e. the semantic preferences of verbs on their arguments: e. g. we expect the object position of the verb *edo* 'eat' to be mostly filled by nouns from the food domain. For this task, we propose a method inspired by Alishahi (2008) and outlined in an earlier version on the IT-TB in McGillivray (2009). SPs are defined as probability distributions over semantic features extracted as sets of LWN nodes. The input data are two subcategorization lexicons automatically extracted from the LDT and the IT-TB (McGillivray and Passarotti, 2009).

Our main contribution is to create a new tool for semantic processing of Latin by adapting computational techniques developed for extant languages to the special case of Latin. A successful adaptation is contingent on overcoming corpus size differences. The way our model combines the syntactic information contained in the treebanks with the lexical semantic knowledge from LWN allows us to overcome some of the difficulties related to the small size of the input corpora. This is the main difference from corpora for modern languages, together with the absence of semantic annotation. Moreover, we face the problem of evaluating our system's ability to generalize over unseen cases by using text occurrences, as access to human linguistic judgements is denied for Latin.

In the rest of the paper we will briefly summarize previous work on SP acquisition and motivate

vado}' (*{progress, come_on, come_along, advance, get_on, get_along, shape_up}*' in the English WN).

3.1 Clustering of frames

The constructions are incrementally built as new frames are included in them; a new frame F is assigned to a construction K if F probabilistically shares some features with the frames in K so that

$$K = \arg \max_k P(k|F) = \arg \max_k P(k)P(F|k),$$

where k ranges over the set of all constructions, including the baseline $k_0 = \{F\}$. The prior probability $P(k)$ is calculated from the number of frames contained in k divided by the total number of frames. Assuming that the frame features are independent, the posterior probability $P(F|k)$ is the product of three probabilities, each one corresponding to the probability that a feature displays in k the same value it displays in F : $P_i(ft_i(F)|k)$ for $i = 1, 2, 3$:

$$P(F|k) = \prod_{i=1,2,3} P_i(ft_i(F)|k)$$

We estimated the probability of a match between the value of ft_1 in k and the value of ft_1 in F as the sum of the *syntactic scores* between F and each frame h contained in k , divided the number n_k of frames in k :

$$P(ft_1(F)|k) = \frac{\sum_{h \in k} \text{synt_score}(h, F)}{n_k}$$

where the syntactic score $\text{synt_score}(h, F) = \frac{|SCS(h) \cap SCS(F)|}{|SCS(F)|}$ calculates the number of syntactic slots shared by h and F over the number of slots in F . $P(ft_1(F)|k)$ is 1 when all the frames in k contain all the syntactic slots of F .

For each argument position a , we estimated the probability $P(ft_2(F)|k)$ as the sum of the *semantic scores* between F and each h in k :

$$P(ft_2(F)|k) = \frac{\sum_{h \in k} \text{sem_score}(h, F)}{n_k}$$

where the semantic score $\text{sem_score}(h, F) = \frac{|S(h) \cap S(F)|}{|S(F)|}$ counts the overlap between the semantic properties $S(h)$ of h (i. e. the LWN hypernyms of the fillers in h) and the semantic properties $S(F)$ of F (for argument a), over $|S(F)|$.

$$P(ft_3(F)|k) = \frac{\sum_{h \in k} \text{syms_score}(h, F)}{n_k}$$

where the synset score $\text{syms_score}(h, F) = \frac{|Synsets(\text{verb}(h)) \cap Synsets(\text{verb}(F))|}{|Synsets(\text{verb}(F))|}$ calculates the overlap between the synsets for the verb in h and the synsets for the verb in F over the number of synsets for the verb in F .³

We introduced the syntactic and synset scores in order to account for a frequent phenomenon in our data: the partial matches between the values of the features in F and in k .

3.2 Selectional preferences

The clustering algorithm defines the set of constructions in which the generalization step over unseen cases is performed. SPs are defined as semantic profiles, that is, probability distributions over the semantic properties, i. e. LWN nodes. For example, we get the probability of the node *actio* 'act' in the position 'A_(in)Obj[acc]' for *eo* 'go'.

If s is a semantic property and a an argument position for a verb v , the semantic profile $P_a(s|v)$ is the sum of $P_a(s, k|v)$ over all constructions k containing v or a WN-synonym of v , i. e. a verb contained in one or more synsets for v . $P_a(s, k|v)$ is approximated as $\frac{P(k,v)P_a(s|k,v)}{P(v)}$, where $P(k, v)$ is estimated as $\frac{n_k \cdot \text{freq}(k,v)}{\sum_{k'} n_{k'} \cdot \text{freq}(k',v)}$

To estimate $P_a(s|k, v)$ we consider each frame h in k and account for: a) the similarity between v and the verb in h ; b) the similarity between s and the fillers of h . This is achieved by calculating a *similarity score* between h, v, a and s , defined as:

$$\text{syms_score}(v, V(h)) \cdot \frac{\sum_f |s \cap S(f)|}{N_{\text{fil}}(h, a)} \quad (1)$$

where $V(h)$ in (1) contains the verbs of h , $N_{\text{fil}}(h, a)$ counts the a -fillers in h , f ranges in the set of a -fillers in h , $S(f)$ contains the semantic properties for f and $|s \cap S(f)|$ is 1 when s appears in $S(f)$ and 0 otherwise.

$P_a(s|k, v)$ is thus obtained by normalizing the sum of these similarity scores over all frames in k , divided by the total number of frames in k containing v or its synonyms.

The similarity scores weight the contributions of the synonyms of v , whose fillers play a role in the generalization step. This is our innovation with respect to Alishahi (2008)'s system. It was introduced because of the sparseness of our data, where

³The algorithm uses smoothed versions of all the previous formulae by adding a very small constant so that the probabilities are never 0.

k	h
1	induco + P_Sb[acc]{forma} introduco + P_Sb{PR} introduco + P_Sb{forma} addo + P_Sb{praesidium}
2	induco + A_Obj[acc]{forma} immitto + A_Obj[acc]{PR}, Obj[dat]{antrum} introduco + A_Obj[acc]{NP}
3	introduco + A_(in)Obj[acc]{finis}, Obj[acc]{copia}, Sb{NP} induco + A_(in)Obj[acc]{effectus}, Obj[acc]{forma}
4	introduco + A_Obj[acc]{forma} induco + A_Obj[acc]{perfectio}, Sb[nom]{PR}
5	induco + A_Obj[acc]{forma}n immitto + A_Obj[acc]{PR}, Obj[dat]{antrum} introduco + A_Obj[acc]{NP}

Table 1: Constructions (k) for the frames (h) containing the verb *introduco* ‘bring in’.

many verbs are hapaxes, which makes the generalization from their fillers difficult.

4 Results and evaluation

The clustering algorithm was run on 15509 frames and it generated 7105 constructions. Table 1 displays the 5 constructions assigned to the 9 frames where the verb *introduco* ‘bring in, introduce’ occurs. Note the semantic similarity between *addo* ‘add to, bring to’, *immitto* ‘send against, insert’, *induco* ‘bring forward, introduce’ and *introduco*, and the similarity between the syntactic patterns and the argument fillers within the same construction. For example, *finis* ‘end, borders’ and *effectus* ‘result’ share the semantic properties ATTRIBUTE, COGNITIO ‘cognition’, CONSCIENTIA ‘conscience’, EVENTUM ‘event’, among others.

The vast majority of constructions contain less than 4 frames. This contrasts with the more general constructions found by Alishahi (2008) and can be explained by several factors. First, the coverage of LWN is quite low with respect to the fillers in our dataset. In fact, 782 fillers out of 2408 could not be assigned to any LWN synset; for these lemmas the semantic scores with all the other nouns are 0, causing probabilities lower than the baseline; this results in assigning the frame to the singleton construction consisting of the frame itself. The same happens for fillers consisting of verbal lemmas, participles, pronouns and named entities, which amount to a third of the total number. Furthermore, the data are not tagged by sense and the system deals with noun ambiguity by listing together all synsets of a word n (and their hypernyms) to form the semantic properties for n : consequently, each sense contributes to the semantic description of n in relation to the number of hypernyms it carries, rather than to its observed

semantic property	probability
<i>actio</i> ‘act’	0.0089
<i>actus</i> ‘act’	0.0089
<i>pars</i> ‘part’	0.0089
<i>object</i>	0.0088
<i>physical object</i>	0.0088
<i>instrumentality</i>	0.0088
<i>instrumentation</i>	0.0088
<i>location</i>	0.0088
<i>populus</i> ‘people’	0.0088
<i>plaga</i> ‘region’	0.0088
<i>regio</i> ‘region’	0.0088
<i>arvum</i> ‘area’	0.0088
<i>orbis</i> ‘area’	0.0088
<i>external body part</i>	0.0088
<i>nympha</i> ‘nymph’, ‘water’	0.0088
<i>latex</i> ‘water’	0.0088
<i>lymph</i> ‘water’	0.0088
<i>intercapedo</i> ‘gap, break’	0.0088
<i>orificium</i> ‘opening’	0.0088

Table 2: Top 20 semantic properties in the semantic profile for *ascendo* ‘ascend’ + A_(de)Obj[abl].

frequency. Finally, a common problem in SP acquisition systems is the noise in the data, including tagging and metaphorical usages. This problem is even greater in our case, where the small size of the data underestimates the variance and therefore overestimates the contribution of noisy observations. Metaphorical and abstract usages are especially frequent in the data from the IT-TB, due to the philosophical domain of the texts.

As to the SP acquisition, we ran the system on all constructions generated by the clustering. We excluded the pronouns occurring as argument fillers, and manually tagged the named entities. For each verb lemma and slot we obtained a probability distribution over the 6608 LWN noun nodes.

Table 2 displays the 20 semantic properties with the highest SP probabilities as ablative arguments of *ascendo* ‘ascend’ introduced by *de* ‘down from’, ‘out of’. This semantic profile was created from the following fillers for the verbs contained in the constructions for *ascendo* and its synonyms: *abyssus* ‘abyss’, *fumus* ‘smoke’, *lacus* ‘lake’, *machina* ‘machine’, *manus* ‘hand’, *negotatio* ‘business’, *mare* ‘sea’, *os* ‘mouth’, *templum* ‘temple’, *terra* ‘land’. These nouns are well represented by the semantic properties related to water and physical places. Note also the high rank of general properties like *actio* ‘act’, which are associated to a large number of fillers and thus generally get a high probability.

Regarding evaluation, we are interested in testing two properties of our model: calibration and discrimination. Calibration is related to the model’s ability to distinguish between high and low probabilities. We verify that our model is

adequately calibrated, since its SP distribution is always very skewed (cf. figure 1). Therefore, the model is able to assign a high probability to a small set of nouns (preferred nouns) and a low probability to a large set of nouns (the rest), thus performing better than the baseline model, defined as the model that assigns the uniform distribution over all nouns (4724 LWN leaf nodes). Moreover, our model’s entropy is always lower than the baseline: 12.2 vs. the 6.9-11.3 range; by the maximum entropy principle, this confirms that the system uses some information for estimating the probabilities: LWN structure, co-occurrence frequency, syntactic patterns. However, we have no guarantee that the model uses this information sensibly. For this, we test the system’s discrimination potential, i. e. its ability to correctly estimate the SP probability of each single LWN node.

noun	SP probability
<i>pars</i> ‘part’	0.0029
<i>locus</i> ‘place’	0.0026
<i>forma</i> ‘form’	0.0023
<i>ratio</i> ‘account’ ‘reason’, ‘opinion’	0.0023
<i>respectus</i> ‘consideration’	0.0022
<i>caput</i> ‘head’, ‘origin’	0.0022
<i>anima</i> ‘soul’	0.0021
<i>animus</i> ‘soul’, ‘spirit’	0.0020
<i>figura</i> ‘form’, ‘figure’	0.0020
<i>spiritus</i> ‘spirit’	0.0020
<i>causa</i> ‘cause’	0.0020
<i>corpus</i> ‘body’	0.0019
<i>sententia</i> ‘judgement’	0.0019
<i>finitio</i> ‘limit’, ‘definition’	0.0019
<i>species</i> ‘sight’, ‘appearance’	0.0019

Table 3: 15 nouns with the highest probabilities as accusative objects of *dico* ‘say’.

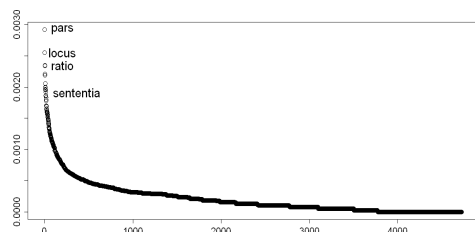


Figure 1: Decreasing SP probabilities of the LWN leaf nodes for the objects of *dico* ‘say’.

Table 3 displays the 15 nouns with the highest probabilities as direct objects for *dico* ‘say’. From table 3 – and the rest of the distribution, represented in figure 1 – we see that the model assigns a high probability to most seen fillers for *dico* in the corpus: *anima* ‘soul’, *corpus* ‘body’, *locus*

‘place’, *pars* ‘part’, etc.

For what concerns evaluating the SP probability assigned to nouns unseen in the training set, Alishahi (2008) follows the approach suggested by Resnik (1993), using human plausibility judgements on verb-noun pairs. Given the absence of native speakers of Latin, we used random occurrences in corpora, considered as positive examples of plausible argument fillers; on the other hand, we cannot extract non-plausible fillers from a corpus unless we use a frequency-based criterion. However, we can measure how well our system predicts the probability of these unseen events.

As a preliminary evaluation experiment, we randomly selected from our corpora a list of 19 high-frequency verbs ($\text{freq.} > 51$) and 7 medium-frequency verbs ($11 < \text{freq.} < 50$), for each of which we chose an interesting argument slot. Then we randomly extracted one filler for each such pair from two collections of Latin texts (*Perseus Digital Library* and *Corpus Thomisticum*), provided that it was not in the training set. The semantic score in equation 1 on page 3 is then calculated between the set of semantic properties of n and that for f , to obtain the probability of finding the random filler n as an argument for a verb v .

For each of the 26 (verb, slot) pairs, we looked at three measures of central tendency: mean, median and the value of the third quantile, which were compared with the probability assigned by the model to the random filler. If this probability was higher than the measure, the outcome was considered a success. The successes were 22 for the mean, 25 for the median and 19 for the third quantile.⁴ For all three measures a binomial test found the success rate to be statistically significant at the 5% level. For example, table 3 and figure 1 show that the filler for *dico*+A_Obj[acc] in the evaluation set – *sententia* ‘judgement’ – is ranked 13th within the verb’s semantic profile.

5 Conclusion and future work

We proposed a method for automatically acquiring probabilistic SP for Latin verbs from a small corpus using the WN hierarchy; we suggested some

⁴The dataset consists of all LWN leaf nodes n , for which we calculated $P_a(n|v)$. By definition, if we divide the dataset in four equal-sized parts (*quartiles*), 25% of the leaf nodes have a probability higher than the value at the third quartile. Therefore, in 19 cases out of 26 the random fillers are placed in the high-probability quarter of the plot, which is a good result, since this is where the preferred arguments gather.

new strategies for tackling the data sparseness in the crucial generalization step over unseen cases. Our work also contributes to the state of the art in semantic processing of Latin by integrating syntactic information from annotated corpora with the lexical resource LWN. This demonstrates the usefulness of the method for small corpora and the relevance of computational approaches for historical linguistics.

In order to measure the impact of the frame clusters for the SP acquisition, we plan to run the system for SP acquisition without performing the clustering step, thus defining all constructions as singleton sets containing one frame each. Finally, an extensive evaluation will require a more comprehensive set, composed of a higher number of unseen argument fillers; from the frequencies of these nouns, it will be possible to directly compare plausible arguments (high frequency) and implausible ones (low frequency). For this, a larger automatically parsed corpus will be necessary.

6 Acknowledgements

We wish to thank Afra Alishahi, Stefano Minozzi and three anonymous reviewers.

References

- E. Agirre and D. Martinez. 2001. Learning class-to-class selectional preferences. In *Proceedings of the ACL/EACL 2001 Workshop on Computational Natural Language Learning (CoNLL-2001)*, pages 1–8.
- A. Alishahi. 2008. *A probabilistic model of early argument structure acquisition*. Ph.D. thesis, Department of Computer Science, University of Toronto.
- D. Bamman and G. Crane. 2006. The design and use of a Latin dependency treebank. In *Proceedings of the Fifth International Workshop on Treebanks and Linguistic Theories*, pages 67–78. ÚFAL MFF UK.
- D. Bamman and G. Crane. 2008. Building a dynamic lexicon from a digital library. In *Proceedings of the 8th ACM/IEEE-CS Joint Conference on Digital Libraries*, pages 11–20.
- L. Bentivogli, P. Forner, and Pianta E. Magnini, B. 2004. Revising wordnet domains hierarchy: Semantics, coverage, and balancing. In *Proceedings of COLING Workshop on Multilingual Linguistic Resources*, pages 101–108.
- S. Clark and D. Weir. 1999. An iterative approach to estimating frequencies over a semantic hierarchy. In *Proceedings of the Joint SIGDAT Conference on Empirical Methods in Natural Language Processing and Very Large Corpora. University of Maryland*, pages 258–265.
- K. Erk. 2007. A simple, similarity-based model for selectional preferences. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics*, pages 216–223.
- D. T. T. Haug and M. L. Jøndal. 2008. Creating a parallel treebank of the old Indo-European Bible translations. In *Proceedings of Language Technologies for Cultural Heritage Workshop*, pages 27–34.
- H. Li and N. Abe. 1998. Generalizing case frames using a thesaurus and the MDL principle. *Computational Linguistics*, 24(2):217–244.
- B. McGillivray and M. Passarotti. 2009. The development of the *Index Thomisticus* Treebank Valency Lexicon. In *Proceedings of the Workshop on Language Technology and Resources for Cultural Heritage, Social Sciences, Humanities, and Education*, pages 33–40.
- B. McGillivray, M. Passarotti, and P. Ruffolo. 2009. The *Index Thomisticus* treebank project: Annotation, parsing and valency lexicon. *TAL*, 50(2):103–127.
- B. McGillivray. 2009. Selectional Preferences from a Latin treebank. In Przepiórkowski A. Passarotti, M., S. Raynaud, and F. van Eynde, editors, *Proceedings of the Eighth International Workshop on Treebanks and Linguistic Theories (TLT8)*, pages 131–136. EDUCatt.
- S. Minozzi. 2009. The Latin Wordnet project. In P. Anreiter and M. Kienpointner, editors, *Proceedings of the 15th International Colloquium on Latin Linguistics (ICLL)*, Innsbrucker Beitrage zur Sprachwissenschaft.
- M. Passarotti and P. Ruffolo. forthcoming. Parsing the *Index Thomisticus* Treebank. some preliminary results. In P. Anreiter and M. Kienpointner, editors, *Proceedings of the 15th International Colloquium on Latin Linguistics*, Innsbrucker Beiträge zur Sprachwissenschaft.
- M. Passarotti. 2007. Verso il Lessico Tomistico Biculturale. La treebank dell’*Index Thomisticus*. In R. Petrilli and D. Femia, editors, *Atti del XIII Congresso Nazionale della Società di Filosofia del Linguaggio*, pages 187–205.
- P. Resnik. 1993. *Selection and Information: A Class-Based Approach to Lexical Relationships*. Ph.D. thesis, University of Pennsylvania.
- M. Rooth, S. Riezler, D. Prescher, G. Carroll, and F. Beil. 1999. Inducing a semantically annotated lexicon via EM-based clustering. In *Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics*, pages 104–111.

Edit Tree Distance alignments for Semantic Role Labelling

Hector-Hugo Franco-Penya

Trinity College Dublin

Dublin, Ireland.

francoph@cs.tcd.ie

Abstract

“Tree SRL system” is a Semantic Role Labelling supervised system based on a tree-distance algorithm and a simple k-NN implementation. The novelty of the system lies in comparing the sentences as tree structures with multiple relations instead of extracting vectors of features for each relation and classifying them. The system was tested with the English CoNLL-2009 shared task data set where 79% accuracy was obtained.

1 Introduction

Semantic Role Labelling (SRL) is a natural language processing task which deals with semantic analysis at sentence-level. SRL is the task of identifying arguments for a certain predicate and labelling them. The predicates are usually verbs. They establish “what happened”. The arguments determine events such as “who”, “whom”, “where”, etc, with reference to one predicate. The possible semantic roles are pre-defined for each predicate. The set of roles depends on the corpora.

SRL is becoming an important tool for information extraction, text summarization, machine translation and question answering (Márquez, et al, 2008).

2 The data

The data set I used is taken from the CoNLL-2009 shared task (Hajič et al., 2009) and is part of Propbank. Propbank (Palmer et al, 2005) is a hand-annotated corpus. It transforms sentences into propositions. It adds a semantic layer to the Penn TreeBank (Marcus et al, 1994) and defines a set of semantic roles for each predicate.

It is difficult to define universal semantic roles for all predicates. That is why PropBank defines a set of semantic roles for each possible sense of each predicate (frame) [See a sample of the frame “raise” on the Figure 1 caption].

	Sentences	predicates	arguments	Predicates per sentence	arguments per sub-tree	File size in Mb
Tra	39279	179014	393699	4.55	2.20	56.2
Dev	1334	6390	13865	4.79	2.17	1.97
Evl	2399	10498	23286	4.38	2.22	3.41

Table 1: The data

The data set is divided into three files: training (Tra), development (Dev) and evaluation (Evl). The following table describes the number of sentences, sub-trees and labels contained in them, and the ratios of sub-trees per sentences and relations per sub-tree.

The core arguments are labelled by numbers. Adjuncts, which are common to all predicates, have their own labels, like: AM-LOC, TMP, NEG, etc. The four most frequent labels in the data set are: A1:35%, A0:20.86%, A2:7.88% and AM-TMP: 7.72%

Propbank was originally built using constituent tree structures, but here only the dependency tree structure version was used. Note that dependency tree structures have labels on the arrows. The tree distance algorithm cannot work with these labelled arrows and so they are moved to the child node as an extra label.

The task performed by the Tree SRL system consists of labelling the relations (predicate arguments) which are assumed to be already identified.

3 Tree Distance

The tree distance algorithm has already been applied to text entailment (Kouylekov & Magnini, 2005) and question answering (Punyakanok et al, 2004; Emms, 2006) with positive results.

The main contribution of this piece of work to the SRL field is the inclusion of the tree distance algorithm into an SRL system, working with tree structures in contrast to the classical “feature extraction” and “classification”. Kim et al (2009) developed a similar system for Information Extraction.

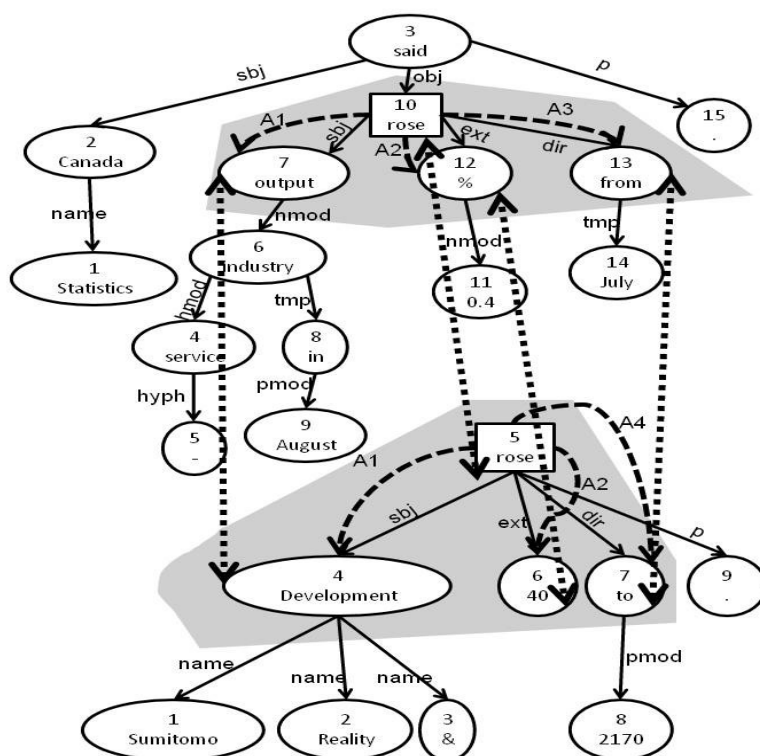


Figure 1: Alignment sample

A two sentence sample, in a dependency tree representation. In each node, the word form and the position of the word in the sentence are shown. Straight arrows represent syntactic dependencies. The label of the dependency is not shown. The square node represent the predicate that is going to be analyzed, (there can be multiple predicates in a single sentence). Semi-dotted arrows between a square node and an ellipse node represent a semantic relation. This arrow has a semantic tag (A1, A2, A3 and A4).

The grey shadow contains all the nodes of the sub tree for the “rose” predicate.

The dotted double arrows between the nodes of both sentences represent the tree distance alignment for both sub-trees. In this particular case every single node is matched.

Both predicate nodes are samples of the frame “raise” sense 01 (which means “go up quantifiably”) where the core arguments are:

- A0:** Agent, causer of motion
- A1:** Logical subject, patient, thing rising
- A2:** EXT, amount raised
- A3:** Start point
- A4:** End point
- AM:** Medium

Tai (1979) introduced a criterion for matching nodes between tree representations (or converting one tree into another one) and (Shasha & Zhang, 1990; Zhang & Shasha, 1989) developed an algorithm that finds an optimal matching tree solution for any given pair of trees. The advantage of this algorithm is that its computational cost is low. The optimal matching depends on the defined atomic cost of matching two nodes.

4 Tree SRL system architecture

For the training and testing data set, all possible sub-trees were extracted. Figure 3 and Figure 5

describe the process. Then, using the tree distance algorithm, the test sub-trees are labelled using the training ones. Finally, the predicted labels get assembled on the original sentence where the test sub-tree came from. Figure 2 describes the process.

A sub-tree extracted from a sentence, contains a predicate node, all its argument nodes and all the ancestors up to the first common ancestor of all nodes. (Figure 1 shows two samples of sub-tree extraction. Figure 3 describes how sub trees are obtained)

Input: training data set (labelled)
Input: testing data set (unlabelled)
Output: testing data set (labelled)
 Load training and testing data;
 Adapt the trees for the tree distance algorithm;
foreach sentence (training & testing data) **do**
 obtain each minimal sub-tree for each pre-
 predicate;
end
foreach sub-tree T from the testing data **do**
 calculate the distance and the alignment
 from T to each training sub-tree;
 sort the list of alignments by ascending
 tree distance;
 use the list to label the sub-tree T;
 Assemble T labels on the original sentence
End

Figure 2: Tree SRL system pseudo code

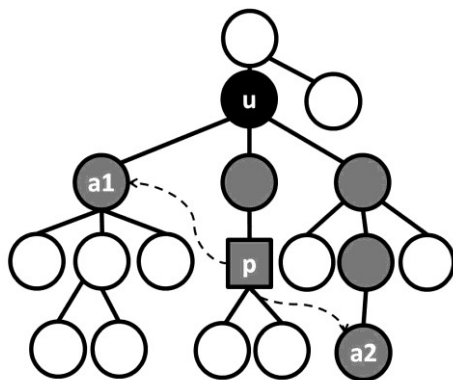


Figure 3: Sub-tree extraction sample. Assuming that “p” (the square node) is a predicate node and the nodes “a1” and “a2” are its arguments (the arguments are defined by the semi-dotted arrows.), the sub-tree extracted from the above sentence will contain the nodes: “a1”, “a2”, “p”, all ancestors of “a1”, “a2” and “p” up to the first common one, in this case node “u”, which is also included in the sub-tree. All of the white nodes are not included in the sub-tree. The straight lines represent syntactic dependency relations.

5 Labelling

Suppose that in Figure 1, the bottom sentence is the query, where the grey shadow contains the sub-tree to be labelled and the top sentence contains the sub-tree sample chosen to label the query. Then, an alignment between the sample sub-tree and the query sub-tree suggests labelling the query sub-tree with A1, A2 and A3, where the first two labels are right but the last label, A4, is predicted as A3, so it is wrong.

Input: A sub-tree to be labelled
Input: list of alignments sorted by ascending tree distance
Output: labelled sub-tree
foreach argument(a) in T **do**
 foreach alignment (ali) in the sorted list **do**
 if there is a semantic relation
 (ali.function(p),ali.function(a))
 Then break loop;
 end
 end
 label relation p-a with the label of the
 relation (ali.function(p),ali.function(a));
end
p is the node predicate.
a is a node argument.
ali is an alignment between the sub-tree that
 has to be labelled and a sub-tree in the train-
 ing dataset.
 The method function is explained in Figure 3.

Figure 4: Labelling a relation. (approach A)

Input: T: tree structure labelled in post order traversal
Input: L: list of nodes to be on the sub-tree in post order traversal
Output: T: Sub-Tree
foreach node x in the list **do**
 mark x as part of the sub-tree;
end
while L contains more than 2 unique values **do**
 [minValue , position]=min(L);
 Value = parent(minValue);
 Mark value as part of the sub-tree;
 L[position] = value;
end
 Remove all nodes that are not marked as part of the sub-tree;

Figure 5: Sub-tree extraction

It is not necessary to label a whole sub-tree (query) using just a single sub-tree sample. However, if the whole query is labelled using a single answer sample, the prediction is guaranteed to be consistent (no repeated argument labels).

Some possible ways to label the semantic relation using a sorted list of alignments (with each sub-tree of the training data set) is discussed ahead. Each sub-tree contains one predicate and several semantic relations, one for each argument node.

5.1 Treating relations independently

In this sub-section, the neighbouring sub-trees for one relation of a sub-tree T refers to the near-

est sub-trees with which the match with T produces a match between two predicate nodes and two argument nodes. A label from the nearest neighbour(s) can be transferred to T for labelling the relation.

The current implementation (**Approach A**), described in more detail in Figure 4, labels a relation using the first nearest neighbour from a list ordered by ascending tree distance. If there are several nearest neighbours, the first one on the list is used. This is a naive implementation of the k-NN algorithm where in case of multiple nearest neighbours only one is used and the others get ignored.

A negative aspect of this strategy is that it can select a different sub-tree based on the input order. This makes the algorithm indeterministic. A way to make it deterministic can be by extending the parameter “k” in case of multiple cases at the same distance or a tie in the voting (**Approach B**).

5.2 Treating relations dependently

In this section, a sample refers to a sub-tree containing all arguments and its labels. The arguments for a certain predicate are related.

Some strategies can lead to non-consistent structures (core argument labels cannot appear twice in the same sub-tree). **Approach B** treats the relations independently. It does not have any mechanism to keep the consistency of the whole predicate structure.

Another way is to find a sample that contains enough information to label the whole sub-tree (**Approach C**). This approach always generates consistent structures. The limitation of this model is that the required sample may not exist or the tree distance may be very high, making those samples poor predictors. The implemented method (**Approach A**) indirectly attempts to find a training sample sub-tree which contains labels for all the arguments of the predicate.

It is expected for tree distances to be smaller than other sub-trees that do not have information to label all the desired relations.

The system tries to get a consistent structure using a simple algorithm. Only in the case when using the nearest tree does not lead to labelling the whole structure, labels are predicted using multiple samples, thereby, risking the structure consistency.

Future implementations will rank possible candidate labels for each relation (probably using multiple samples).

A “joint scoring algorithm”, which is commonly used (Marquez et al, 2008), can be applied for consistency checking after finding the rank probability for all the argument labels for the same predicate (**Approach D**).

6 Experiments: the matching cost

The cost of matching two nodes is crucial to the performance of the system. Different atomic measures (ways to measure the cost of matching two nodes) that were tested are explained ahead. Results for experiments using these atomic measures are given in Table 2.

6.1 Binary system

For Binary system, the atomic cost of matching two nodes is one if label POS or dependency relations are different, otherwise the cost is zero. The atomic cost of inserting or deleting a node is always one. Note that the measure is totally based on the syntactic structure (words are not used).

6.2 Ternary system

The next intuitive measure is how the system would perform in case of a ternary cost (ternary system). The atomic cost is half if POS or dependency relation is different, one if POS and dependency relation are different or zero in all other case. For this system, Table 2 shows a very similar accuracy to the binary one.

6.3 Hamming system

The atomic cost of matching two nodes is the sum of the following sub costs:

- 0.25 if POS is different.
- 0.25 if dependency relation is different.
- 0.25 if Lemma is different.
- 0.25 if one node is a predicate but the other is not or if both nodes are predicates but with different lemma.

The cost to create or delete nodes is one. Note that the sum of all costs cannot be greater than one.

6.4 Predicate match system

The analysis of results for the previous systems shows that the accuracy is higher for the sub-trees that are labelled using sub-trees with the same predicate node. Consequently, this strategy attempts to force the predicate to be the same.

In this system, the atomic cost of matching two nodes is the sum of the following sub costs:

- 0.3 if POS is different.
- 0.3 if dependency relation is different.
- 1 if one is a predicate and the other node is not or both nodes are predicates but with different lemma.

The cost to create or delete nodes is one.

6.5 Complex system

This strategy attempts to improve the accuracy by adding an extra label to the argument nodes and using it.

The atomic cost of matching two nodes is the sum of the following sub costs:

- 0.1 for each different label (dependency relation or POS or lemma).
- 0.1 for each pair of different labels (dependency relation or POS or lemma).
- 0.4 if one node is a predicate and the other is not.
- 0.4 if both nodes are predicates and lemma is different.
- 2 if one node is marked as an argument and the other is not or one node is marked as a predicate and the other is not.

The atomic cost of deleting or inserting a node is: two if the node is an argument or predicate node and one in any other case.

7 Results

Table 2 shows the accuracy of all the systems. The validation data set is added to the training data set when the system is labelling the evaluation data set. This is a common methodology followed in CoNLL2009 (Li et al, 2009).

System	Evaluation	Development
Binary	64.36%	61.12%
Ternary	64.88%	61.28%
Hamming	78.01%	
Predicate Match	76.98%	
Complex	78.98%	

Table 2: System accuracy

Accuracy is measured as the percentage of semantic labels correctly predicted.

The implementation of the Tree SRL system takes several days to run a single experiment. It makes non viable the idea of using the development data set for adjusting parameters and that is why, for the last three systems (Hamming, Predicate Match and Complex), the accuracy over the development data set is not measured. The same reason supports adding the development data set

to the training data set without over fitting the system, because the development data set is not really used for adjusting parameters.

However, the observations of the system on the development data set shows:

1. If the complexity gets increased (Ternary), the number of cases having the multiple nearest sub-trees gets reduced.
2. The output of the system only contains five per cent of inconsistent structures (Binary and Ternary), which is lower than expected. 0.5% of inconsistent sub-trees were detected in the training data-set.
3. Higher accuracy for the relations where a sub-tree is labelled using a sub-tree sample which has the same predicate node. This has led to the design of the “predicate match” and the “complex” systems.
4. Some sub-trees are very small (just one node). This resulted in low accuracy for they predicted labels due to multiple nearest neighbours.

It is surprising that the hamming measure reaches higher accuracy than the “predicate match”, which uses more information, and is also surprising that the accuracies for “Hamming”, “Predicate Match” and “Complex” systems are very similar.

The CoNLL-2009 SRL shared task was evaluated on multiple languages: Catalan, Chinese, Czech, English, German, Japanese and Spanish. Some results for those languages using “Tree SRL System Binary” are shown in Table 3.

Language	Accuracy on evaluation	Training data set size in Mb
English	64.36%	56
Spanish	57.86%	46
Catalan	58.49%	43
Japanese	50.71%	8
German	These languages had been excluded from the experiments because some of the sentences did not follow a dependency tree structure.	
Czech		
Chinese		

Table 3: Accuracy for other languages (Binary system)

The accuracy results for multiple languages suggest that the size of the corpora has a strong influence on the results of the system performance.

The results are not comparable with the rest of the CoNLL-2009 systems because the task is different. This system does not identify arguments and does not perform predicate sense disambiguation.

8 Conclusion

The tree distance algorithm has been applied successfully to build a SRL system. Future work will focus on improving the performance of the system by: a) trying to extend the sub-trees which will contain more contextual information, b) using different approaches to label semantic relations discussed in Section 5. Also, the system will be expanded to identify arguments using a tree distance algorithm.

Evaluating the task of identifying the arguments and labelling the relations separately will assist in determining which systems to combine to create an hybrid system with better performance.

Acknowledgments

This research is supported by the Science Foundation Ireland (Grant 07/CE/I1142) as part of the Centre for Next Generation Localisation (www.cngl.ie) at Trinity College Dublin.

Thanks are due to Dr Martin Emms for his support on the development of this project.

References

- Martin Emms. 2006. Variants of Tree Similarity in a Question Answering Task. *In Proceedings of the Workshop on Linguistic Distances, held in conjunction with COLING 2006*, 100–108, Sydney, Australia, Association for Computational Linguistics.
- Jan Hajič, Massimiliano Ciaramita, Richard Johansson, Daisuke Kawahara, Maria Antonia Martí, Luis Màrquez, Adam Meyers, Joakim Nivre, Sebastian Padó, Jan Štěpánek, Pavel Stravnák, Mihai Surdeanu, Nianwen Xue and Yi Zhang. 2009. The CoNLL-2009 shared task: syntactic and semantic dependencies in multiple languages. *In CoNLL '09: Proceedings of the Thirteenth Conference on Computational Natural Language Learning* (pp. 1-18). Morristown, NJ, USA: Association for Computational Linguistics.
- Seokhwan Kim, Minwoo Jeong and Gary Geunbae Lee. 2009. A Local Tree Alignment-based Soft Pattern Matching Approach for Information Extraction. *Proceedings of NAAACL HLT*, 169-172. Boulder, Colorado, June 2009
- Milen Kouylekov and Bernardo Magnini. 2005. Recognizing textual entailment with tree edit distance algorithms. *In Recognizing Textual Entailment* (pp. 17-20). Southampton, U.K.
- Baoli Li, Martin Emms, Saturnino Luz and Carl Vogel. 2009. Exploring multilingual semantic role labeling. *In CoNLL '09: Proceedings of the Thirteenth Conference on Computational Natural Language Learning* (pp. 73-78). Morristown, NJ, USA: Association for Computational Linguistics.
- Mitchell Marcus, Beatrice Santorini and Mary Ann Marcinkiewicz. 1994. Building a large annotated corpus of English: The Penn Treebank. *Computational linguistics*, 19(2), 313–330.
- Alessandro Moschitti, Daniele Pighin and Roberto Basili. 2008. Tree kernels for semantic role labeling. *Computational Linguistics*, 34(2), 193-224. Cambridge, MA, USA: MIT Press.
- Lluís Màrquez, Xavier Carreras, Kenneth C. Litkowski and Suzanne Stevenson. 2008. Semantic Role Labeling: An Introduction to the Special Issue. *Computational Linguistics*, 34(2), 145-159.
- Martha Palmer, Paul Kingsbury and Daniel Gildea. 2005. The Proposition Bank: An Annotated Corpus of Semantic Roles. *Computational Linguistics*, 31(1), 71-106.
- Vasin Punyakanok, Dan Roth and Wen-tau Yih. 2004. Mapping dependencies trees: An application to question answering. *In Proceedings of AI&Math 2004* (pp. 1-10). Ford.
- Dennis Shasha and Kaizhong Zhang. 1990. Fast algorithms for the unit cost editing distance between trees. *J. Algorithms*, 11(4), 581-621. Duluth, MN, USA: Academic Press, Inc.
- Kuo-Chung Tai. 1979. The Tree-to-Tree Correction Problem. *J. ACM*, 26(3), 422-433. New York, NY, USA: ACM.
- Kaizhong Zhang and Dennis Shasha. 1989. Simple fast algorithms for the editing distance between trees and related problems. *SIAM J. Comput.*, 18(6), 1245-1262. Philadelphia, PA, USA: Society for Industrial and Applied Mathematics.

Automatic Sanskrit Segmentizer Using Finite State Transducers

Vipul Mittal

Language Technologies Research Center, IIT-H,
Gachibowli, Hyderabad, India.

vipulmittal@research.iiit.ac.in

Abstract

In this paper, we propose a novel method for automatic segmentation of a Sanskrit string into different words. The input for our segmentizer is a Sanskrit string either encoded as a Unicode string or as a Roman transliterated string and the output is a set of possible splits with weights associated with each of them. We followed two different approaches to segment a Sanskrit text using *sandhi*¹ rules extracted from a parallel corpus of manually sandhi split text. While the first approach augments the finite state transducer used to analyze Sanskrit morphology and traverse it to segment a word, the second approach generates all possible segmentations and validates each constituent using a morph analyzer.

1 Introduction

Sanskrit has a rich tradition of oral transmission of texts and this process causes the text to undergo euphonic changes at the word boundaries. In oral transmission, the text is predominantly spoken as a continuous speech. However, continuous speech makes the text ambiguous. To overcome this problem, there is also a tradition of reciting the pada-pāṭha (recitation of words) in addition to the recitation of a saṃhitā (a continuous sandhied text). In the written form, because of the dominance of oral transmission, the text is written as a continuous string of letters rather than a sequence of words. Thus, the Sanskrit texts consist of a very

¹*Sandhi* means euphony transformation of words when they are consecutively pronounced. Typically when a word w_1 is followed by a word w_2 , some terminal segment of w_1 merges with some initial segment of w_2 to be replaced by a “smoothed” phonetic interpolation, corresponding to minimizing the energy necessary to reconfigure the vocal organs at the juncture between the words.

long sequence of phonemes, with the word boundaries having undergone euphonic changes. This makes it difficult to split a continuous string into words and process the text automatically.

Sanskrit words are mostly analyzed by building a finite state transducer (Beesley, 1998). In the first approach, this transducer was modified by linking the final states to appropriate intermediate states incorporating the sandhi rules. This approach then allows one to traverse the string from left to right and generate all and only possible splits that are morphologically valid. The second approach is very closely based on the *Optimality Theory* (Prince and Smolensky, 1993) where we generate all the possible splits for a word and validate each using a morphological analyzer. We use one of the fastest morphological analyzers available viz. the one developed by *Apertium* group². The splits that are not validated are pruned out. Based on the number of times the first answer is correct, we achieved an accuracy of around 92% using the second approach while the first approach performed with around 71% accuracy.

2 Issues involved in Sanskrit Processing

The segmentizer is an important component of an NLP system. Especially, languages such as Chinese (Badino, 2004), Japanese, Thai (Haruechaiyasak, 2008) or Vietnamese (Thang et al., 2008) which do not mark word boundaries explicitly or highly agglutinative languages like Turkish need segmentizers. In all these languages, there are no explicit delimiters to specify the word boundaries. In Thai, each syllable is transcribed using several characters and there is no space in the text between syllables. So the problem of segmentation is basically twofold: (1) syllable segmentation followed by (2) word segmentation itself. A sentence in these languages

²<http://wiki.apertium.org/wiki/Ittoolbox>; It processes around 50,000 words per sec.

is segmented by predicting the word boundaries, where euphonic changes do not occur across the word boundaries and it is more like mere concatenation of words. So the task here is just to choose between various combinations of the words in a sentence.

However, in Sanskrit, euphonic changes occur across word boundaries leading to addition and deletion of some original part of the combining words. These euphonic changes in Sanskrit introduce non-determinism in the segmentation. This makes the segmentation process in Sanskrit more complex than in Chinese or Japanese. In case of highly agglutinative languages like Turkish, the components are related to each other semantically involving dependency analysis. Whereas in Sanskrit, only the compounds involve a certain level of dependency analysis, while sandhi is just gluing of words together, without the need for words to be related semantically. For example, consider the following part of a verse,

San: *nāradam paripapraccha*
vālmīkirmunipuṅgavam
 gloss: to_the_Narada asked Valmiki-
 to_the_wisest_among_sages
 Eng: Valmiki asked the Narada, the wisest among
 the sages.

In the above verse, the words *vālmīkiḥ* and *munipuṅgavam* (wisest among the sages - an adjective of Narada) are not related semantically, but still undergo euphonic change and are glued together as *vālmīkirmunipuṅgavam*.

Further, the split need not be unique. Here is an example, where a string *māturājñāmparipālaya* may be decomposed in two different ways after undergoing euphonic changes across word boundaries.

- *mātuh ājñām paripālaya* (obey the order of mother) and,
- *mā āturājñām paripālaya* (do not obey the order of the diseased).

There are special cases where the sandhied forms are not necessarily written together. In such cases, the white space that physically marks the boundary of the words, logically refers to a single sandhied form. Thus, the white space is deceptive, and if treated as a word boundary, the morphological analyzer fails to recognize the

word. For example, consider

śrutvā ca nārado vacaḥ.

In this example, the space between *śrutvā* and *ca* represent a proper word boundary and the word *śrutvā* is recognized by the morphological analyzer whereas the space between *nārado* and *vacaḥ* does not mark the word boundary making it deceptive. Because of the word *vacaḥ*, *nāradaḥ* has undergone a phonetic change and is rendered as *nārado*. In unsandhied form, it would be written as,

San: *śrutvā ca nāradaḥ vacaḥ.*
 gloss: after listening and Narada's speech
 Eng: And after listening to Narada's speech

The third factor aggravating Sanskrit segmentation is productive compound formation. Unlike English, where either the components of a compound are written as distinct words or are separated by a hyphen, the components of compounds in Sanskrit are always written together. Moreover, before these components are joined, they undergo the euphonic changes. The components of a compound typically do not carry inflection or in other words they are the bound morphemes used only in compounds. This forces a need of a special module to recognize compounds.

Assuming that a sandhi handler to handle the sandhi involving spaces is available and a bound morpheme recognizer is available, we discuss the development of sandhi splitter or a segmentizer that splits a continuous string of letters into meaningful words. To illustrate this point, we give an example. Consider the text,

śrutvā caitatrilokajñō vālmīkernārado vacaḥ.

We assume that the sandhi handler handling the sandhi involving spaces is available and it splits the above string as,

śrutvā caitatrilokajñāḥ vālmīkernāradaḥ vacaḥ.

The sandhi splitter or segmentizer is supposed to split this into

śrutvā ca etat triloka-jñāḥ vālmīkeḥ nāradaḥ vacaḥ.

This presupposes the availability of rules corresponding to euphonic changes and a good coverage morphological analyzer that can also analyze the bound morphemes in compounds.

A segmentizer for Sanskrit developed by Huet (Huet, 2009), decorates the final states of its finite state transducer handling Sanskrit morphology with the possible sandhi rules. However, it is still not clear how one can prioritize various splits with this approach. Further, this system in current state demands some more work before the sandhi splitter of this system can be used as a standalone system allowing plugging in of different morphological analyzers. With a variety of morphological analyzers being developed by various researchers³, at times with complementary abilities, it would be worth to experiment with various morphological analyzers for splitting a sandhied text. Hence, we thought of exploring other alternatives and present two approaches, both of which assume the existence of a good coverage morphological analyzer. Before we describe our approaches, we first define the scoring matrix used to prioritize various analyses followed by the baseline system.

3 Scoring Matrix

Just as in the case of any NLP systems, with the sandhi splitter being no exception, it is always desirable to produce the most likely output when a machine produces multiple outputs. To ensure that the correct output is not deeply buried down the pile of incorrect answers, it is natural to prioritize solutions based on some frequencies. A Parallel corpus of Sanskrit text in sandhied and sandhi split form is being developed as a part of the Consortium project in India. The corpus contains texts from various fields ranging from children stories, dramas to Ayurveda texts. Around 100K words of such a parallel corpus is available from which around 25,000 parallel strings of unsandhied and corresponding sandhied texts were extracted. The same corpus was also used to extract a total of 2650 sandhi rules including the cases of mere concatenation, and the frequency distribution of these sandhi rules. Each sandhi rule is a triple (x, y, z)

³<http://sanskrit.uohyd.ernet.in>,
<http://www.sanskritlibrary.org>, <http://sanskrit.jnu.ernet.in>

where y is the last letter of the first primitive, z is the first letter of the second primitive, and x is the letter sequence created by euphonic combination. We define the estimated probability of the occurrence of a sandhi rule as follows:

Let R_i denote the i^{th} rule with f_{R_i} as the frequency of occurrence in the manually split parallel text. The probability of rule R_i is:

$$P_{R_i} = \frac{f_{R_i}}{\sum_{i=1}^n f_{R_i}}$$

where n denotes the total number of sandhi rules found in the corpus.

Let a word be split into a candidate S_j with k constituents as $\langle c_1, c_2, \dots, c_k \rangle$ by applying $k - 1$ sandhi rules $\langle R_1, R_2, \dots, R_{k-1} \rangle$ in between the constituents. It should be noted here that the rules R_1, \dots, R_{k-1} and the constituents c_1, \dots, c_k are interdependent since a different rule sequence will result in a different constituents sequence. Also, except c_1 and c_k , all intermediate constituents take part in two segmentations, one as the right word and one as the left.

The weight of the split S_j is defined as:

$$W_{S_j} = \frac{\prod_{x=1}^{k-1} (P_{c_x} + P_{c_{x+1}}) * P_{R_x}}{k}$$

where P_{c_x} is the probability of occurrence of the word c_x in the corpus. The factor of k was introduced to give more preference to the split with less number of segments than the one with more segments.

4 Baseline System

We define our own baseline system which assumes that each Sanskrit word can be segmented only in two constituents. A word is traversed from left to right and is segmented by applying the first applicable rule provided both the constituents are valid morphs. Using the 2,650 rules, on a test data of 2,510 words parallel corpus, the baseline performance of the system was around 52.7% where the first answer was correct.

5 Two Approaches

We now present the two approaches we explored for sandhi splitting.

5.1 Augmenting FST with Sandhi rules

In this approach, we build an FST, using OpenFst (Allauzen et al., 2007) toolkit, incorporating

sandhi rules in the FST itself and traverse it to find the sandhi splittings.

We illustrate the augmentation of a sandhi rule with an example. Let the two strings be $xaXi$ (dadhi)⁴ and $awra$ (atra). The initial FST without considering any sandhi rules is shown in Figure 1.

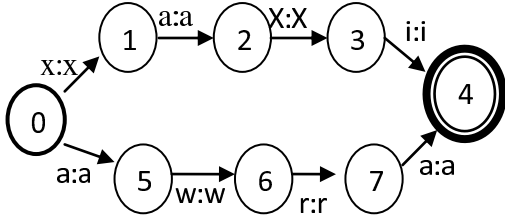


Figure 1: Initial FST accepting only two words $xaXi$ and $awra$.

As the figure depicts, 0 is the start state and 4 is the final state. Each transition is a 4-tuple $\langle c, n, i, o \rangle$ where c is current state, n is the next state, i is the input symbol and o is the output. The FST marks word boundaries by flushing out certain features about the words whenever it encounters a valid word. Multiple features are separated by a '|'. E.g., the output for $xaXi$ is $lc,s|vc,s$ and for $awra$ it is vc,s where lc,s stands for *locative, singular* and vc,s is *vocative, singular*. The FST in Figure 1 recognize exactly two words $xaXi$ and $awra$.

One of the sandhi rule states that $i+a \rightarrow ya$ which will be represented as a triple (ya, i, a) . Applying the sandhi rule, we get: $xaXi + awra \rightarrow xaXyawra$. After adding this sandhi rule to the FST, we get the modified FST that is represented in Figure 2.

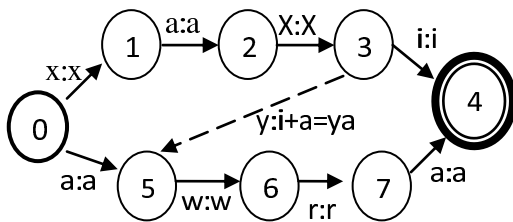


Figure 2: Modified FST after inserting the rule. - - - indicates the newly added transition.

Here, a transition arc is added depicting the rule which says that on receiving an input symbol ya at state 3, go to state 5 with an output $i+a \rightarrow ya$.

⁴A Roman transliteration scheme called WX transliteration is used, which is one-to-one phoneme level representation of Devanāgarī script.

Thus the new FST accepts $xaXyawra$ in addition to $xaXi$ and $awra$.

Thus, we see that the original transducer gets modified with all possible transitions at the end of a final phoneme, and hence, also explodes the number of transitions leading to a complex transducer.

The basic outline of the algorithm to split the given string into sub-strings is:

Algorithm 1 To split a string into sub-strings

- 1: Let the FST for morphology be f .
 - 2: Add sandhi rules to the final states of f linking them to the intermediary states to get f' .
 - 3: Traverse f' to find all possible splits for a word. If a sandhi rule is encountered, split the word and continue with the remaining part.
 - 4: Calculate the weights of the possible outputs with the formula discussed in section 3.
-

The pseudo-code of the algorithm used to insert sandhi rules in the FST is illustrated here:

Algorithm 2 To insert sandhi rules in the FST

- 1: I = Input Symbol; X = last character of the result of the rule.
 - 2: **for** each transition in the FST transition table **do**
 - 3: **if** next state is a final state **then**
 - 4: **for** all rules where I is the last character of first word **do**
 - 5: S = next state from the start state on encountering X;
 - 6: Y = first character of the result of the rule;
 - 7: transition T = current state, S, Y, rule;
 - 8: Add T into the FST;
 - 9: **end for**
 - 10: **end if**
 - 11: **end for**
-

The main problem with this approach is that every finite state can have as many transitions as the number of euphonic rules resulting in phoneme change. This increases the size of the FST considerably. It should be noted that, we have not included the cases, where there is just a concatenation. In such cases, if the input string is not exhausted, but the current state is a final state, we go back to the start state with the remaining string as the input.

5.1.1 Results

The performance of this system measured in terms of the number of times the highest ranked segmentation is correct, with around 500 sandhi rules, and only noun morphology tested on the same test data used for testing baseline system gave the following rank-wise distribution presented in Table 1.

Rank	% of output
1	71.2509
2	5.64543
3	3.85324
4	3.35651
5	1.56123
>5	14.33268

Table 1: Rank-wise Distribution for Approach-1.

The system was slow consuming, on an average, around 10 seconds per string of 15 letters.⁵

With the increase in the sandhi rules, though system's performance was better, it slowed down the system further. Moreover, this was tested only with the inflection morphology of nouns. The verb inflection morphology and the derivational morphology were not used at all. Since, the system is supposed to be part of a real time application viz. machine translation, we decided to explore other possibilities.

5.2 Approach based on Optimality Theory

Our second approach follows optimality theory(OT) which proposes that the observed forms of a language are a result of the interaction between the conflicting constraints. The three basic components of the theory are:

1. GEN - generates all possible outputs, or candidates.
2. CON - provides the criteria and the constraints that will be used to decide between candidates.
3. EVAL - chooses the optimal candidate based on the conflicts on the constraints.

OT assumes that these components are universal and the grammars differ in the way they rank the universal constraint set, CON. The grammar of

⁵Tested on a system with 2.93GHz Core 2 Duo processor and 2GB RAM

each language ranks the constraints in some dominance order in such a way that every constraint must have outperformed every lower ranked constraint. Thus a candidate A is optimal if it performs better than some other candidate B on a higher ranking constraint even if A has more violations of a lower ranked constraint than B.

The GEN function produces every possible segmentation by applying the rules wherever applicable. The rules tokenize the input surface form into individual constituents. This might contain some insignificant words that will be eventually pruned out using the morphological analyser in the EVAL function thus leaving the winning candidate. Therefore, the approach followed is very closely based on optimality theory. The morph analyser has no role in the generation of the candidates but only during their validation thus composing the back-end of the segmentizer. In original OT, the winning candidate need not satisfy all the constraints but it must outperform all the other candidates on some higher ranked constraint. While in our scenario, the winning candidate must satisfy all the constraints and therefore there could be more than one winning candidates.

Currently we are applying only two constraints. We are planning to introduce some more constraints. The constraints applied are:

- C1 : All the constituents of a split must be valid morphs.
- C2 : Select the split with maximum weight, as defined in section 3.

The basic outline of the algorithm is:

-
- 1: Recursively break a word at every possible position applying a sandhi rule and generate all possible candidates for the input.
 - 2: Pass the constituents of all the candidates through the morph analyzer.
 - 3: Declare the candidate as a valid candidate, if all its constituents are recognized by the morphological analyzer.
 - 4: Assign weights to the accepted candidates and sort them based on the weights.
 - 5: The optimal solution will be the one with the highest salience.
-

5.2.1 Results

The current morphological analyzer can recognize around 140 million words. Using the 2650 rules

and the same test data used for previous approach, we obtained the following results:

- Almost 93% of the times, the highest ranked segmentation is correct. And in almost 98% of the cases, the correct split was among the top 3 possible splits.
- The system consumes around 0.04 seconds per string of 15 letters on an average.

The complete rank wise distribution is given in Table 2.

Rank	% of output	
	Approach-1	Approach-2
1	71.2509	92.8771
2	5.64543	5.44693
3	3.85324	1.07076
4	3.35651	0.41899
5	1.56123	0.09311
>5	14.33268	0.0931

Table 2: Complete rank-wise Distribution.

6 Conclusion

We presented two methods to automatically segment a Sanskrit word into its morphologically valid constituents. Though both the approaches outperformed the baseline system, the approach that is close to optimality theory gives better results both in terms of time consumption and segmentations. The results are encouraging. But the real test of this system will be when it is integrated with some real application such as a machine translation system. This sandhi splitter being modular, wherein one can plug in different morphological analyzer and different set of sandhi rules, the splitter can also be used for segmentization of other languages.

Future Work The major task would be to explore ways to shift rank 2 and rank 3 segmentations more towards rank 1. We are also exploring the possibility of including some semantic information about the words while defining weights. The sandhi with white spaces also needs to be handled.

Acknowledgments

I would like to express my gratitude to Amba Kulkarni and Rajeev Sangal for their guidance and support.

References

- Akshar Bharati, Amba P. Kulkarni, and V Sheeba. 2006. *Building a wide coverage Sanskrit morphological analyzer: A practical approach*. The First National Symposium on Modelling and Shallow Parsing of Indian Languages, IIT-Bombay.
- Alan Prince and Paul Smolensky. 1993. *Optimality Theory: Constraint Interaction in Generative Grammar*. RuCCS Technical Report 2 at Center for Cognitive Science, Rutgers University, Piscataway.
- Amba Kulkarni and Devanand Shukla. 2009. *Sanskrit Morphological analyzer: Some Issues*. To appear in Bh.K Festschrift volume by LSI.
- Choochart Haruechaiyasak, Sarawoot Kongyoung, and Matthew N. Dailey. 2008. *A Comparative Study on Thai Word Segmentation Approaches*. ECTI-CON, Krabi.
- Cyril Allauzen, Michael Riley, Johan Schalkwyk, Wojciech Skut, and Mehryar Mohri. 2007. *OpenFst: A General and Efficient Weighted Finite-State Transducer Library*. CIAA'07, Prague, Czech Republic.
- Deniz Yuret and Ergun Biçici. 2009. *Modeling Morphologically Rich Languages Using Split Words and Unstructured Dependencies*. ACL-IJCNLP'09, Singapore.
- DINH Q. Thang, LE H. Phuong, NGUYEN T. M. Huyen, NGUYEN C. Tu, Mathias Rossignol, and VU X. Luong. 2008. *Word Segmentation of Vietnamese Texts: a Comparison of Approaches*. LREC'08, Marrakech, Morocco.
- G rard Huet. 2009. *Formal structure of Sanskrit text: Requirements analysis for a mechanical Sanskrit processor*. Sanskrit Computational Linguistics 1 & 2, pages 266-277, Springer-Verlag LNAI 5402.
- John C. J. Hoeks and Petra Hendriks. 2005. *Optimality Theory and Human Sentence Processing: The Case of Coordination*. Proceedings of the 27th Annual Meeting of the Cognitive Science Society, Erlbaum, Mahwah, NJ, pp. 959-964.
- Kenneth R. Beesley. 1998. *Arabic morphology using only finite-state operations* Proceedings of the ACL Workshop on Computational Approaches to Semitic Languages, Montr al, Qu bec.
- Leonardo Badino. 2004. *Chinese Text Word-Segmentation Considering Semantic Links among Sentences*. INTERSPEECH 2004 - ICSLP, Jeju, Korea.

Adapting Self-training for Semantic Role Labeling

Rasoul Samad Zadeh Kaljahi
FCSIT, University of Malaya
50406, Kuala Lumpur, Malaysia.
rsk7945@perdana.um.edu.my

Abstract

Supervised semantic role labeling (SRL) systems trained on hand-crafted annotated corpora have recently achieved state-of-the-art performance. However, creating such corpora is tedious and costly, with the resulting corpora not sufficiently representative of the language. This paper describes a part of an ongoing work on applying bootstrapping methods to SRL to deal with this problem. Previous work shows that, due to the complexity of SRL, this task is not straight forward. One major difficulty is the propagation of classification noise into the successive iterations. We address this problem by employing *balancing* and *preselection* methods for self-training, as a bootstrapping algorithm. The proposed methods could achieve improvement over the base line, which do not use these methods.

1 Introduction

Semantic role labeling has been an active research field of computational linguistics since its introduction by Gildea and Jurafsky (2002). It reveals the event structure encoded in the sentence, which is useful for other NLP tasks or applications such as information extraction, question answering, and machine translation (Surdeanu et al., 2003). Several CoNLL shared tasks (Carreras and Marquez, 2005; Surdeanu et al., 2008) dedicated to semantic role labeling affirm the increasing attention to this field.

One important supportive factor of studying *supervised* statistical SRL has been the existence of hand-annotated semantic corpora for training SRL systems. *FrameNet* (Baker et al., 1998) was the first such resource, which made the emergence of this research field possible by the seminal work of Gildea and Jurafsky (2002). However, this corpus only exemplifies the semantic role assignment by selecting some illustrative examples for annotation. This questions its suitability

for statistical learning. *Propbank* was started by Kingsbury and Palmer (2002) aiming at developing a more representative resource of English, appropriate for statistical SRL study.

Propbank has been used as the learning framework by the majority of SRL work and competitions like CoNLL shared tasks. However, it only covers the newswire text from a specific genre and also deals only with verb predicates.

All state-of-the-art SRL systems show a dramatic drop in performance when tested on a new text domain (Punyakanok et al., 2008). This evince the infeasibility of building a comprehensive hand-crafted corpus of natural language useful for training a robust semantic role labeler.

A possible relief for this problem is the utility of *semi-supervised* learning methods along with the existence of huge amount of natural language text available at a low cost. Semi-supervised methods compensate the scarcity of labeled data by utilizing an additional and much larger amount of unlabeled data via a variety of algorithms.

Self-training (Yarowsky, 1995) is a semi-supervised algorithm which has been well studied in the NLP area and gained promising result. It iteratively extend its training set by labeling the unlabeled data using a base classifier trained on the labeled data. Although the algorithm is theoretically straightforward, it involves a large number of parameters, highly influenced by the specifications of the underlying task. Thus to achieve the best-performing parameter set or even to investigate the usefulness of these algorithms for a learning task such as SRL, a thorough experiment is required. This work investigates its application to the SRL problem.

2 Related Work

The algorithm proposed by Yarowsky (1995) for the problem of word sense disambiguation has been cited as the origination of self-training. In that work, he bootstrapped a ruleset from a

Feature Name	Description
Phrase Type	Phrase type of the constituent
Position+Predicate Voice	Concatenation of constituent position relative to verb and verb voice
Predicate Lemma	Lemma of the predicate
Predicate POS	POS tag of the predicate
Path	Tree path of non-terminals from predicate to constituent
Head Word Lemma	Lemma of the head word of the constituent
Content Word Lemma	Lemma of the content word of the constituent
Head Word POS	POS tag of the head word of the constituent
Content Word POS	POS tag of the head word of the constituent
Governing Category	The first VP or S ancestor of a NP constituent
Predicate Subcategorization	Rule expanding the predicate's parent
Constituent Subcategorization *	Rule expanding the constituent's parent
Clause+VP+NP Count in Path	Number of clauses, NPs and VPs in the path
Constituent and Predicate Distance	Number of words between constituent and predicate
Compound Verb Identifier	Verb predicate structure type: simple, compound, or discontinuous compound
Head Word Location in Constituent *	Location of head word inside the constituent based on the number of words in its right and left

Table 1: Features

small number of seed words extracted from an online dictionary using a corpus of unannotated English text and gained a comparable accuracy to fully supervised approaches.

Subsequently, several studies applied the algorithm to other domains of NLP. Reference resolution (Ng and Cardie 2003), POS tagging (Clark et al., 2003), and parsing (McClosky et al., 2006) were shown to be benefited from self-training. These studies show that the performance of self-training is tied with its several parameters and the specifications of the underlying task.

In SRL field, He and Gildea (2006) used self-training to address the problem of unseen frames when using FrameNet as the underlying training corpus. They generalized FrameNet frame ele-

ments to 15 thematic roles to control the complexity of the process. The improvement gained by the progress of self-training was small and inconsistent. They reported that the NULL label (non-argument) had often dominated other labels in the examples added to the training set.

Lee et al. (2007) attacked another SRL learning problem using self-training. Using Propbank instead of FrameNet, they aimed at increasing the performance of supervised SRL system by exploiting a large amount of unlabeled data (about 7 times more than labeled data). The algorithm variation was similar to that of He and Gildea (2006), but it only dealt with core arguments of the Propbank. They achieved a minor improvement too and credited it to the relatively poor performance of their base classifier and the insufficiency of the unlabeled data.

3 SRL System

To have enough control over entire the system and thus a flexible experimental framework, we developed our own SRL system instead of using a third-party system. The system works with PropBank-style annotation and is described here.

Syntactic Formalism: A Penn Treebank constituent-based approach for SRL is taken. Syntactic parse trees are produced by the reranking parser of Charniak and Johnson (2005).

Architecture: A two-stage pipeline architecture is used, where in the first stage less-probable argument candidates (samples) in the parse tree are pruned, and in the next stage, final arguments are identified and assigned a semantic role. However, for unlabeled data, a preprocessing stage identifies the verb predicates based on the POS tag assigned by the parser. The joint argument identification and classification is chosen to decrease the complexity of self-training process.

Features: Features are listed in table 1. We tried to avoid features like named entity tags to less depend on extra annotation. Features marked with * are used in addition to common features in the literature, due to their impact on the performance in feature selection process.

Classifier: We chose a *Maximum Entropy* classifier for its efficient training time and also its built-in multi-classification capability. Moreover, the probability score that it assigns to labels is useful in selection process in self-training. The *Maxent Toolkit*¹ was used for this purpose.

¹http://homepages.inf.ed.ac.uk/lzhang10/maxent_toolkit.html

- 1- Add the seed example set L to currently empty training set T .
- 2- Train the base classifier C with training set T .
- 3- Iterate the following steps until the stop criterion S is met.
 - a- **Select** p examples from U into pool P .
 - b- Label pool P with classifier C
 - c- **Select** n labeled examples with the highest confidence score whose score meets a certain threshold t and **add** to training set T .
 - d- Retrain the classifier C with new training set.

Figure 1: Self-training Algorithm

	WSJ Test			Brown Test		
	P	R	F1	P	R	F1
Cur	77.43	68.15	72.50	69.14	57.01	62.49
Pun	82.28	76.78	79.44	73.38	62.93	67.75

Table 2: Performances of the current system (Cur) and the state-of-the-art (Punyakano et al., 2008)

4 Self-training

4.1 The Algorithm

While the general theme of the self-training algorithm is almost identical in different implementations, variations of it are developed based on the characteristics of the task in hand, mainly by customizing several involved parameters. Figure 1 shows the algorithm with highlighted parameters.

The size of seed labeled data set L and unlabeled data U , and their ratio are the fundamental parameters in any semi-supervised learning. The data used in this work is explained in section 5.1.

In addition to performance, efficiency of the classifier (C) is important for self-training, which is computationally expensive. Our classifier is a compromise between performance and efficiency. Table 2 shows its performance compared to the state-of-the-art (Punyakano et al. 2008) when trained on the whole labeled training set.

Stop criterion (S) can be set to a predetermined number of iterations, finishing all of the unlabeled data, or convergence of the process in terms of improvement. We use the second option for all experiments here.

In each iteration, one can label entire the unlabeled data or only a portion of it. In the latter case, a number of unlabeled examples (p) are

selected and loaded into a *pool* (P). The selection can be based on a specific strategy, known as *preselection* (Abney, 2008) or simply done according to the original order of the unlabeled data. We investigate preselection in this work.

After labeling the p unlabeled data, training set is augmented by adding the newly labeled data. Two main parameters are involved in this step: *selection* of labeled examples to be added to training set and *addition* of them to that set.

Selection is the crucial point of self-training, in which the propagation of labeling noise into upcoming iterations is the major concern. One can select all of labeled examples, but usually only a number of them (n), known as *growth size*, based on a quality measure is selected. This measure is often the confidence score assigned by the classifier. To prevent poor labelings diminishing the quality of training set, a threshold (t) is set on this confidence score. Selection is also influenced by other factors, one of which being the balance between selected labels, which is explored in this study and explained in detail in the section 4.3.

The selected labeled examples can be retained in unlabeled set to be labeled again in next iterations (*delibility*) or moved so that they are labeled only once (*indelibility*). We choose the second approach here.

4.2 Preselection

While using a pool can improve the efficiency of the self-training process, there can be two other motivations behind it, concerned with the performance of the process.

One idea is that when all data is labeled, since the growth size is often much smaller than the labeled size, a uniform set of examples preferred by the classifier is chosen in each iteration. This leads to a biased classifier like the one discussed in previous section. Limiting the labeling size to a pool and at the same time (pre)selecting divergence examples into it can remedy the problem.

The other motivation is originated from the fact that the base classifier is relatively weak due to small seed size, thus its predictions, as the measure of confidence in selection process, may not be reliable. Preselecting a set of unlabeled examples more probable to be correctly labeled by the classifier in initial steps seems to be a useful strategy against this fact.

We examine both ideas here, by a random preselection for the first case and a measure of simplicity for the second case. Random preselection is built into our system, since we use randomized

training data. As the measure of simplicity, we propose the number of samples extracted from each sentence; that is we sort unlabeled sentences in ascending order based on the number of samples and load the pool from the beginning.

4.3 Selection Balancing

Most of the previous self-training problems involve a binary classification. Semantic role labeling is a multi-class classification problem with an unbalanced distribution of classes in a given text. For example, the frequency of A1 as the most frequent role in CoNLL training set is 84,917, while the frequency of 21 roles is less than 20. The situation becomes worse when the dominant label NULL (for non-arguments) is added for argument identification purpose in a joint architecture. This biases the classifiers towards the frequent classes, and the impact is magnified as self-training proceeds.

In previous work, although they used a reduced set of roles (yet not balanced), He and Gildea (2006) and Lee et al. (2007), did not discriminate between roles when selecting high-confidence labeled samples. The former study reports that the majority of labels assigned to samples were NULL and argument labels appeared only in last iterations.

To attack this problem, we propose a natural way of balancing, in which instead of labeling and selection based on argument samples, we perform a sentence-based selection and labeling. The idea is that argument roles are distributed over the sentences. As the measure for selecting a labeled sentence, the average of the probabilities assigned by the classifier to all argument samples extracted from the sentence is used.

5 Experiments and Results

In these experiments, we target two main problems addressed by semi-supervised methods: the performance of the algorithm in exploiting unlabeled data when labeled data is scarce and the domain-generalizability of the algorithm by using an out-of-domain unlabeled data.

We use the CoNLL 2005 shared task data and setting for testing and evaluation purpose. The evaluation metrics include *precision*, *recall*, and their harmonic mean, *F1*.

5.1 The Data

The labeled data are selected from Propbank corpus prepared for CoNLL 2005 shared task. Our learning curve experiments on varying size

of labeled data shows that the steepest increase in F1 is achieved by 1/10th of CoNLL training data. Therefore, for training a base classifier as high-performance as possible, while simulating the labeled data scarcity with a reasonably small amount of it, 4000 sentence are selected randomly from the total 39,832 training sentences as seed data (L). These sentences contain 71,400 argument samples covering 38 semantic roles out of 52 roles present in the total training set.

We use one unlabeled training set (U) for in-domain and another for out-of-domain experiments. The former is the remaining portion of CoNLL training data and contains 35,832 sentences (698,567 samples). The out-of-domain set was extracted from Open American National Corpus² (OANC), a 14-million words multi-genre corpus of American English. The whole corpus was preprocessed to prune some problematic sentences. We also excluded the *biomed* section due to its large size to retain the domain balance of the data. Finally, 304,711 sentences with the length between 3 and 100 were parsed by the syntactic parser. Out of these, 35,832 sentences were randomly selected for the experiments reported here (832,795 samples).

Two points are worth noting about the results in advance. First, we do not exclude the argument roles not present in seed data when evaluating the results. Second, we observed that our predicate-identification method is not reliable, since it is solely based on POS tags assigned by parser which is error-prone. Experiments with gold predicates confirmed this conclusion.

5.2 The Effect of Balanced Selection

Figures 2 and 3 depict the results of using unbalanced and balanced selection with WSJ and OANC data respectively. To be comparable with previous work (He and Gildea, 2006), the growth size (n) for unbalanced method is 7000 samples and for balanced method is 350 sentences, since each sentence roughly contains 20 samples. A probability threshold (t) of 0.70 is used for both cases. The F1 of base classifier, best-performed classifier, and final classifier are marked.

When trained on WSJ unlabeled set, the balanced method outperforms the other in both WSJ (68.53 vs. 67.96) and Brown test sets (59.62 vs. 58.95). A two-tail t-test based on different random selection of training data confirms the statistical significance of this improvement at $p \leq 0.05$ level. Also, the self-training trend is

² <http://www.americannationalcorpus.org/OANC>

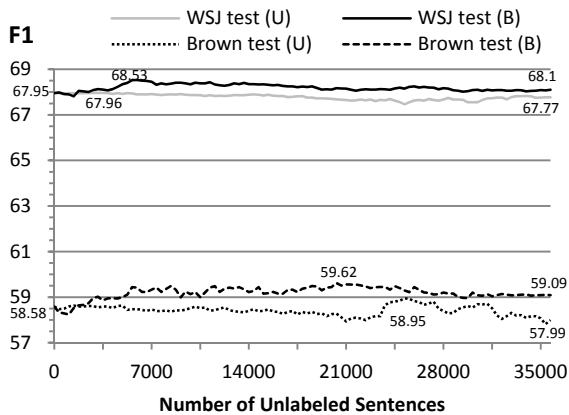


Figure 2: Balanced (B) and Unbalanced (U) Selection with WSJ Unlabeled Data

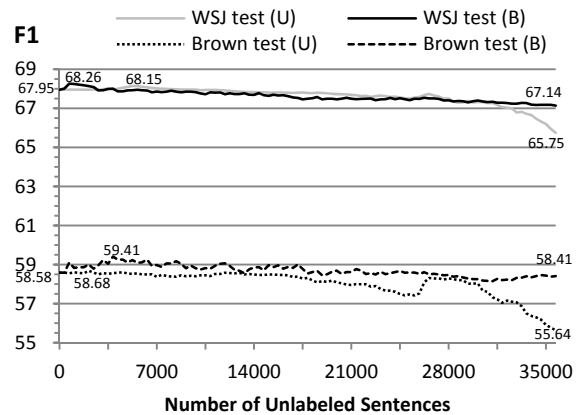


Figure 3: Balanced (B) and Unbalanced (U) Selection with OANC Unlabeled Data

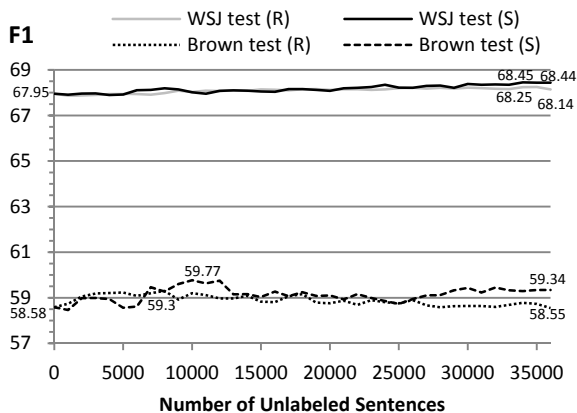


Figure 4: Random (R) and Simplicity (S) Pre-selection with WSJ Unlabeled Data

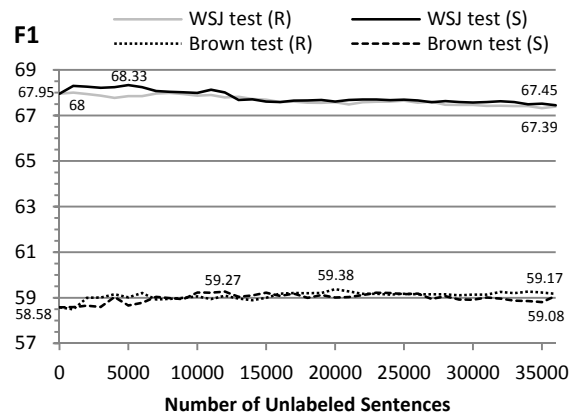


Figure 5: Random (R) and Simplicity (S) Pre-selection with OANC Unlabeled Data

more promising with both test sets. When trained on OANC, the F1 degrades with both methods as self-training progress. However, for both test sets, the best classifier is achieved by the balanced selection (68.26 vs. 68.15 and 59.41 vs. 58.68). Moreover, balanced selection shows a more normal behavior, while the other degrades the performance sharply in the last iterations (due to a swift drop of recall).

Consistent with previous work, with unbalanced selection, non-NULL-labeled unlabeled samples are selected only after the middle of the process. But, with the balanced method, selection is more evenly distributed over the roles.

A comparison between the results on Brown test set with each of unlabeled sets shows that in-domain data generalizes even better than out-of-domain data (59.62 vs. 59.41 and also note the trend). One apparent reason is that the classifier cannot accurately label the out-of-domain unlabeled data successively used for training. The lower quality of our out-of-domain data can be another reason for this behavior. Furthermore,

the parser we used was trained on WSJ, so it negatively affected the OANC parses and consequently its SRL results.

5.3 The Effect of Preselection

Figures 4 and 5 show the results of using pool with random and simplicity-based preselection with WSJ and OANC data respectively. The pool size (p) is 2000, and growth size (n) is 1000 sentences. The probability threshold (t) used is 0.5.

Comparing these figures with the previous figures shows that preselection improves the self-training trend, so that more unlabeled data can still be useful. This observation was consistent with various random selection of training data.

Between the two strategies, simplicity-based method outperforms the random method in both self-training trend and best classifier F1 (68.45 vs. 68.25 and 59.77 vs. 59.3 with WSJ and 68.33 vs. 68 with OANC), though the t-test shows that the F1 difference is not significant at $p \leq 0.05$. This improvement does not apply to the case of using OANC data when tested with Brown data

(59.27 vs. 59.38), where, however, the difference is not statistically significant. The same conclusion to the section 5.2 can be made here.

6 Conclusion and Future Work

This work studies the application of self-training in learning semantic role labeling with the use of unlabeled data. We used a balancing method for selecting newly labeled examples for augmenting the training set in each iteration of the self-training process. The idea was to reduce the effect of unbalanced distribution of semantic roles in training data. We also used a pool and examined two preselection methods for loading unlabeled data into it.

These methods showed improvement in both classifier performance and self-training trend. However, using out-of-domain unlabeled data for increasing the domain generalization ability of the system was not more useful than using in-domain data. Among possible reasons are the low quality of the used data and the poor parses of the out-of-domain data.

Another major factor that may affect the self-training behavior here is the poor performance of the base classifier compared to the state-of-the-art (see Table 2), which exploits more complicated SRL architecture. Due to high computational cost of self-training approach, bootstrapping experiments with such complex SRL approaches are difficult and time-consuming.

Moreover, parameter tuning process shows that other parameters such as pool-size, growth number and probability threshold are very effective. Therefore, more comprehensive parameter tuning experiments than what was done here is required and may yield better results.

We are currently planning to port this setting to co-training, another bootstrapping algorithm. One direction for future work can be adapting the architecture of the SRL system to better match with the bootstrapping process. Another direction can be adapting bootstrapping parameters to fit the semantic role labeling complexity.

References

Abney, S. 2008. *Semisupervised Learning for Computational Linguistics*. Chapman and Hall, London.

Baker, F., Fillmore, C. and Lowe, J. 1998. The Berkeley FrameNet project. In *Proceedings of COLING-ACL*, pages 86-90.

Charniak, E. and Johnson, M. 2005. Coarse-to-fine n-best parsing and MaxEnt discriminative reranking.

In *Proceedings of the 43rd Annual Meeting of the ACL*, pages 173-180.

- Carreras, X. and Marquez, L. 2005. Introduction to the CoNLL-2005 shared task: Semantic role labeling. In *Proceedings of the 9th Conference on Natural Language Learning (CoNLL)*, pages. 152-164.
- Clark S., Curran, R. J. and Osborne M. 2003. Bootstrapping POS taggers using Unlabeled Data. In *Proceedings of the 7th Conference on Natural Language Learning At HLT-NAACL 2003*, pages 49-55.
- Gildea, D. and Jurafsky, D. 2002. Automatic labeling of semantic roles. *CL*, 28(3):245-288.
- He, S. and Gildea, H. 2006. Self-training and Co-training for Semantic Role Labeling: Primary Report. TR 891, University of Colorado at Boulder
- Kingsbury, P. and Palmer, M. 2002. From Treebank to PropBank. In *Proceedings of the 3rd International Conference on Language Resources and Evaluation (LREC-2002)*.
- Lee, J., Song, Y. and Rim, H. 2007. Investigation of Weakly Supervised Learning for Semantic Role Labeling. In *Proceedings of the Sixth international Conference on Advanced Language Processing and Web information Technology (ALPIT 2007)*, pages 165-170.
- McClosky, D., Charniak, E., and Johnson, M. 2006. Effective self-training for parsing. In *Proceedings of the Main Conference on Human Language Technology Conference of the North American Chapter of the ACL*, pages 152-159.
- Ng, V. and Cardie, C. 2003. Weakly supervised natural language learning without redundant views. In *Proceedings of the 2003 Conference of the North American Chapter of the ACL on Human Language Technology*, pages 94-101.
- Punyakankok, V., Roth, D. and Yi, W. 2008. The Importance of Syntactic Parsing and Inference in Semantic Role Labeling. *CL*, 34(2):257-287.
- Surdeanu, M., Harabagiu, S., Williams, J. and Aarseth, P. 2003. Using predicate argument structures for information extraction. In *Proceedings of the 41st Annual Meeting of the ACL*, pages 8-15.
- Surdeanu, M., Johansson, R., Meyers, A., Marquez, L. and Nivre, J. 2008. The CoNLL 2008 shared task on joint parsing of syntactic and semantic dependencies. In *Proceedings of the 12th Conference on Natural Language Learning (CoNLL)*, pages 159-177.
- Yarowsky, E. 1995. Unsupervised Word Sense Disambiguation Rivaling Supervised Methods. In *proceeding of the 33rd Annual Meeting of ACL*, pages 189-196.

Weakly Supervised Learning of Presupposition Relations between Verbs

Galina Tremper

Department of Computational Linguistics
Heidelberg University, Germany
tremper@cl.uni-heidelberg.de

Abstract

Presupposition relations between verbs are not very well covered in existing lexical semantic resources. We propose a weakly supervised algorithm for learning presupposition relations between verbs that distinguishes five semantic relations: presupposition, entailment, temporal inclusion, antonymy and other/no relation. We start with a number of seed verb pairs selected manually for each semantic relation and classify unseen verb pairs. Our algorithm achieves an overall accuracy of 36% for type-based classification.

1 Introduction

A main characteristic of natural language is that significant portions of content conveyed in a message may not be overtly realized. This is the case for presuppositions: e.g. the utterance *Columbus didn't manage to reach India* presupposes that *Columbus had tried to reach India*. This presupposition does not need to be stated, but is implicitly understood. Determining the presuppositions of events reported in texts can be exploited to improve the quality of many natural language processing applications, such as information extraction, text understanding, text summarization, question-answering or machine translation.

The phenomenon of presupposition has been thoroughly investigated by philosophers and linguists (i.a. Stalnaker, 1974; van der Sandt, 1992). There are only few attempts for practical implementations of presupposition in computational linguistics (e.g. Bos, 2003). Especially, presupposition is understudied in the field of corpus-based learning of semantic relations. Machine learning methods have been previously applied to determine semantic relations such as *is-a* and *part-of*, also *succession*, *reaction* and *production* (Pantel

and Pennacchiotti, 2006). Chklovski and Pantel (2004) explored classification of fine-grained verb semantic relations, such as *similarity*, *strength*, *antonymy*, *enablement* and *happens-before*. For the task of entailment recognition, learning of entailment relations was attempted (Pekar, 2008). None of the previous work investigated subclassifying semantic relations including presupposition and entailment, two relations that are closely related, but behave differently in context.

In particular, the inferential behaviour of presuppositions and entailments crucially differs in special semantic contexts. E.g., while presuppositions are preserved under negation (as in *Columbus managed/didn't manage to reach India* the presupposition *tried to*), entailments do not survive under negation (*John F. Kennedy has been/has not been killed*). Here the entailment *died* only survives in the positive sentence. Such differences are crucial for both analysis and generation-oriented NLP tasks.

This paper presents a weakly supervised algorithm for learning presupposition relations between verbs cast as a discriminative classification problem. The structure of the paper is as follows: Section 2 reviews state of the art. Section 3 introduces our task and the learning algorithm. Section 4 reports on experiment organization; the results are presented in Section 5. Finally, we summarise and present objectives for future work.

2 Related Work

One of the existing semantic resources related to our paper is WordNet (Fellbaum, 1998). It comprises lexical semantic information about English nouns, verbs, adjectives and adverbs. Among the semantic relations defined specifically for verbs are entailment, hyponymy, troponymy, antonymy and cause. However, not all of them are well covered, for example, there are only few entries for presupposition and entailment in WordNet.

One attempt to acquire fine-grained semantic relations from corpora is VerbOcean (Chklovski and Pantel, 2004). Chklovski and Pantel used a semi-automatic approach for extracting semantic relations between verbs using a list of patterns. The selection of the semantic relations was inspired by WordNet. VerbOcean showed good accuracy values for the *antonymy* (50%), *similarity* (63%) and *strength* (75%) relations. However, VerbOcean doesn't distinguish between *entailment* and *presupposition*; they are conflated in the classes *enablement* and *happens-before*.

A distributional method for extracting highly associated verbs was proposed by Lin and Pantel (2001). This method extracts semantically related words with good precision, but it does not determine the type and symmetry of the relation. However, the method is able to recognize the existence of semantic relations holding between verbs and hence can be used as a basis for finding and further discriminating more detailed semantic relations.

3 A Weakly Supervised Approach to Learning Presupposition Relations

We describe a weakly supervised approach for learning semantic relations between verbs including implicit relations such as presupposition. Our aim is to perform a type-based classification of verb pairs. I.e., we determine the class of a verb-pair relation by observing co-occurrences of these verbs in contexts that are indicative for their intrinsic meaning relation. This task differs from a token-based classification, which aims at classifying each verb pair instance as it occurs in context.

Classified relations. We distinguish between the five classes of semantic relations presented in Table 1. We chose *entailment*, *temporal inclusion* and *antonymy*, because these relations may be confounded with the *presupposition* relation. A special class *other/no* comprises semantic relations not discussed in this paper (e.g. synonymy) and verb pairs that are not related by a semantic relation. The relations can be subdivided into symmetric and asymmetric relations, and relations that involve temporal sequence, or those that do not involve a temporal order, as displayed in Table 1.

A Weakly Supervised Learning Approach. Our algorithm starts with a small number of seed verb pairs selected manually for each relation and iteratively classifies a large set of unseen and un-

Semantic Relation	Example	Symmetry	Temporal Sequence
Presupposition	<i>find - seek,</i> <i>answer - ask</i>	asymmetric	yes
Entailment	<i>look - see,</i> <i>buy - own</i>	asymmetric	yes
Temporal Inclusion	<i>walk - step,</i> <i>talk - whisper</i>	symmetric	no
Antonymy	<i>win - lose,</i> <i>love - hate</i>	symmetric	no
Other/no	<i>have - own,</i> <i>sing - jump</i>	undefined	undefined

Table 1: Selected Semantic Relations

labeled verb pairs. Each iteration has two phases:

- 1. Training the Classifiers** We independently train binary classifiers for each semantic relation using both shallow and deep features.
- 2. Ensemble Learning and Ranking** Each of the five classifiers is applied to each sentence from an unlabeled corpus. The predictions of the classifiers are combined using ensemble learning techniques to determine the most confident classification. The obtained list of the classified instances is ranked using pattern scores, in order to select the most reliable candidates for extension of the training set.

Features. Both shallow lexical-syntactic and deep syntactic features are used for the classification of semantic relations. They include:

1. the distance between two analyzed verbs and the order of their appearance
2. verb form (tense, aspect, modality, voice), presence of negation and polarity verbs¹
3. coordinating/subordinating conjunctions
4. adverbial adjuncts
5. PoS-tag-contexts (two words preceding and two words following each verb)
6. the length of the path of grammatical functions relating the two verbs
7. co-reference relation holding between the subjects and objects of the verbs (both verbs have the same subject/object, subject of one verb corresponds to the object of the second or there is no relation between them).

In order to extract these features the training corpus is parsed using a deep parser.

¹Polarity verbs are taken from the polarity lexicon of Nairn et al. (2006). It encodes whether the complement of proposition embedding verbs is true or false. We used the verbs themselves as a feature without their polarity-tags.

4 Experimental Setting

Initial Subset of Verb Pair Candidates. Unlike other semi-supervised approaches, we don't use patterns for acquiring new candidates for classification. Candidate verb pairs are obtained from a previously compiled list of highly associated verbs. We use the DIRT Collection (Lin and Pantel, 2001) from which we further extract pairs of highly associated verbs as candidates for classification. The advantage of this resource is that it consists of pairs of verbs which stand in a semantic relation (cf. Section 2). This considerably reduces the number of verb pairs that need to be processed as candidates in our classification task.

DIRT contains 5,604 verb types and 808,764 verb pair types. This still represents a huge number of verb pairs to be processed. We therefore filtered the extracted set by checking verb pair frequency in the first three parts of the ukWAC corpus (Baroni et al., 2009) (UKWAC_1...3) and by applying the PMI test with threshold 2.0. This reduces the number of verb pairs to 199,393.

For each semantic relation we select three verb pairs as seeds. The only exception is *temporal inclusion* for which we selected six verb pairs, due to the low frequency of such verb pairs within a single sentence. These verb pairs were used for building an initial training corpus of verb pairs in context. The remaining verb pairs are used to build the corpus of unlabeled verb pairs in context in the iterative classification process.

Preprocessing. Given these verb pairs, we extracted sentences for training and for unlabeled data set from the first three parts of the UKWAC corpus (Baroni et al., 2009). We compiled a set of CQP queries (Evert, 2005) to find sentences that contain both verbs of a verb pair and applied them on UKWAC_1...3 to build the training and unlabeled subcorpora. We filter out sentences with more than 60 words and sentences with a distance between verbs exceeding 20 words. To avoid growing complexity, only sentences with exactly one occurrence of each verb pair are retained. We also remove sentences that trigger wrong candidates, in which the auxiliaries *have* or *do* appear in a candidate verb pair.

The corpus is parsed using the XLE parser (Crouch et al., 2008). Its output contains both the structural and functional information we need to extract the shallow and deep features used in the

classification, and to generate patterns.

Training Corpus. From this preprocessed corpus, we created a training corpus that contains three different components:

1. *Manually annotated training set.* All sentences containing seed verb pairs extracted from UKWAC_1 are annotated manually with two values *true/false* in order to separate the negative training data.
2. *Automatically annotated training set.* We build an extended, heuristically annotated training set for the seed verb pairs, by extracting further instances from the remaining corpora (UKWAC_2 and UKWAC_3). Using the manual annotations of step 1., we manually compiled a small stoplist of patterns that are used to filter out wrong instances. The constructed stoplist serves as an elementary disambiguation step. For example, the verbs *look* and *see* can stand in an entailment relation if *look* is followed by the prepositions *at*, *on*, *in*, but not in case of prepositions *after* or *forward* (e.g. *looking forward to*).
3. *Synonymous verb pairs.* To further enrich the training set of data, synonyms of the verb pairs are manually selected from WordNet. The corresponding verb pairs were extracted from UKWAC_1...3. In order to avoid adding noise, we used only synonyms of unambiguous verbs. The problem of ambiguity of the target verbs wasn't considered at this step.

The overall size of the training set for the first classification step is 15,717 sentences from which 5,032 are manually labeled, 9,918 sentences are automatically labeled and 757 sentences contain synonymous verb pairs. The distribution is unbalanced: *temporal inclusion* e.g. covers only 2%, while *entailment* covers 39% of sentences. We balanced the training set by undersampling *entailment* and *other/no* by 20% and correspondingly oversampling the *temporal inclusion* class.

Patterns. Similar to other pattern-based approaches we use a set of seed verb pairs to induce indicative patterns for each semantic relation. We use the induced patterns to restrict the number of the verb pair candidates and to rank the labelled instances in the iterative classification step.

The patterns use information about the verb forms of analyzed verb pairs, modal verbs and the

polarity verbs (only if they are related to the analyzed verbs) and coordinating/subordinating conjunctions connecting two verbs. The analyzed verbs in the sentence are substituted with V1 and V2 placeholders in the pattern. For example, for the sentence: *Here we should be careful for there are those who seek and do not find.* and the verb pair (*find,seek*) we induce the following pattern: *V2 and do [not|n't] V1*. The patterns are extracted automatically from deep parses of the training corpus. Examples of the best patterns we determined for semantic relations are presented in Table 2.

Semantic Relation	Patterns
Presupposition	V2-ed * though * was * V1-ed, V2-ed * but was [not n't] V1-ed, V2-ing * might V1
Entailment	if * V1 * V2, V1-ing * [shall will 'll] V2, V2 * by V1-ing
Temporal Inclusion	V2 * V1-ing, V1-ing and V2-ing, when V2 * V1
Antonymy	V1 or * V2, either * V1 or * V2, V1-ed * but V2-ed
Other/no	V1 * V2, V1-ing * V2-ing, V2-ed * and * V1-ed

Table 2: Patterns for Selected Semantic Relations

Pattern ranks are used to compute the reliability score for instances, as proposed by Pantel and Pennacchiotti (2006). The pattern reliability is calculated as follows:

$$r_{\pi}(p) = \frac{1}{|I|} \sum_{i \in I} \frac{pmi(i,p)}{max_{pmi}} \times r_i(i) \quad (1)$$

where:

$pmi(i,p)$ - pointwise mutual information (PMI) between the instance i and the pattern p ;

max_{pmi} - maximum PMI between all patterns and all instances;

$r_i(i)$ - reliability of an instance i . For seeds $r_i(i) = 1$ (they are selected manually), for the next iterations the instance reliability is:

$$r_i(i) = \frac{1}{|P|} \sum_{p \in P} \frac{pmi(i,p)}{max_{pmi}} \times r_{\pi}(p) \quad (2)$$

We also consider using the patterns as a feature for classification, in case they turn out to be sufficiently discriminative.

Training Binary Classifiers. We independently train 5 binary classifiers, one for each semantic relation, using the J48 decision tree algorithm (Witten and Frank, 2005).

Data Sets. As the primary goal of this paper is to classify semantic relations on the type level, we elaborated a first gold standard dataset for type-based classification. We used a small sample of 100 verb pairs randomly selected from the automatically labeled corpus. This sample was manually annotated by two judges after we had eliminated the system annotations in order not to influence the judges' decisions. The judges had the possibility to select more than one annotation, if necessary. We measured inter-annotator agreement was 61% ($k \approx 0.21$). The low agreement shows the difficulty of decision in the annotation of fine-grained semantic relations.²

While the first gold standard dataset of verb pairs was annotated *out of context*, we constructed a second gold standard of verb pairs annotated at the token level, i.e. in context. This second data set can be used to evaluate a token-based classifier (a task not attempted in the present paper). It also offers a ground truth for type-based classification, in that it controls for contextual ambiguity effects. I.e., we can extract a type-based gold standard on the basis of the token-annotated data.³ We proposed to one judge to annotate the same 100 verb pair types as in the previous annotation task, this time in context. For this purpose we randomly selected 10 instances for each verb pair type (for rare verb pair types only 5). We compared the gold standards elaborated by the same judge for type-based and token-based classification:

- 62% of verb pair types were annotated with the same labels on both levels, indicating correct annotation
- 10% of verb pair types were assigned conflicting labels, indicating wrong annotation
- 28% of verb pair types were assigned labels not present on the type level, or the type level label was not assigned in context

The figures show that for the most part the type-based annotation conforms with the ground truth obtained from token-based annotation. Only 10% of verb pair types were established as conflicting with the ground truth. The remaining 28% can be considered as potentially correct: either the annotated data does not contain the appropriate context for a given type label or the type-level anno-

²Data inspection revealed that one annotator was more experienced in semantic annotation tasks. We evaluate our system using the annotations of only one judge.

³This option was not pursued in the present paper.

tation, performed without context, does not foresee an existing relation. This points to a general difficulty, namely to acquire representative data sets for token-level annotation, and also to perform type-level annotations without context for the present task.

Combining Classifiers in Ensemble Learning.

Both token-based and type-based classification starts with determining of the most confident classification for instances. Each instance of the corpus of unlabeled verb pairs is classified by the individual binary classifiers. In order to select the most confident classification we compare the votes of the individual classifiers as follows:

1. If an instance is classified by one of the classifiers as *true* with confidence less than 0.75, we discard this classification.
2. If an instance is classified as *true* by more than one classifier, we consider only the classification with the highest confidence.⁴

In contrast to token-based classification that accepts only one semantic relation, for type-based classification we allow the existence of more than one semantic relation for a verb pair. To avoid the unreliable classifications, we apply several filters:

1. If less than 10% of the instances for a verb pair are classified with some specific semantic relation, this classification is considered to be unconfident and is discarded.
2. If a verb pair is classified as positive for more than three semantic relations, this verb pair remains unclassified.
3. If a verb pair is classified with up to three semantic relations and if more than 10% of the examples are classified with any of these relations, the verb pair is labeled with all of them.

Iteration and Stopping Criterion. After determining the most confident classification we rank the instances, following the ranking procedure of Pantel and Pennacchiotti (2006). Instances that exceed a reliability threshold (0.3 for our experiment) are selected for the extended training set. The remaining instances are returned to the unlabeled set. The algorithm stops if the average reliability score is smaller than a threshold value. In our paper we concentrate on the first iteration. Extension of the training set and re-ranking of patterns will be reported in future work.

⁴We assume that within a given context a verb pair can exhibit only one relation.

Semantic relation (Count1/Count2)	Majority	Without NONE	Baseline
Presupposition (12/22)	67%	36%	18%
Entailment (9/20)	67%	35%	8%
Temp. Inclusion (7/11)	71%	36%	19%
Antonymy (11/24)	72%	42%	12%
NONE (61/29)	49%	31%	43%
Macro-Average	56%	36%	
Micro-Average	65%	36%	

Table 3: Accuracy for type-based classification

5 Evaluation Results

Results for type-based classification. We evaluate the accuracy of classification based on two alternative measures:

1. *Majority* - the semantic relation with which the majority of the sentences containing a verb pair have been annotated.
2. *Without NONE* - as in 1., but after removing the label *NONE* from all relation assignments except for those cases where *NONE* is the only label assigned to a verb pair.⁵

We computed accuracy as the number of verb pairs which were correctly labeled by the system divided by the total number of system labels. We compare our results against a baseline of random assignment, taking the distribution found in the manually labeled gold standard as the underlying verb relation distribution. Table 3 shows the accuracy results for each semantic relation⁶.

Results for token-based classification. We also evaluate the accuracy of classification for token-based classification as the number of instances which were correctly labeled by the system divided by the total number of system labels. As the baseline we took the relation distribution on the token level. Table 4 shows the accuracy results for each semantic relation.

Discussion. The results obtained for type-based classification are well above the baseline with one exception. The best performance is achieved by *antonymy* (72% and 42% respectively for both

⁵The second measure was used because in many cases the relation *NONE* has been determined to be the majority class.

⁶Count1 is the total number of system labels for the Majority measure and Count2 is the total number of system labels for the Without *NONE* measure.

Semantic relation	Count	Accuracy	Baseline
Presupposition	43	21%	8%
Entailment	39	15%	5%
Temp. Inclusion	15	13%	3%
Antonymy	34	29%	5%
NONE	511	81%	79%
Macro-Average		61%	
Micro-Average		31%	

Table 4: Accuracy for token-based classification

measures), followed by *temporal inclusion*, *presupposition* and *entailment*. Accuracy scores for token-based classification (excluding NONE) are lower at 29% to 13%. Error analysis of randomly selected false positives shows that the main reason for lower accuracy on the token level is that the context is not always significant enough to determine the correct relation.

Comparison to Related Work. Other projects such as VerbOcean (Chklovski and Pantel, 2004) report higher accuracy: the average accuracy is 65.5% if at least one tag is correct and 53% for the correct preferred tag. However, we cannot objectively compare the results of VerbOcean to our system because of the difference in the set of relation classes and evaluation procedures. Similar to us, Chklovski and Pantel (2004) evaluated VerbOcean using a small sample of data which was presented to two judges for manual evaluation. In contrast to our setup, they didn't remove the system annotations from the evaluation data set. Given the difficulty of the classification we suspect that correction of system output relations for establishing a gold standard bears a strong risk in favouring system classifications.

6 Conclusion and Future Work

The results achieved in our experiment show that weakly supervised methods can be applied for learning presupposition relations between verbs. Our work also shows that they are more difficult to classify than other typical lexical semantic relations, such as antonymy. Error analysis suggests that many errors can be avoided if verbs are disambiguated in context. It would be interesting to test our algorithm with different amounts of manually annotated training sets and different combinations of manually and automatically annotated training sets to determine the minimal amount of

data needed to assure good accuracy.

In future work we will integrate word sense disambiguation as well as information about predicate-argument structure. Also, we are going to analyze the influence of single features on the classification and determining optimal feature sets, as well as the question of including patterns in the feature set. In this paper we used the same combination of features for all classifiers.

7 Acknowledgements

I would like to thank Anette Frank for supervision of this work, Dekang Lin and Patrick Pantel for sharing the DIRT resource and Carina Silberer and Christine Neupert creation of the gold standard.

References

- Baroni, M., Bernardini, S., Ferraresi, A., Zanchetta, E.: The WaCky Wide Web: a collection of very large linguistically processed web-crawled corpora. *Journal of Language Resources and Evaluation*, Vol.43 (3), 209–226 (2009)
- Bos, J.: Implementing the Binding and Accommodation Theory for Anaphora Resolution and Presupposition Projection. *Computational Linguistics*, Vol.29 (2), 179–210 (2003)
- Chklovski, T., Pantel, P.: Verbocean: Mining the web for fine-grained semantic verb relations. *Proceedings of EMNLP 2004*, 33–40, Barcelona (2004)
- Crouch, D., Dalrymple, M., Kaplan, R., King, T., Maxwell, J., Newman, P.: XLE Documentation. Palo Alto Research Center (2008)
- Evert, S.: The CQP Query Language Tutorial (CWB Version 2.2.b90). IMS, Stuttgart (2005)
- Fellbaum, C.: *WordNet: An Electronic Lexical Database*. 1st edition, MIT Press (1998)
- Lin, D., Pantel, P.: Discovery of Inference Rules for Question Answering. *Natural Language Engineering*, Vol.7, 343–360 (2001)
- Nairn, R., Condoravdi, C., Karttunen, L.: Computing Relative Polarity for Textual Inference. *Proc. of ICoS-5*, Buxton, UK (2006)
- Pantel, P., Pennacchiotti, M.: Espresso: Leveraging Generic Patterns for Automatically Harvesting Semantic Relations. *COLING 2006*, 113–120 (2006)
- Pekar, V.: Discovery of event entailment knowledge from text corpora. *Computer Speech & Language*, Vol.22 (1), 1–16 (2008)
- Stalnaker, R.C.: *Pragmatic Presuppositions*. *Semantics and Philosophy*, New York: Univ. Press (1974)
- van der Sandt, R.: Presupposition Projection as Anaphora Resolution. *Journal of Semantics*, Vol.9, 333–377 (1992)
- Witten, I., Frank, E.: *Data Mining: Practical Machine Learning Tools and Techniques*. (2005)

Importance of linguistic constraints in statistical dependency parsing

Bharat Ram Ambati

Language Technologies Research Centre, IIT-Hyderabad,
Gachibowli, Hyderabad, India – 500032.

ambati@research.iit.ac.in

Abstract

Statistical systems with high accuracy are very useful in real-world applications. If these systems can capture basic linguistic information, then the usefulness of these statistical systems improve a lot. This paper is an attempt at incorporating linguistic constraints in statistical dependency parsing. We consider a simple linguistic constraint that a verb should not have multiple subjects/objects as its children in the dependency tree. We first describe the importance of this constraint considering Machine Translation systems which use dependency parser output, as an example application. We then show how the current state-of-the-art dependency parsers violate this constraint. We present two new methods to handle this constraint. We evaluate our methods on the state-of-the-art dependency parsers for Hindi and Czech.

1 Introduction

Parsing is one of the major tasks which helps in understanding the natural language. It is useful in several natural language applications. Machine translation, anaphora resolution, word sense disambiguation, question answering, summarization are few of them. This led to the development of grammar-driven, data-driven and hybrid parsers. Due to the availability of annotated corpora in recent years, data driven parsing has achieved considerable success. The availability of phrase structure treebank for English (Marcus et al., 1993) has seen the development of many efficient parsers. Using the dependency analysis, a similar large scale annotation effort for Czech, has been the Prague Dependency Treebank (Hajicova, 1998). Unlike English, Czech is a free-word-order language and is also morphologically very rich. It has been suggested that free-word-order languages can be handled better using the dependency based framework than the constituency based one (Hudson, 1984; Shieber, 1985; Mel'čuk, 1988, Bharati et al., 1995). The basic

difference between a constituent based representation and a dependency representation is the lack of nonterminal nodes in the latter. It has also been noted that use of appropriate edge labels gives a level of semantics. It is perhaps due to these reasons that the recent past has seen a surge in the development of dependency based treebanks.

Due to the availability of dependency treebanks, there are several recent attempts at building dependency parsers. Two CoNLL shared tasks (Buchholz and Marsi, 2006; Nivre et al., 2007a) were held aiming at building state-of-the-art dependency parsers for different languages. Recently in NLP Tools Contest in ICON-2009 (Husain, 2009 and references therein), rule-based, constraint based, statistical and hybrid approaches were explored towards building dependency parsers for three Indian languages namely, Telugu, Hindi and Bangla. In all these efforts, state-of-the-art accuracies are obtained by two data-driven parsers, namely, Malt (Nivre et al., 2007b) and MST (McDonald et al., 2006). The major limitation of both these parsers is that they won't take linguistic constraints into account explicitly. But, in real-world applications of the parsers, some basic linguistic constraints are very useful. If we can make these parsers handle linguistic constraints also, then they become very useful in real-world applications.

This paper is an effort towards incorporating linguistic constraints in statistical dependency parser. We consider a simple constraint that a verb should not have multiple subjects/objects as its children. In section 2, we take machine translation using dependency parser as an example and explain the need of this linguistic constraint. In section 3, we propose two approaches to handle this case. We evaluate our approaches on the state-of-the-art dependency parsers for Hindi and Czech and analyze the results in section 4. General discussion and future directions of the work are presented in section 5. We conclude our paper in section 6.

2 Motivation

In this section we take Machine Translation (MT) systems that use dependency parser output as an example and explain the need of linguistic constraints. We take a simple constraint that a verb should not have multiple subjects/objects as its children in the dependency tree. Indian Language to Indian Language Machine Translation System¹ is one such MT system which uses dependency parser output. In this system the general framework has three major components. a) dependency analysis of the source sentence. b) transfer from source dependency tree to target dependency tree, and c) sentence generation from the target dependency tree. In the transfer part several rules are framed based on the source language dependency tree. For instance, for Telugu to Hindi MT system, based on the dependency labels of the Telugu sentence post-positions markers that need to be added to the words are decided. Consider the following example,

(1)
Telugu: raamu oka pamdu tinnaadu
'Ramu' 'one' 'fruit' 'ate'

Hindi: raamu ne eka phala khaayaa
'Ramu' 'ERG' 'one' 'fruit' 'ate'

English: "Ramu ate a fruit".

In the above Telugu sentence, 'raamu' is the subject of the verb 'tinnaadu'. While translating this sentence to Hindi, the post-position marker 'ne' is added to the subject. If the dependency parser marks two subjects, both the words will have 'ne' marker. This affects the comprehensibility. If we can avoid such instances, then the output of the MT system will be improved.

This problem is not due to morphological richness or free-word-order nature of the target language. Consider an example of free-word-order language to fixed-word-order language MT system like Hindi to English MT system. The dependency labels help in identifying the position of the word in the target sentence. Consider the example sentences given below.

(2a) raama seba khaatha hai
'Ram' 'apple' 'eats' 'is'
'Ram eats an apple'

(2b) seba raama khaatha hai
'apple' 'Ram' 'eats' 'is'
'Ram eats an apple'

Though the source sentence is different, the target sentence is same. Even though the source sentences are different, the dependency tree is same for both the sentences. In both the cases, 'raama' is the subject and 'seba' is the object of the verb 'khaatha'. This information helps in getting the correct translation. If the parser for the source sentence assigns the label 'subject' to both 'raama' and 'seba', the MT system can not give the correct output.

There were some attempts at handling these kind of linguistic constraints using integer programming approaches (Riedel et al., 2006; Bharati et al., 2008). In these approaches dependency parsing is formulated as solving an integer program as McDonald et al. (2006) has formulated dependency parsing as MST problem. All the linguistic constraints are encoded as constraints while solving the integer program. In other words, all the parses that violate these constraints are removed from the solution list. The parse which satisfies all the constraints is considered as the dependency tree for the sentence. In the following section, we describe two new approaches to avoid multiple subjects/objects for a verb.

3 Approaches

In this section, we describe the two different approaches for avoiding the cases of a verb having multiple subjects/objects as its children in the dependency tree.

3.1 Naive Approach (NA)

In this approach we first run a parser on the input sentence. Instead of first best dependency label, we extract the k-best labels for each token in the sentence. For each verb in the sentence, we check if there are multiple children with the dependency label 'subject'. If there are any such cases, we extract the list of all the children with label 'subject'. We find the node in this list which appears left most in the sentence with respect to other nodes. We assign 'subject' to this node. For the rest of the nodes in this list we assign the second best label and remove the first best label from their respective k-best list of labels. We check recursively, till all such instances are

¹ <http://sampark.iiit.ac.in/>

avoided. We repeat the same procedure for ‘object’.

Main criterion to avoid multiple subjects/objects in this approach is position of the node in the sentence. Consider the following example,

Eg. 3: *raama seba khaatha hai*
‘Ram’ ‘apple’ ‘eats’ ‘is’
‘Ram eats an apple’

Suppose the parser assigns the label ‘subject’ to both the nouns, ‘*raama*’ and ‘*seba*’. Then naive approach assigns the label subject to ‘*raama*’ and second best label to ‘*seba*’ as ‘*raama*’ precedes ‘*seba*’.

In this manner we can avoid a verb having multiple children with dependency labels subject/object.

Limitation to this approach is word-order. The algorithm described here works well for fixed word order languages. For example, consider a language with fixed word order like English. English is a SVO (Subject, Verb, Object) language. Subject always occurs before the object. So, if a verb has multiple subjects, based on position we can say that the node that occurs first will be the subject. But if we consider a free-word order language like Hindi, this approach wouldn’t work always.

Consider (2a) and (2b). In both these examples, ‘*raama*’ is the subject of the verb ‘*khaatha*’ and ‘*seba*’ is the object of the verb ‘*khaatha*’. The only difference in these two sentences is the order of the word. In (2a), subject precedes object. Whereas in (2b), object precedes subject. Suppose the parser identifies both ‘*raama*’ and ‘*seba*’ as subjects. NA can correctly identify ‘*raama*’ as the subject in case of (2a). But in case of (2b), ‘*seba*’ is identified as the subject. To handle these kind of instances, we use a probabilistic approach.

3.2 Probabilistic Approach (PA)

The probabilistic approach is similar to naive approach except that the main criterion to avoid multiple subjects/objects in this approach is probability of the node having a particular label. Whereas in naive approach, position of the node is the main criterion to avoid multiple subjects/objects. In this approach, for each node in the sentence, we extract the k-best labels along with their probabilities. Similar to NA, we first check for each verb if there are multiple children with the dependency label ‘subject’. If there are any such cases, we extract the list of all the

children with label ‘subject’. We find the node in this list which has the highest probability value. We assign ‘subject’ to this node. For the rest of the nodes in this list we assign the second best label and remove the first best label from their respective k-best list of labels. We check recursively, till all such instances are avoided. We repeat the same procedure for ‘object’.

Consider (2a) and (2b). Suppose the parser identifies both ‘*raama*’ and ‘*seba*’ as subjects. Probability of ‘*raama*’ being a subject will be more than ‘*seba*’ being a subject. So, the probabilistic approach correctly marks ‘*raama*’ as subject in both (2a) and (2b). But, NA couldn’t identify ‘*raama*’ as subject in (2b).

4 Experiments

We evaluate our approaches on the state-of-the-art parsers for two languages namely, Hindi and Czech. First we calculate the instances of multiple subjects/objects in the output of the state-of-the-art parsers for these two languages. Then we apply our approaches and analyze the results.

4.1 Hindi

Recently in NLP Tools Contest in ICON-2009 (Husain, 2009 and references herein), rule-based, constraint based, statistical and hybrid approaches were explored for parsing Hindi. All these attempts were at finding the inter-chunk dependency relations, given gold-standard POS and chunk tags. The state-of-the-art accuracy of 74.48% LAS (Labeled Attachment Score) is achieved by Ambati et al. (2009) for Hindi. They used two well-known data-driven parsers, Malt² (Nivre et al., 2007b), and MST³ (McDonald et al., 2006) for their experiments. As the accuracy of the labeler of MST parser is very low, they used maximum entropy classification algorithm, MAXENT⁴ for labeling.

For Hindi, dependency annotation is done using paninian framework (Begum et al., 2008; Bharati et al., 1995). So, in Hindi, the equivalent labels for subject and object are ‘*karta (k1)*’ and ‘*karma (k2)*’. ‘*karta*’ and ‘*karma*’ are syntactico-semantic labels which have some properties of both grammatical roles and thematic roles. *k1* behaves similar to subject and agent. *k2* behaves similar to object and patient (Bharati et al., 1995; Bharati et al., 2009). Here, by object we mean

² Malt Version 1.3.1

³ MST Version 0.4b

⁴ http://homepages.inf.ed.ac.uk/lzhang10/maxent_toolkit.htm

only direct object. Thus we consider only *k1* and *k2* labels which are equivalent of subject and direct object. Annotation scheme is such that there wouldn't be multiple subjects/objects for a verb in any case (Bharati et al., 2009). For example, even in case of coordination, coordinating conjunction is the head and conjuncts are children of the coordinating conjunction. The coordinating conjunction is attached to the verb with *k1/k2* label and the conjuncts get attached to the coordinating conjunction with a dependency label '*ccof*'.

We replicated the experiments of Ambati et al. (2009) on test set (150 sentences) of Hindi and analyzed the outputs of Malt and MST+MaxEnt. We consider this as the baseline. In the output of Malt, there are 39 instances of multiple subjects/objects. There are 51 such instances in the output of MST+MAXENT.

Malt is good at short distance labeling and MST is good at long distance labeling (McDonald and Nivre, 2007). As '*k1*' and '*k2*' are short distance labels, Malt could able predict these labels more accurately than MST. Because of this output of MST has higher number of instances of multiple subjects/objects than Malt.

	<i>Total Instances</i>
<i>Malt</i>	39
<i>MST + MAXENT</i>	51

Table 1: Number of instances of multiple subjects or objects in the output of the state-of-the-art parsers for Hindi

Both the parsers output first best label for each node in the sentence. In case of Malt, we modified the implementation to extract all the possible dependency labels with their scores. As Malt uses libsvm for learning, we couldn't able to get the probabilities. Though interpreting the scores provided by libsvm as probabilities is not the correct way, that is the only option currently available with Malt. In case of MST+MAXENT, labeling is performed by MAXENT. We used a java version of MAXENT⁵ to extract all possible tags with their scores. We applied both the naive and probabilistic approaches to avoid multiple subjects/objects. We evaluated our experiments based on unlabeled attachment score (UAS), labeled attachment score (LAS) and labeled score

(LS) (Nivre et al., 2007a). Results are presented in Table 2.

As expected, PA performs better than NA. With PA we got an improvement of 0.26% in LAS over the previous best results for Malt. In case of MST+MAXENT we got an improvement of 0.61% in LAS over the previous best results. Note that in case of MST+MAXENT, the slight difference between state-of-the-art results of Ambati et al. (2009) and our baseline accuracy is due different MAXENT package used.

	<i>Malt</i>			<i>MST+MAXENT</i>		
	<i>UAS</i>	<i>LAS</i>	<i>LS</i>	<i>UAS</i>	<i>LAS</i>	<i>LS</i>
<i>Baseline</i>	90.14	74.48	76.38	91.26	72.75	75.26
<i>NA</i>	90.14	74.57	76.38	91.26	72.84	75.26
<i>PA</i>	90.14	74.74	76.56	91.26	73.36	75.87

Table 2: Comparison of NA and PA with previous best results for Hindi

Improvement in case of MST+MAXENT is greater than that of Malt. One reason is because of more number of instances of multiple subjects/objects in case of MST+MAXENT. Other reason is use of probabilities in case MST+MAXENT. Whereas in case of Malt, we interpreted the scores as probabilities which is not a good way to do. But, in case of Malt, that is the only option available.

4.2 Czech

In case of Czech, we replicated the experiments of Hall et al. (2007) using latest version of Malt (version 1.3.1) and analyzed the output. We consider this as the baseline. The minor variation of the baseline results from the results of CoNLL-2007 shared task is due to different version Malt parser being used. Due to practical reasons we couldn't use the older version. In the output of Malt, there are 39 instances of multiple subjects/objects out of 286 sentences in the testing data. In case of Czech, the equivalent labels for subject and object are 'agent' and 'theme'.

Czech is a free-word-order language similar to Hindi. So as expected, PA performed better than NA. Interestingly, accuracy of PA is lower than the baseline. Main reason for this is scores of libsvm of Malt. We explain the reason for this using the following example, consider a verb 'V' has two children 'C1' and 'C2' with dependency label subject. Assume that the label for 'C1' is subject and the label of 'C2' is object in the gold-data. As the parser marked 'C1' with subject, this

⁵ <http://maxent.sourceforge.net/>

adds to the accuracy of the parser. While avoiding multiple subjects, if ‘C1’ is marked as subject, then the accuracy doesn’t drop. If ‘C2’ is marked as object then the accuracy increases. But, if ‘C2’ is marked as subject and ‘C1’ is marked as object then the accuracy drops. This could happen if probability of ‘C1’ having subject as label is lower than ‘C1’ having subject as the label. This is because of two reasons, (a) parser itself wrongly predicted the probabilities, and (b) parser predicted correctly, but due to the limitation of libsvm, we couldn’t get the scores correctly.

	<i>UAS</i>	<i>LAS</i>	<i>LS</i>
<i>Baseline</i>	82.92	76.32	83.69
<i>NA</i>	82.92	75.92	83.35
<i>PA</i>	82.92	75.97	83.40

Table 3: Comparison of NA and PA with previous best results for Czech

5 Discussion and Future Work

Results show that the probabilistic approach performs consistently better than the naive approach. For Hindi, we could able to achieve an improvement 0.26% and 0.61% in *LAS* over the previous best results using Malt and MST respectively. We couldn’t able to achieve any improvement in case of Czech due to the limitation of libsvm learner used in Malt.

We plan to evaluate our approaches on all the data-sets of CoNLL-X and CoNLL-2007 shared tasks using Malt. Settings of MST parser are available only for CoNLL-X shared task data sets. So, we plan to evaluate our approaches on CoNLL-X shared task data using MST also. Malt has the limitation for extracting probabilities due to libsvm learner. Latest version of Malt (version 1.3.1) provides option for liblinear learner also. Liblinear provides option for extracting probabilities. So we can also use liblinear learning algorithm for Malt and explore the usefulness of our approaches. Currently, we are handling only two labels, subject and object. Apart from subject and object there can be other labels for which multiple instances for a single verb is not valid. We can extend our approaches to handle such labels also. We tried to incorporate one simple linguistic constraint in the statistical dependency parsers. We can also explore the ways of incorporating other useful linguistic constraints.

6 Conclusion

Statistical systems with high accuracy are very useful in practical applications. If these systems can capture basic linguistic information, then the usefulness of the statistical system improves a lot. In this paper, we presented a new method of incorporating linguistic constraints into the statistical dependency parsers. We took a simple constraint that a verb should not have multiple subjects/objects as its children. We proposed two approaches, one based on position and the other based on probabilities to handle this. We evaluated our approaches on state-of-the-art dependency parsers for Hindi and Czech.

Acknowledgments

I would like to express my gratitude to Prof. Joakim Nivre and Prof. Rajeev Sangal for their guidance and support. I would also like to thank Mr. Samar Husain for his valuable suggestions.

References

- B. R. Ambati, P. Gadde and K. Jindal. 2009. Experiments in Indian Language Dependency Parsing. In *Proceedings of the ICON09 NLP Tools Contest: Indian Language Dependency Parsing*, pp 32-37.
- R. Begum, S. Husain, A. Dhvaj, D. Sharma, L. Bai, and R. Sangal. 2008. Dependency annotation scheme for Indian languages. In *Proceedings of IJCNLP-2008*.
- A. Bharati, V. Chaitanya and R. Sangal. 1995. *Natural Language Processing: A Paninian Perspective*, Prentice-Hall of India, New Delhi, pp. 65-106.
- A. Bharati, S. Husain, D. M. Sharma, and R. Sangal. 2008. A Two-Stage Constraint Based Dependency Parser for Free Word Order Languages. In *Proceedings of the COLIPS International Conference on Asian Language Processing 2008 (IALP)*. Chiang Mai, Thailand.
- S. Buchholz and E. Marsi. 2006. CoNLL-X shared task on multilingual dependency parsing. In *Proc. of the Tenth Conf. on Computational Natural Language Learning (CoNLL)*.
- E. Hajicova. 1998. Prague Dependency Treebank: From Analytic to Tectogrammatical Annotation. In *Proc. TSD ’98*.
- J. Hall, J. Nilsson, J. Nivre, G. Eryigit, B. Megyesi, M. Nilsson and M. Saers. 2007. Single Malt or Blended? A Study in Multilingual Parser Optimization. In *Proceedings of the CoNLL Shared Task Session of EMNLP-CoNLL*.
- R. Hudson. 1984. *Word Grammar*, Basil Blackwell, 108 Cowley Rd, Oxford, OX4 1JF, England.

- S. Husain. 2009. Dependency Parsers for Indian Languages. In *Proceedings of ICON09 NLP Tools Contest: Indian Language Dependency Parsing*. Hyderabad, India.
- M. Marcus, B. Santorini, and M.A. Marcinkiewicz. 1993. Building a large annotated corpus of English: The Penn Treebank, *Computational Linguistics* 1993.
- I. A. Mel'čuk. 1988. *Dependency Syntax: Theory and Practice*, State University, Press of New York.
- R. McDonald, K. Lerman, and F. Pereira. 2006. Multilingual dependency analysis with a two-stage discriminative parser. In *Proceedings of the Tenth Conference on Computational Natural Language Learning (CoNLL-X)*, pp. 216–220.
- R. McDonald and J. Nivre. 2007. Characterizing the errors of data-driven dependency parsing models. In *Proc. of EMNLP-CoNLL*.
- J. Nivre, J. Hall, S. Kubler, R. McDonald, J. Nilsson, S. Riedel and D. Yuret. 2007a. The CoNLL 2007 Shared Task on Dependency Parsing. In *Proceedings of EMNLP/CoNLL-2007*.
- J. Nivre, J. Hall, J. Nilsson, A. Chanev, G. Eryigit, S. Kübler, S. Marinov and E Marsi. 2007b. *MaltParser: A language-independent system for data-driven dependency parsing*. *Natural Language Engineering*, 13(2), 95-135.
- S. Riedel, Ruket Çakıcı and Ivan Meza-Ruiz. 2006. Multi-lingual Dependency Parsing with Incremental Integer Linear Programming. In *Proceedings of the Tenth Conference on Computational Natural Language Learning (CoNLL-X)*.
- S. M. Shieber. 1985. Evidence against the context-freeness of natural language. In *Linguistics and Philosophy*, p. 8, 334–343.

The use of formal language models in the typology of the morphology of Amerindian languages

Andrés Osvaldo Porta
Universidad de Buenos Aires
hugporta@yahoo.com.ar

Abstract

The aim of this work is to present some preliminary results of an investigation in course on the typology of the morphology of the native South American languages from the point of view of the formal language theory. With this object, we give two contrasting examples of descriptions of two Aboriginal languages finite verb forms morphology: Argentinean Quechua (quichua santiagueño) and Toba. The description of the morphology of the finite verb forms of Argentinean quechua, uses finite automata and finite transducers. In this case the construction is straightforward using two level morphology and then, describes in a very natural way the Argentinean Quechua morphology using a regular language. On the contrary, the Toba verbs morphology, with a system that simultaneously uses prefixes and suffixes, has not a natural description as regular language. Toba has a complex system of causative suffixes, whose successive applications determinate the use of prefixes belonging different person marking prefix sets. We adopt the solution of Creider et al. (1995) to naturally deal with this and other similar morphological processes which involve interactions between prefixes and suffixes and then we describe the toba morphology using linear context-free languages.¹

1 Introduction

It has been proved (Johnson, 1972; Kaplan and Kay, 1994) that regular models have an expressive

¹This work is part of the undergraduate thesis Finite state morphology: The Koskenniemi's two level morphology model and its application to describing the morphosyntax of two native Argentinean languages

power equal to the noncyclic components of generative grammars representing the morphophonology of natural languages. However, these works make no considerations about what class of formal languages is the natural for describing the morphology of one particular language. On the other hand, the criteria of classification of Amerindian languages, do not involve complexity criteria. In order to establish criteria that take into account the complexity of the description we present two contrasting examples in two Argentinean native languages: toba and quichua santiagueño. While the quichua has a natural representation in terms of a regular language using two level morphology, we will show that the Toba morphology has a more natural representation in terms of linear context-free languages.

2 Quichua Santiagueño

The quichua santiagueño is a language of the Quechua language family. It is spoken in the Santiago del Estero state, Argentina. Typologically is an agglutinative language and its structure is almost exclusively based on the use of suffixes and is extremely regular. The morphology takes a dominant part in this language with a rich set of validation suffixes. The quichua santiagueño has a much simpler phonologic system than other languages of this family: for example it has no series of aspirated or glottalized stops.

Since the description of the verbal morphology is rich enough for our aim to expose the regular nature of quichua santiagueño morphology, we have restricted our study to the morphology of finite verbs forms. We use the two level morphology paradigm to express with finite regular transducers the rules that clearly illustrate how naturally this language phonology is regular. The construction uses the descriptive works of Alderetes (2001) and Nardi (Albarracín et. al, 2002)

2.1 Phonological two level rules for the quichua santiagueño

In this section we present the alphabet of the quichua santiagueño with which we have implemented the quichua phonological rules in the paradigm of two level morphology. The subsets abbreviations are: V (vowel), Vlng (underlying vowel), Valt (high vowel), VMed (median vowel), Vbaj (bass vowel), Ftr (transparent to medialization phonema), Cpos (posterior consonant).

ALPHABET

a e i o u p t ch k q s sh ll m
n l r w y b d g gg f h x r rr
A E I O U W N Q Y + '

NULL 0

ANY @

BOUNDARY #

SUBSET C p t ch k q s sh ll m
n l r w y b d f h
x r rr h Q

SUBSET V i e a o u A E I O U

SUBSET Vlng I E A O U

SUBSET Valt u i U I

SUBSET VMed e o E O

SUBSET Vbaj a A

SUBSET Ftr n y r Y N

SUBSET Cpos gg q Q

With the aim of showing the simplicity of the phonologic rules we transcribe the two-level rules we have implemented with the transducers in the thesis. R1-R4 model the medialization vowels processes, R5-R7 are elision and epenthesis processes with very specific contexts and R7 represents a diachornic phonological process with a subjacent form present in others quechua dialects.

Rules

R1 i:i /<= CPos:@ _

R2 i:i /<= _ Ftr:@ CPos:@

R3 u:u /<= CPos:@ _

R4 u:u /<= _ Ftr:@ CPos:@

R5 W:w <=> a:a a:a +:0 _a:a +:0

R6 U:0 <=> m:m _+:0 p:p u:u +:0

R7 N:0 <=> _+:0 r:@ Q:@ a:a +:0

2.2 Quichua Santigueño morphology

The grammar that models the agglutination order is showed with a non deterministic finite automata. This implemented automata is presented in Figure 1. This description of the morphophonology was implemented using PC-KIMMO (Antworth, 1990)

3 The Toba morphology

The Toba language belongs, with the languages pilaga, mocovi and kaduveo, to the guaycuru language family (Messineo, 2003; Klein, 1978). The toba is spoken in the Gran Chaco region (which is comprised between Argentina, Bolivia and Paraguay) and in some reduced settlements near Buenos Aires, Argentina. From the point of view of the morphologic typology it presents characteristics of a polysynthetic agglutinative language. In this language the verb is the morphologically more complex wordclass. The grammatical category of person is prefixed to the verbal theme. There are suffixes to indicate plurals and other grammatical categories as aspect, location-direction, reflexive and reciprocal and desiderative mode. The verb has no mark of time. As an example of a typical verb we can considerate the *sanadatema*:

Example 1 .

s- anat(a) -d -em -a
IAct- advice 2 dat ben ²
" I advice you"

One of the characteristics of the toba verb morphology is a system of markation active-inactive on the verbal prefixes (Messineo, 2003; Klein, 1978). There are in this language two sets or verbal prefixes that mark action:

1. Class I (In): codifies inactive participants, objects of transitive verbs and patients of intransitive verbs. .
2. Class II (Act): codifies active participants, subjects of transitive and intransitive verbs.

²abrev: Act:active, ben:benefactive, dat:dative, inst: intrumental, Med: Median voice, pos: Possessor, refl: reflexive

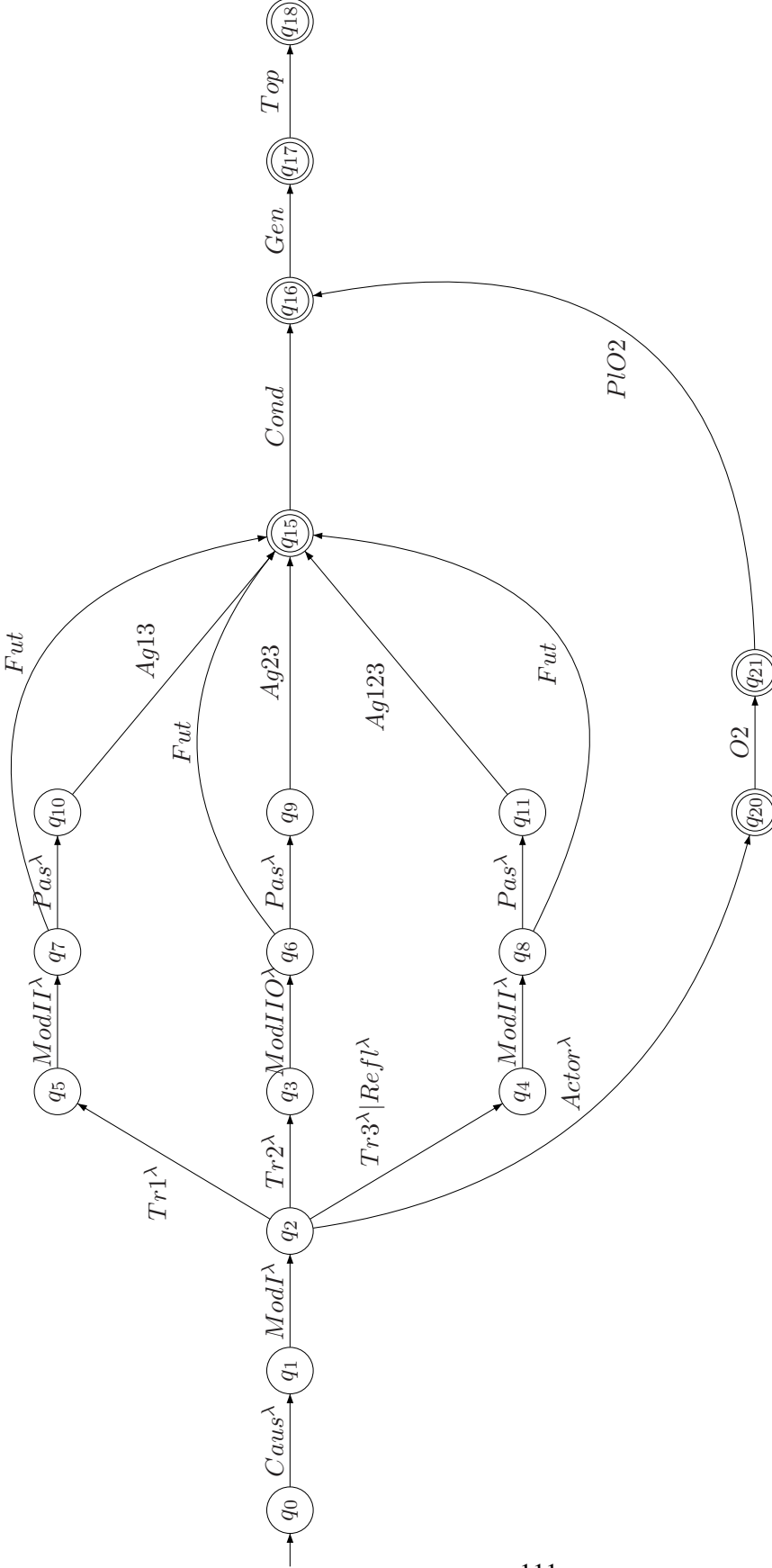


Figure 1: Schema of the verbal morphology of the quichua santiagueño. The supra indices λ indicate possible null transitions.
 Abrev.: Caus: Causative suffixes. ModI: Set I of Modal Suffixes. Tri : i th. person transition Suffixes. ModII: Set II of modal suffixes .Pas: Past suffixes. Ag: AgentSuffix , in this case, for example, Ag1, indicates the agent suffix for the 1st person. Ag12 is an abbreviator for A1 U A2. Cond : Conditional Suffixes. Gen : General Suffixes. Fut: future suffixes. Top : Topicaliser suffixes. PIO2: Plural of Object 2nd person Suffixes.

Active affected(Medium voice, Med): codifies the presence of an active participant affected by the action that the verb codifies. .

The toba has a great quantity of morphological processes that involve interactions between suffixes and prefixes. In the next example the suffixation of the reflexive (*-l'at*) forces the use of the active person with prefixes of the voice medium class because the agent is affected by the action.

Example 2 .

- (a) *y-* *alawat*
 3Activa *-kill*
 " *He* *kills*"
- (b) *n-* *alawat* *-l'at*
 3Med- *kill* *-refl*
 " *He kills* *himself*"

The agglutination of this suffix occurs in the last suffix box (after locatives, directional and other derivational suffixes). Then, if we model this process using finite automata we will add many items to the lexicon (Sproat, 1992). The derivation of names from verbs is very productive. There are many nominalizer suffixes. The resulting names use obligatory possessing person nominal prefixes.

Example 3 .

- l-* *edaGan* *-at*
- 3pos *write* *instr*
- "his pencil"

The toba language also presents a complex system of causative suffixes that act as switching the transitivity of the verb. Transitivity is specially appreciable in the switching of the 3rd person prefix mark. In section 3.2 we will use this process to show how linear context free grammars are a better than regular grammars for modeling agglutination in this language, but first we will present the former class of languages and its acceptor automata.

3.1 Linear context free languages and two-taped nondeterministic finite-state automata

A linear context-free language is a language generated by a grammar context-free grammars *G* in which every production has one of the three forms (Creider et al., 1995):

1. $A \rightarrow a$, with *a* terminal symbol
2. $A \rightarrow aB$, with *B* a non terminal symbol and *a* a terminal symbol.

3. $A \rightarrow Ba$, with *B* a non terminal symbol and *a* a terminal symbol.

Linear context-free grammars have been studied by Rosenberg (1967) who showed that there is an equivalence between them and two-taped nondeterministic finite-state automata. Informally, a two-head nondeterministic finite-state automata could be thought as a generalization of a usual nondeterministic finite-state automata which has two read heads that independently reads in two different tapes, and at each transition only one tape moves. When both tapes have been processed, if the automata is at a final state, the parsing is successful. In the ambit that we are studying we can think that if a word is a string of prefixes, a stem and suffixes, one automata head will read will the prefixes and the other the suffixes. Taking into account that linear grammars are in Rosenberg's terms: "The lowest class of nonregular context-free grammars", Creider et al. (1995) have taken this formal language class as the optimal to model morphological processes that involve interaction between prefixes and suffixes.

3.2 Analysis of the third person verbal paradigm

In this section we model the morphology of the third person of transitive verbs using two-taped finite nondeterministics automata. The modeling of this person is enough to show this description advantages with respect to others in terms of regular languages. The transitivity of the verb plays an important role in the selection of the person marker Class. The person markers are (Messineo, 2003):

1. *i-/y-* for transitive verbs *y* and some intransitive subjects (Pr ActT).
2. *d(Vowel)* for verbs typically intransitives (Pr ActI).
3. *n*: subjects of medium voice (Pr ActM).

The successive application of the causative seems to act here, as was interpreted by Buckwalter (2001), like making the switch in the original verb transitivity as is shown en Example 4 in the next page.

Example 4 .

IV	de-	que'e		he eats
TV	i-	qui'	-aGan	he eats(something)
IV	de-	qui'	-aGanataGan	he feeds
TV	i-	qui'	-aGanataGanaGan	he feeds(a person)
IV	de	qui'	-aGanaGanataGan	he command to feed

If we want to model this morphological process using finite automata again we must enlarge the lexicon size. The resulting grammar, although capable of modeling the morphology of the toba, would not work effectively. The effectiveness of a grammar is a measure of their productivity (Heintz, 1991). Taking into account the productivity of causative and reflexive verbal derivation we will prefer a description in terms of a context-free linear grammar with high effectivity than another using regular languages with low effectivity.

To model the behavior of causative agglutination and the interaction with person prefixes using the two-head automata, we define two paths determined by the parity of the causative suffixes which have been agglutinated to the verb. We have also to take into consideration the optative posterior agglutination of reflexive and reciprocal suffixes which forces the use of medium voice person prefix. From the third person is also formed the third person indefinite actor from a prefix, *qa* -, which is at left and adjacent to the usual mark of the third person and after the mark of negation *sa*-. Therefore, their agglutination is reserved to the last transitions. The resulting two-typed automata showed in Figure 2 also takes into account the relative order of the boxes and so the mutual restrictions between them (Klein, 1978).

4 Future Research

It is interesting to note that phonological rules in toba can be naturally expressed by regular Finite Transducers. There are, however, many South American native languages that presents morphological processes analogous to the Toba and some can present phonological processes that will have a more natural expression using Linear Finite Transducers. For example the Guarani language presents nasal harmony which expands from the root to both suffixes and prefixes (Krivoshein, 1994). This kind of characterization can have some value in language classification and the modeling of the great diversity of South American languages morphology can allow to obtain a formal concept of natural description of a language.

References

- Lelia Albarracín, Mario Tebes y Jorge Alderetes(eds.) 2002. *Introducción al quichua santiagueño por Ricardo L.J. Nardi*. Editorial DUNKEN: Buenos Aires, Argentina.
- Jorge Ricardo Alderetes 2002. *El quichua de Santiago del Estero. Gramática y vocabulario..* Tucumán: Facultad de Filosofía y Letras, UNT:Buenos Aires, Argentina.
- Evan L. Antworth 1990. *PC-KIMMO: a two-level processor for morphological analysis.No. 16 in Occasional publications in academic computing*. No. 16 in Occasional publications in academic computing. Dallas: Summer Institute of Linguistics.
- Alberto Buckwalter 2001. *Vocabulario toba*. Formosa / Indiana, Equipo Menonita.
- Chet Creider, Jorge Hankamer, and Derick Wood. 1995. Preset two-head automata and morphological analysis of natural language . *International Journal of Computer Mathematics*, Volume 58, Issue 1, pp. 1-18.
- Joos Heintz y Claus Schönig 1991. Turcic Morphology as Regular Language. *Central Asiatic Journal*, 1-2, pp 96-122.
- C. Douglas Johnson 1972. *Formal Aspects of Phonological Description*. The Hague:Mouton.
- Ronald M. Kaplan and Martin Kay. 1994. Regular models of phonological rule systems . *Computational Linguistics*,20(3):331-378.
- Harriet Manelis Klein 1978. *Una gramática de la lengua toba: morfología verbal y nominal*. Universidad de la República, Montevideo, Uruguay.
- Natalia Krivoshein de Canese 1994. *Gramática de la lengua guaraní*. Colección Nemity, Asunción, Paraguay.
- María Cristina Messineo 2003. *Lengua toba (guaycurú). Aspectos gramaticales y discursivos*. LINCOM Studies in Native American Linguistics 48. München: LINCOM EUROPA Academic Publisher.
- A.L Rosenberg 1967 A Machine Realization of the linear Context-Free Languages. *Information and Control*, 10: 177-188.
- Richard Sproat 1992. *Morphology and Computation*. The MIT Press.

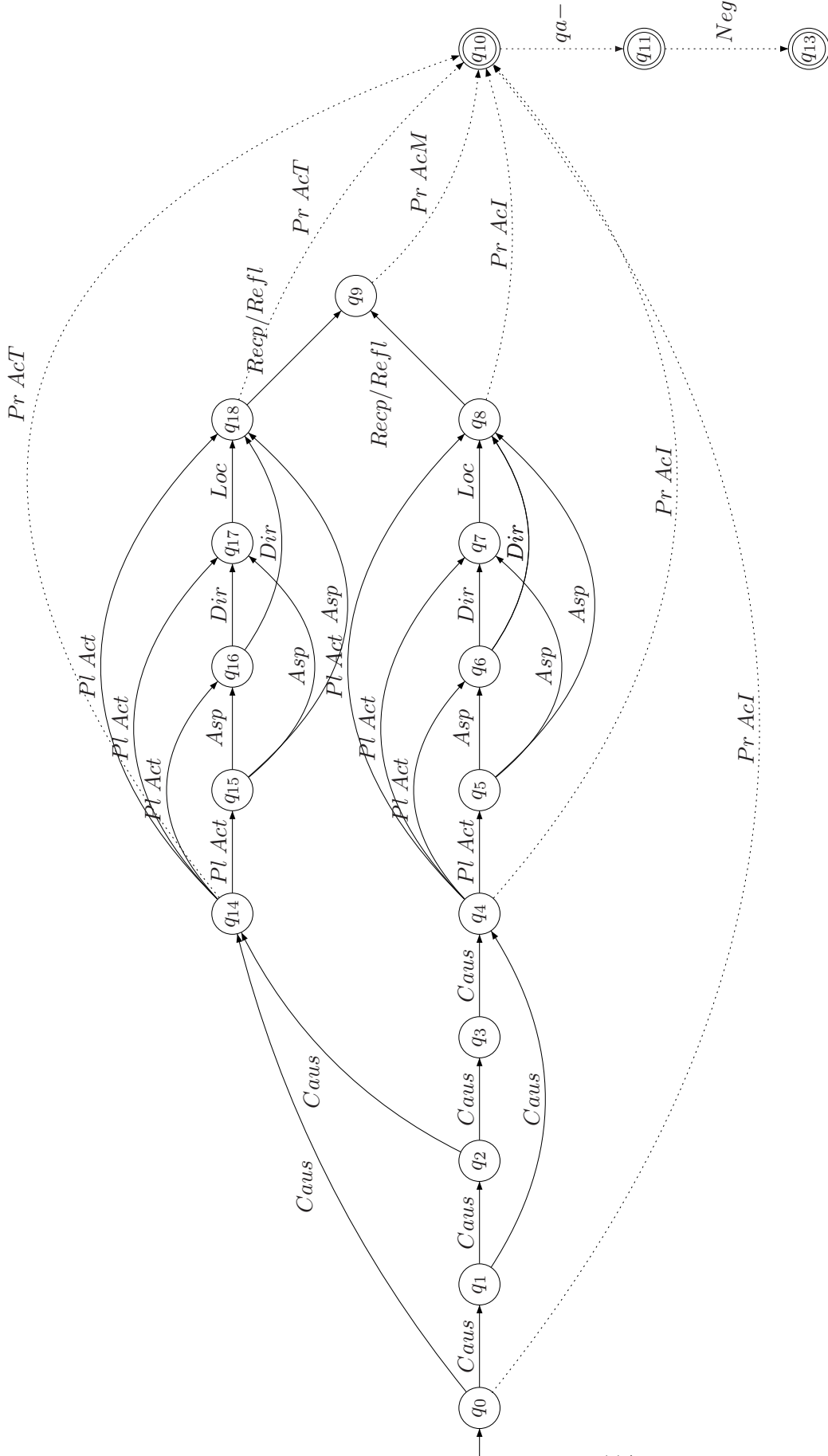


Figure 2: Schema of the 3rd person intransitive verb morphology of the toba .

The entire and dotted lines indicating transitions of the suffix and prefix tape, respectively

Abrev: Caus: Causative suffix. Pl Act: plural actors suffix. Asp: aspectual suffix. Dir: directive suffix. Loc: locative suffix. Recept: reciprocal action suffix. Refl: reflexive prefix. Pr:Pr: acting person prefix(T: transitive, I: intransitive, M: medium) $qa-$: indeterminate person prefix. Neg: negation prefix

Author Index

Akın, Ahmet Afşın, 31
Ambati, Bharat Ram, 103

Bogdanova, Daria, 67

Franco-Penya, Hector-Hugo, 79

Haulrich, Martin, 55

Inumella, Abhilash, 13

Lison, Pierre, 7
Liu, Zhiyuan, 49

McFate, Clifton, 61
McGillivray, Barbara, 73
Mermer, Coşkun, 31
Mittal, Vipul, 85

Nguyen, Thin, 43

Plüss, Brian, 1
Porta, Andres Osvaldo, 109

Reddy, Siva, 13

Samad Zadeh Kaljahi, Rasoul, 91
Sangati, Federico, 19, 37
Scheible, Christian, 25

Tremper, Galina, 97

Umanski, Daniil, 37

Xie, Lixing, 49

Zheng, Yabin, 49

ACL 2010

**48th Annual Meeting of the
Association for Computational Linguistics**

Proceedings of System Demonstrations

13 July 2010
Uppsala University
Uppsala, Sweden

Production and Manufacturing by
Taberg Media Group AB
Box 94, 562 02 Taberg
Sweden

©2010 The Association for Computational Linguistics

Order copies of this and other ACL proceedings from:

Association for Computational Linguistics (ACL)
209 N. Eighth Street
Stroudsburg, PA 18360
USA
Tel: +1-570-476-8006
Fax: +1-570-476-0860
acl@aclweb.org

Introduction

Welcome to the proceedings of the system demonstration session. This volume contains the papers of the system demonstrations presented at the 48th Annual Meeting of the Association for Computational Linguistics, held in Uppsala, Sweden, on July 13 2010.

The system demonstrations program offers the presentation of early research prototypes as well as interesting mature systems. The system demonstration chair and the members of the program committee received 65 submissions, 14 of which were selected for inclusion in the program after review by two members of the program committee.

I would like to thank the members of the program committee for their excellent job in reviewing the submissions and providing their support for the final decision.

Chair

Sandra Kübler, Indiana University

Program Committee:

Thorsten Brants (Google, USA)
Sabine Buchholz (Toshiba Research Europe, U.K.)
Paul Davis (Motorola, USA)
Thierry Declerck (DFKI, Germany)
Markus Dickinson (Indiana University, USA)
Michael Gasser (Indiana University, USA)
Günther Görz (Universität Erlangen, Germany)
Iryna Gurevych (Technische Universität Darmstadt, Germany)
Ahmed Hassan (University of Michigan, USA)
Julia Hockenmaier (University of Illinois, USA)
Wolfgang Hoepfner (Universität Duisburg, Germany)
Angela Kluttsch (Universität Duisburg, Germany)
Yuji Matsumoto (Nara Institute of Science and Technology, Japan)
Emad Mohamed (Indiana University, USA)
Preslav Nakov (National University of Singapore)
Nicolas Nicolov (JD Power and Associates, USA)
Petya Osenova (Bulgarian Academy of Sciences, Bulgaria)
Arzucan Ozgur (University of Michigan, USA)
Vahed Qazvinian (University of Michigan, USA)
Paul Rodrigues (University of Maryland, USA)
Nathan Sanders (Indiana University, USA)
Martin Scholz (Universität Erlangen, Germany)
Kiril Simov (Bulgarian Academy of Sciences, Bulgaria)
Antal van den Bosch (Tilburg University, The Netherlands)
Holger Wunsch (Universität Tübingen, Germany)
Annie Zaenen (Palo Alto Research Center, USA)
Desislava Zhekova (Universität Bremen, Germany)

Additional Reviewers:

Daniel Bär (Technische Universität Darmstadt, Germany)
Georgi Georgiev (Ontotext, Bulgaria)
Niklas Jakob (Technische Universität Darmstadt, Germany)
Jason Kessler (JD Power and Associates, USA)
Michael Matuschek (Technische Universität Darmstadt, Germany)
Bjørn Zenker (Universität Erlangen, Germany)

Table of Contents

<i>Grammar Prototyping and Testing with the LinGO Grammar Matrix Customization System</i> Emily M. Bender, Scott Drellishak, Antske Fokkens, Michael Wayne Goodman, Daniel P. Mills, Laurie Poulson and Safiyah Saleem	1
<i>cdec: A Decoder, Alignment, and Learning Framework for Finite-State and Context-Free Translation Models</i> Chris Dyer, Adam Lopez, Juri Ganitkevitch, Jonathan Weese, Ferhan Ture, Phil Blunsom, Hendra Setiawan, Vladimir Eidelman and Philip Resnik	7
<i>Beetle II: A System for Tutoring and Computational Linguistics Experimentation</i> Myroslava O. Dzikovska, Johanna D. Moore, Natalie Steinhauser, Gwendolyn Campbell, Elaine Farrow and Charles B. Callaway	13
<i>GernEdiT - The GermaNet Editing Tool</i> Verena Henrich and Erhard Hinrichs	19
<i>WebLicht: Web-Based LRT Services for German</i> Erhard Hinrichs, Marie Hinrichs and Thomas Zastrow	25
<i>The S-Space Package: An Open Source Package for Word Space Models</i> David Jurgens and Keith Stevens	30
<i>Talking NPCs in a Virtual Game World</i> Tina Klüwer, Peter Adolphs, Feiyu Xu, Hans Uszkoreit and Xiwen Cheng	36
<i>An Open-Source Package for Recognizing Textual Entailment</i> Milen Kouylekov and Matteo Negri	42
<i>Personalising Speech-To-Speech Translation in the EMIME Project</i> Mikko Kurimo, William Byrne, John Dines, Philip N. Garner, Matthew Gibson, Yong Guan, Teemu Hirsimäki, Reima Karhila, Simon King, Hui Liang, Keiichiro Oura, Lakshmi Saheer, Matt Shannon, Sayaki Shiota and Jilei Tian	48
<i>Hunting for the Black Swan: Risk Mining from Text</i> Jochen Leidner and Frank Schilder	54
<i>Speech-Driven Access to the Deep Web on Mobile Devices</i> Taniya Mishra and Srinivas Bangalore	60
<i>Tools for Multilingual Grammar-Based Translation on the Web</i> Aarne Ranta, Krasimir Angelov and Thomas Hallgren	66
<i>Demonstration of a Prototype for a Conversational Companion for Reminiscing about Images</i> Yorick Wilks, Roberta Catizone, Alexiei Dingli and Weiwei Cheng	72
<i>It Makes Sense: A Wide-Coverage Word Sense Disambiguation System for Free Text</i> Zhi Zhong and Hwee Tou Ng	78

Conference Program

Tuesday, July 13, 2010

- 15:00–17:35 *Grammar Prototyping and Testing with the LinGO Grammar Matrix Customization System*
Emily M. Bender, Scott Drellishak, Antske Fokkens, Michael Wayne Goodman, Daniel P. Mills, Laurie Poulson and Safiyah Saleem
- 15:00–17:35 *cdec: A Decoder, Alignment, and Learning Framework for Finite-State and Context-Free Translation Models*
Chris Dyer, Adam Lopez, Juri Ganitkevitch, Jonathan Weese, Ferhan Ture, Phil Blunsom, Hendra Setiawan, Vladimir Eidelman and Philip Resnik
- 15:00–17:35 *Beetle II: A System for Tutoring and Computational Linguistics Experimentation*
Myroslava O. Dzikovska, Johanna D. Moore, Natalie Steinhauser, Gwendolyn Campbell, Elaine Farrow and Charles B. Callaway
- 15:00–17:35 *GernEdiT - The GermaNet Editing Tool*
Verena Henrich and Erhard Hinrichs
- 15:00–17:35 *WebLicht: Web-Based LRT Services for German*
Erhard Hinrichs, Marie Hinrichs and Thomas Zastrow
- 15:00–17:35 *The S-Space Package: An Open Source Package for Word Space Models*
David Jurgens and Keith Stevens
- 15:00–17:35 *Talking NPCs in a Virtual Game World*
Tina Klüwer, Peter Adolphs, Feiyu Xu, Hans Uszkoreit and Xiwen Cheng
- 15:00–17:35 *An Open-Source Package for Recognizing Textual Entailment*
Milen Kouylekov and Matteo Negri
- 15:00–17:35 *Personalising Speech-To-Speech Translation in the EMIME Project*
Mikko Kurimo, William Byrne, John Dines, Philip N. Garner, Matthew Gibson, Yong Guan, Teemu Hirsimäki, Reima Karhila, Simon King, Hui Liang, Keiichiro Oura, Lakshmi Saheer, Matt Shannon, Sayaki Shiota and Jilei Tian
- 15:00–17:35 *Hunting for the Black Swan: Risk Mining from Text*
Jochen Leidner and Frank Schilder
- 15:00–17:35 *Speech-Driven Access to the Deep Web on Mobile Devices*
Taniya Mishra and Srinivas Bangalore
- 15:00–17:35 *Tools for Multilingual Grammar-Based Translation on the Web*
Aarne Ranta, Krasimir Angelov and Thomas Hallgren

Tuesday, July 13, 2010 (continued)

15:00–17:35 *Demonstration of a Prototype for a Conversational Companion for Reminiscing about Images*

Yorick Wilks, Roberta Catizone, Alexiei Dingli and Weiwei Cheng

15:00–17:35 *It Makes Sense: A Wide-Coverage Word Sense Disambiguation System for Free Text*

Zhi Zhong and Hwee Tou Ng

Grammar Prototyping and Testing with the LinGO Grammar Matrix Customization System

Emily M. Bender, Scott Drellishak, Antske Fokkens, Michael Wayne Goodman,
Daniel P. Mills, Laurie Poulson, and Safiyyah Saleem

University of Washington, Seattle, Washington, USA

{ebender, sfd, goodmami, dpmills, lpoulson, ssaleem}@uw.edu,
afokkens@coli.uni-saarland.de

Abstract

This demonstration presents the LinGO Grammar Matrix grammar customization system: a repository of distilled linguistic knowledge and a web-based service which elicits a typological description of a language from the user and yields a customized grammar fragment ready for sustained development into a broad-coverage grammar. We describe the implementation of this repository with an emphasis on how the information is made available to users, including in-browser testing capabilities.

1 Introduction

This demonstration presents the LinGO Grammar Matrix grammar customization system¹ and its functionality for rapidly prototyping grammars. The LinGO Grammar Matrix project (Bender et al., 2002) is situated within the DELPH-IN² collaboration and is both a repository of reusable linguistic knowledge and a method of delivering this knowledge to a user in the form of an extensible precision implemented grammar. The stored knowledge includes both a cross-linguistic core grammar and a series of “libraries” containing analyses of cross-linguistically variable phenomena. The core grammar handles basic phrase types, semantic compositionality, and general infrastructure such as the feature geometry, while the current set of libraries includes analyses of word order, person/number/gender, tense/aspect, case, coordination, pro-drop, sentential negation, yes/no questions, and direct-inverse marking, as well as facilities for defining classes (types) of lexical entries and lexical rules which apply to those types. The grammars produced are compatible with both the grammar development tools and the

grammar-based applications produced by DELPH-IN. The grammar framework used is Head-driven Phrase Structure Grammar (HPSG) (Pollard and Sag, 1994) and the grammars map bidirectionally between surface strings and semantic representations in the format of Minimal Recursion Semantics (Copestake et al., 2005).

The Grammar Matrix project has three goals—one engineering and two scientific. The engineering goal is to reduce the cost of creating grammars by distilling the solutions developed in existing DELPH-IN grammars and making them easily available for new projects. The first scientific goal is to support grammar engineering for linguistic hypothesis testing, allowing users to quickly customize a basic grammar and use it as a medium in which to develop and test analyses of more interesting phenomena.³ The second scientific goal is to use computational methods to combine the results of typological research and formal syntactic analysis into a single resource that achieves both typological breadth (handling the known range of realizations of the phenomena analyzed) and analytical depth (producing analyses which work together to map surface strings to semantic representations) (Drellishak, 2009).

2 System Overview

Grammar customization with the LinGO Grammar Matrix consists of three primary activities: filling out the questionnaire, preliminary testing of the grammar fragment, and grammar creation.

2.1 Questionnaire

Most of the linguistic phenomena supported by the questionnaire vary across languages along multiple dimensions. It is not enough, for example,

¹<http://www.delph-in.net/matrix/customize/>

²<http://www.delph-in.net>

³Research of this type based on the Grammar Matrix includes (Crysmann, 2009) (tone change in Hausa) and (Fokkens et al., 2009) (Turkish suspended affixation).

simply to know that the target language has coordination. It is also necessary to know, among other things, what types of phrases can be coordinated, how those phrases are marked, and what patterns of marking appear in the language. Supporting a linguistic phenomenon, therefore, requires eliciting the answers to such questions from the user. The customization system elicits these answers using a detailed, web-based, typological questionnaire, then interprets the answers without human intervention and produces a grammar in the format expected by the LKB (Copestake, 2002), namely TDL (type description language).

The questionnaire is designed for linguists who want to create computational grammars of natural languages, and therefore it freely uses technical linguistic terminology, but avoids, when possible, mentioning the internals of the grammar that will be produced, although a user who intends to extend the grammar will need to become familiar with HPSG and TDL before doing so.

The questionnaire is presented to the user as a series of connected web pages. The first page the user sees (the “main page”) contains some introductory text and hyperlinks to direct the user to other sections of the questionnaire (“subpages”). Each subpage contains a set of related questions that (with some exceptions) covers the range of a single Matrix library. The actual questions in the questionnaire are represented by HTML form fields, including: text fields, check boxes, radio buttons, drop-downs, and multi-select drop-downs. The values of these form fields are stored in a “choices file”, which is the object passed on to the grammar customization stage.

2.1.1 Unbounded Content

Early versions of the customization system (Bender and Flickinger, 2005; Drellishak and Bender, 2005) only allowed a finite (and small) number of entries for things like lexical types. For instance, users were required to provide exactly one transitive verb type and one intransitive verb type. The current system has an iterator mechanism in the questionnaire that allows for repeated sections, and thus unlimited entries. These repeated sections can also be nested, which allows for much more richly structured information.

The utility of the iterator mechanism is most apparent when filling out the Lexicon subpage. Users can create an arbitrary number of lexical rule “slots”, each with an arbitrary number of

morphemes which each in turn bear any number of feature constraints. For example, the user could create a tense-agreement morphological slot, which contains multiple portmanteau morphemes each expressing some combination of tense, subject person and subject number values (e.g., French *-ez* expresses 2nd person plural subject agreement together with present tense).

The ability provided by the iterators to create unbounded content facilitates the creation of substantial grammars through the customization system. Furthermore, the system allows users to expand on some iterators while leaving others unspecified, thus modeling complex rule interactions even when it cannot cover features provided by these rules. A user can correctly model the morphotactic framework of the language using “skeletal” lexical rules—those that specify morphemes’ forms and their co-occurrence restrictions, but perhaps not their morphosyntactic features. The user can then, post-customization, augment these rules with the missing information.

2.1.2 Dynamic Content

In earlier versions of the customization system, the questionnaire was static. Not only was the number of form fields static, but the questions were the same, regardless of user input. The current questionnaire is more dynamic. When the user loads the customization system’s main page or subpages, appropriate HTML is created on the fly on the basis of the information already collected from the user as well as language-independent information provided by the system.

The questionnaire has two kinds of dynamic content: expandable lists for unbounded entry fields, and the population of drop-down selectors. The lists in an iterated section can be expanded or shortened with “Add” and “Delete” buttons near the items in question. Drop-down selectors can be automatically populated in several different ways.⁴ These dynamic drop-downs greatly lessen the amount of information the user must remember while filling out the questionnaire and can prevent the user from trying to enter an invalid value. Both of these operations occur without refreshing the page, saving time for the user.

⁴These include: the names of currently-defined features, the currently-defined values of a feature, or the values of variables that match a particular regular expression.

2.2 Validation

It makes no sense to attempt to create a consistent grammar from an empty questionnaire, an incomplete questionnaire, or a questionnaire containing contradictory answers, so the customization system first sends a user's answers through "form validation". This component places a set of arbitrarily complex constraints on the answers provided. The system insists, for example, that the user not state the language contains no determiners but then provide one in the Lexicon subpage. When a question fails form validation, it is marked with a red asterisk in the questionnaire, and if the user hovers the mouse cursor over the asterisk, a pop-up message appears describing how form validation failed. The validation component can also produce warnings (marked with red question marks) in cases where the system can generate a grammar from the user's answers, but we have reason to believe the grammar won't behave as expected. This occurs, for example, when there are no verbal lexical entries provided, yielding a grammar that cannot parse any sentences.

2.3 Creating a Grammar

After the questionnaire has passed validation, the system enables two more buttons on the main page: "Test by Generation" and "Create Grammar". "Test by Generation" allows the user to test the performance of the current state of the grammar without leaving the browser, and is described in §3. "Create Grammar" causes the customization system to output an LKB-compatible grammar that includes all the types in the core Matrix, along with the types from each library, tailored appropriately, according to the specific answers provided for the language described in the questionnaire.

2.4 Summary

This section has briefly presented the structure of the customization system. While we anticipate some future improvements (e.g., visualization tools to assist with designing type hierarchies and morphotactic dependencies), we believe that this system is sufficiently general to support the addition of analyses of many different linguistic phenomena. The system has been used to create starter grammars for more than 40 languages in the context of a graduate grammar engineering course.

To give sense of the size of the grammars produced by the customization system, Table 1

compares the English Resource Grammar (ERG) (Flickinger, 2000), a broad-coverage precision grammar in the same framework under development since 1994, to 11 grammars produced with the customization system by graduate students in a grammar engineering class at the University of Washington. The students developed these grammars over three weeks using reference materials and the customization system. We compare the grammars in terms of the number types they define, as well as the number of lexical rule and phrase structure rule instances.⁵ We separate types defined in the Matrix core grammar from language-specific types defined by the customization system. Not all of the Matrix-provided types are used in the definition of the language-specific rules, but they are nonetheless an important part of the grammar, serving as the foundation for further hand-development. The Matrix core grammar includes a larger number of types whose function is to provide disjunctions of parts of speech. These are given in Table 1, as "head types". The final column in the table gives the number of "choices" or specifications that the users gave to the customization system in order to derive these grammars.

3 Test-by-generation

The purpose of the test-by-generation feature is to provide a quick method for testing the grammar compiled from a choices file. It accomplishes this by generating sentences the grammar deems grammatical. This is useful to the user in two main ways: it quickly shows whether any ungrammatical sentences are being licensed by the grammar and, by providing an exhaustive list of licensed sentences for an input template, allows users to see if an expected sentence is not being produced.

It is worth emphasizing that this feature of the customization system relies on the bidirectionality of the grammars; that is, the fact that the same grammar can be used for both parsing and generation. Our experience has shown that grammar developers quickly find generation provides a more stringent test than parsing, especially for the ability of a grammar to model ungrammaticality.

3.1 Underspecified MRS

Testing by generation takes advantage of the generation algorithm include in the LKB (Carroll et al.,

⁵Serious lexicon development is taken as a separate task and thus lexicon size is not included in the table.

Language	Family	Lg-specific types	Matrix types	Head types	Lex rules	Phrasal rules	Choices
ERG	Germanic	3654	N/A	N/A	71	226	N/A
Breton	Celtic	220	413	510	57	49	1692
Cherokee	Iroquoian	182	413	510	95	27	985
French	Romance	137	413	510	29	22	740
Jamamadí	Arauan	188	413	510	87	11	1151
Lushootseed	Salish	95	413	510	20	8	391
Nishnaabemwin	Algonquian	289	413	510	124	50	1754
Pashto	Iranian	234	413	510	86	19	1839
Pali	Indo-Aryan	237	413	510	92	55	1310
Russian	Slavic	190	413	510	56	35	993
Shona	Bantu	136	413	510	51	9	591
Vietnamese	Austro-Asiatic	105	413	510	2	26	362
Average		182.9	413	510	63.5	28.3	1073.5

Table 1: Grammar sizes in comparison to ERG

1999). This algorithm takes input in the form of Minimal Recursion Semantics (MRS) (Copestake et al., 2005): a bag of elementary predications, each bearing features encoding a predicate string, a label, and one or more argument positions that can be filled with variables or with labels of other elementary predications.⁶ Each variable can further bear features encoding “variable properties” such as tense, aspect, mood, sentential force, person, number or gender.

In order to test our starter grammars by generation, therefore, we must provide input MRSSs. The shared core grammar ensures that all of the grammars produce and interpret valid MRSSs, but there are still language-specific properties in these semantic representations. Most notably, the predicate strings are user-defined (and language-specific), as are the variable properties. In addition, some coarser-grained typological properties (such as the presence or absence of determiners) lead to differences in the semantic representations. Therefore, we cannot simply store a set of MRSSs from one grammar to use as input to the generator.

Instead, we take a set of stored template MRSSs and generalize them by removing all variable properties (allowing the generator to explore all possible values), leaving only the predicate strings and links between the elementary predications. We then replace the stored predicate strings with ones selected from among those provided by the user. Figure 1a shows an MRS produced by a grammar fragment for English. Figure 1b shows the MRS with the variable properties removed and the predicate strings replaced with generic place-holders. One such template is needed for every sentence type (e.g., intransitive, transitive,

- a. $\langle h1,e2, \{h7:_{cat_n_rel}(x4:SG:THIRD), h3:exist_q_rel(x4, h5, h6), h1:_{sleep_v_rel}(e2:PRES, x4)\}, \{h5 \text{ req } h7\} \rangle$
- b. $\langle h1,e2, \{h7:#NOUN1\#(x4), h3:#DET1\#(x4, h5, h6), h1:#VERB\#(e2, x4)\}, \{h5 \text{ req } h7\} \rangle$

Figure 1: Original and underspecified MRS

negated-intransitive, etc.). In order to ensure that the generated strings are maximally informative to the user testing a grammar, we take advantage of the lexical type system. Because words in lexical types as defined by the customization system differ only in orthography and predicate string, and not in syntactic behavior, we need only consider one word of each type. This allows us to focus the range of variation produced by the generator on (a) the differences between lexical types and (b) the variable properties.

3.2 Test by generation process

The first step of the test-by-generation process is to compile the choices file into a grammar. Next, a copy of the LKB is initialized on the web server that is hosting the Matrix system, and the newly-created grammar is loaded into this LKB session.

We then construct the underspecified MRSSs in order to generate from them. To do this, the process needs to find the proper predicates to use for verbs, nouns, determiners, and any other parts of speech that a given MRS template may require. For nouns and determiners, the choices file is searched for the predicate for one noun of each lexical noun type, all of the determiner predicates, and whether or not each noun type needs a determiner or not. For verbs, the process is more complicated, requiring valence information as well as predicate strings in order to select the correct MRS template. In order to get this information, the process traverses the type hierarchy above the verbal lexical

⁶This latter type of argument encodes scopal dependencies. We abstract away here from the MRS approach to scope underspecification which is nonetheless critical for its computational tractability.

types until it finds a type that gives valence information about the verb. Once the process has all of this information, it matches verbs to MRS templates and fills in appropriate predicates.

The test-by-generation process then sends these constructed MRSS to the LKB process and displays the generation results, along with a brief explanation of the input semantics that gave rise to them, in HTML for the user.⁷

4 Related Work

As stated above, the engineering goal of the Grammar Matrix is to facilitate the rapid development of large-scale precision grammars. The starter grammars output by the customization system are compatible in format and semantic representations with existing DELPH-IN tools, including software for grammar development and for applications including machine translation (Oepen et al., 2007) and robust textual entailment (Bergmair, 2008).

More broadly, the Grammar Matrix is situated in the field of multilingual grammar engineering, or the practice of developing linguistically-motivated grammars for multiple languages within a consistent framework. Other projects in this field include ParGram (Butt et al., 2002; King et al., 2005) (LFG), the CoreGram project⁸ (e.g., (Müller, 2009)) (HPSG), and the MetaGrammar project (de la Clergerie, 2005) (TAG).

To our knowledge, however, there is only one other system that elicits typological information about a language and outputs an appropriately customized implemented grammar. The system, described in (Black, 2004) and (Black and Black, 2009), is called PAWS (Parser And Writer for Syntax) and is available for download online.⁹ PAWS is being developed by SIL in the context of both descriptive (prose) grammar writing and “computer-assisted related language adaptation”, the practice of writing a text in a target language by starting with a translation of that text in a related source language and mapping the words from target to source. Accordingly, the output of PAWS consists of both a prose descriptive grammar

⁷This set-up scales well to multiple users, as the user’s interaction with the LKB is done once per customized grammar, providing output for the user to peruse as his or her leisure. The LKB process does not persist, but can be started again by reinvoking test-by-generation, such as when the user has updated the grammar definition.

⁸<http://hpsg.fu-berlin.de/Projects/core.html>

⁹http://www.sil.org/computing/catalog/show_software.asp?id=85

and an implemented grammar. The latter is in the format required by PC-PATR (McConnel, 1995), and is used primarily to disambiguate morphological analyses of lexical items in the input string.

Other systems that attempt to elicit linguistic information from a user include the Expedition (McShane and Nirenburg, 2003) and Avenue projects (Monson et al., 2008), which are specifically targeted at developing machine translation for low-density languages. These projects differ from the Grammar Matrix customization system in eliciting information from native speakers (such as paradigms or translations of specifically tailored corpora), rather than linguists. Further, unlike the Grammar Matrix customization system, they do not produce resources meant to sustain further development by a linguist.

5 Demonstration Plan

Our demonstration illustrates how the customization system can be used to create starter grammars and test them by invoking test-by-generation. We first walk through the questionnaire to illustrate the functionality of libraries and the way that the user interacts with the system to enter information. Then, using a sample grammar for English, we demonstrate how test-by-generation can expose both overgeneration (ungrammatical generated strings) and undergeneration (gaps in generated paradigms). Finally, we return to the questionnaire to address the bugs in the sample grammar and retest to show the result.

6 Conclusion

This paper has presented an overview of the LinGO Grammar Matrix Customization System, highlighting the ways in which it provides access to its repository of linguistic knowledge. The current customization system covers a sufficiently wide range of phenomena that the grammars it produces are non-trivial. In addition, it is not always apparent to a user what the implications will be of selecting various options in the questionnaire, nor how analyses of different phenomena will interact. The test-by-generation methodology allows users to interactively explore the consequences of different linguistic analyses within the platform. We anticipate that it will, as a result, encourage users to develop more complex grammars within the customization system (before moving on to hand-editing) and thereby gain more benefit.

Acknowledgments

This material is based upon work supported by the National Science Foundation under Grant No. 0644097. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation.

References

- Emily M. Bender and Dan Flickinger. 2005. Rapid prototyping of scalable grammars: Towards modularity in extensions to a language-independent core. In *Proc. of IJCNLP-05 (Posters/Demos)*.
- Emily M. Bender, Dan Flickinger, and Stephan Oepen. 2002. The grammar matrix: An open-source starter-kit for the rapid development of cross-linguistically consistent broad-coverage precision grammars. In *Proc. of the Workshop on Grammar Engineering and Evaluation at COLING 2002*, pages 8–14.
- Richard Bergmair. 2008. Monte Carlo semantics: McPIET at RTE4. In *Text Analysis Conference (TAC 2008) Workshop-RTE-4 Track. National Institute of Standards and Technology*, pages 17–19.
- Cheryl A. Black and H. Andrew Black. 2009. PAWS: Parser and writer for syntax: Drafting syntactic grammars in the third wave. In *SIL Forum for Language Fieldwork*, volume 2.
- Cheryl A. Black. 2004. Parser and writer for syntax. Paper presented at the International Conference on Translation with Computer-Assisted Technology: Changes in Research, Teaching, Evaluation, and Practice, University of Rome “La Sapienza”, April 2004.
- Miriam Butt, Helge Dyvik, Tracy Holloway King, Hiroshi Masuichi, and Christian Rohrer. 2002. The parallel grammar project. In *Proc. of the Workshop on Grammar Engineering and Evaluation at COLING 2002*, pages 1–7.
- John Carroll, Ann Copestake, Dan Flickinger, and Victor Poznański. 1999. An efficient chart generator for (semi-) lexicalist grammars. In *Proc. of the 7th European workshop on natural language generation (EWNLG99)*, pages 86–95.
- Ann Copestake, Dan Flickinger, Carl Pollard, and Ivan A. Sag. 2005. Minimal recursion semantics: An introduction. *Research on Language & Computation*, 3(4):281–332.
- Ann Copestake. 2002. *Implementing Typed Feature Structure Grammars*. CSLI, Stanford.
- Berthold Crysmann. 2009. Autosegmental representations in an HPSG for Hausa. In *Proc. of the Workshop on Grammar Engineering Across Frameworks 2009*.
- Éric Villemonte de la Clergerie. 2005. From meta-grammars to factorized TAG/TIG parsers. In *Proc. of IWPT’05*, pages 190–191.
- Scott Drellishak and Emily M. Bender. 2005. A coordination module for a crosslinguistic grammar resource. In Stefan Müller, editor, *Proc. of HPSG 2005*, pages 108–128, Stanford. CSLI.
- Scott Drellishak. 2009. *Widespread But Not Universal: Improving the Typological Coverage of the Grammar Matrix*. Ph.D. thesis, University of Washington.
- Dan Flickinger. 2000. On building a more efficient grammar by exploiting types. *Natural Language Engineering*, 6:15–28.
- Antske Fokkens, Laurie Poulson, and Emily M. Bender. 2009. Inflectional morphology in Turkish VP-coordination. In Stefan Müller, editor, *Proc. of HPSG 2009*, pages 110–130, Stanford. CSLI.
- Tracy Holloway King, Martin Forst, Jonas Kuhn, and Miriam Butt. 2005. The feature space in parallel grammar writing. *Research on Language & Computation*, 3(2):139–163.
- Stephen McConnel. 1995. *PC-PATR Reference Manual*. Summer Institute for Linguistics. <http://www.sil.org/pcpatr/manual/pcpatr.html>.
- Marjorie McShane and Sergei Nirenburg. 2003. Parameterizing and eliciting text elements across languages for use in natural language processing systems. *Machine Translation*, 18:129–165.
- Christian Monson, Ariadna Font Llitjts, Vamshi Ambati, Lori Levin, Alon Lavie, Alison Alvarez, Roberto Aranovich, Jaime Carbonell, Robert Frederick, Erik Peterson, and Katharina Probst. 2008. Linguistic structure and bilingual informants help induce machine translation of lesser-resourced languages. In *LREC’08*.
- Stefan Müller. 2009. Towards an HPSG analysis of Maltese. In Bernard Comrie, Ray Fabri, Beth Hume, Manwel Mifsud, Thomas Stolz, and Martine Vanhove, editors, *Introducing Maltese linguistics. Papers from the 1st International Conference on Maltese Linguistics*, pages 83–112. Benjamins, Amsterdam.
- Stephan Oepen, Erik Velldal, Jan Tore Lning, Paul Meurer, Victoria Rosn, and Dan Flickinger. 2007. Towards hybrid quality-oriented machine translation. On linguistics and probabilities in MT. In *11th International Conference on Theoretical and Methodological Issues in Machine Translation*.
- Carl Pollard and Ivan A. Sag. 1994. *Head-Driven Phrase Structure Grammar*. The University of Chicago Press, Chicago, IL.

cdec: A Decoder, Alignment, and Learning Framework for Finite-State and Context-Free Translation Models

Chris Dyer

University of Maryland
redpony@umd.edu

Adam Lopez

University of Edinburgh
alopez@inf.ed.ac.uk

Juri Ganitkevitch

Johns Hopkins University
juri@cs.jhu.edu

Jonathan Weese

Johns Hopkins University
jweese@cs.jhu.edu

Ferhan Ture

University of Maryland
fture@cs.umd.edu

Phil Blunsom

Oxford University
pblunsom@comlab.ox.ac.uk

Hendra Setiawan

University of Maryland
hendra@umiacs.umd.edu

Vladimir Eidelman

University of Maryland
vlad@umiacs.umd.edu

Philip Resnik

University of Maryland
resnik@umiacs.umd.edu

Abstract

We present `cdec`, an open source framework for decoding, aligning with, and training a number of statistical machine translation models, including word-based models, phrase-based models, and models based on synchronous context-free grammars. Using a single unified internal representation for translation forests, the decoder strictly separates model-specific translation logic from general rescoring, pruning, and inference algorithms. From this unified representation, the decoder can extract not only the 1- or k -best translations, but also alignments to a reference, or the quantities necessary to drive discriminative training using gradient-based or gradient-free optimization techniques. Its efficient C++ implementation means that memory use and runtime performance are significantly better than comparable decoders.

1 Introduction

The dominant models used in machine translation and sequence tagging are formally based on either weighted finite-state transducers (FSTs) or weighted synchronous context-free grammars (SCFGs) (Lopez, 2008). Phrase-based models (Koehn et al., 2003), lexical translation models (Brown et al., 1993), and finite-state conditional random fields (Sha and Pereira, 2003) exemplify the former, and hierarchical phrase-based models the latter (Chiang, 2007). We introduce a software package called `cdec` that manipulates both

classes in a unified way.¹

Although open source decoders for both phrase-based and hierarchical translation models have been available for several years (Koehn et al., 2007; Li et al., 2009), their extensibility to new models and algorithms is limited by two significant design flaws that we have avoided with `cdec`. First, their implementations tightly couple the translation, language model integration (which we call *rescoring*), and pruning algorithms. This makes it difficult to explore alternative translation models without also re-implementing rescoring and pruning logic. In `cdec`, model-specific code is only required to construct a translation forest (§3). *General* rescoring (with language models or other models), pruning, inference, and alignment algorithms then apply to the unified data structure (§4). Hence *all* model types benefit immediately from new algorithms (for rescoring, inference, etc.); new models can be more easily prototyped; and controlled comparison of models is made easier.

Second, existing open source decoders were designed with the traditional phrase-based parameterization using a very small number of dense features (typically less than 10). `cdec` has been designed from the ground up to support any parameterization, from those with a handful of dense features up to models with millions of sparse features (Blunsom et al., 2008; Chiang et al., 2009). Since the inference algorithms necessary to compute a training objective (e.g. conditional likelihood or expected BLEU) and its gradient operate on the unified data structure (§5), any model type can be trained using with any of the supported training

¹The software is released under the Apache License, version 2.0, and is available from <http://cdec-decoder.org/>.

criteria. The software package includes general function optimization utilities that can be used for discriminative training (§6).

These features are implemented without compromising on performance. We show experimentally that `cdec` uses less memory and time than comparable decoders on a controlled translation task (§7).

2 Decoder workflow

The decoding pipeline consists of two phases. The first (Figure 1) transforms input, which may be represented as a source language sentence, lattice (Dyer et al., 2008), or context-free forest (Dyer and Resnik, 2010), into a translation forest that has been rescored with all applicable models.

In `cdec`, the only model-specific logic is confined to the first step in the process where an input string (or lattice, etc.) is transduced into the unified hypergraph representation. Since the model-specific code need not worry about integration with rescoring models, it can be made quite simple and efficient. Furthermore, prior to language model integration (and distortion model integration, in the case of phrase based translation), pruning is unnecessary for most kinds of models, further simplifying the model-specific code. Once this unscored translation forest has been generated, any non-coaccessible states (i.e., states that are not reachable from the goal node) are removed and the resulting structure is rescored with language models using a user-specified intersection/pruning strategy (§4) resulting in a rescored translation forest and completing phase 1.

The second phase of the decoding pipeline (depicted in Figure 2) computes a value from the rescored forest: 1- or k -best derivations, feature expectations, or intersection with a target language reference (sentence or lattice). The last option generates an *alignment forest*, from which a word alignment or feature expectations can be extracted. Most of these values are computed in a time complexity that is linear in the number of edges and nodes in the translation hypergraph using `cdec`'s semiring framework (§5).

2.1 Alignment forests and alignment

Alignment is the process of determining if and how a translation model generates a $\langle source, target \rangle$ string pair. To compute an alignment under a translation model, the phase 1 translation hypergraph is reinterpreted as a synchronous context-

free grammar and then used to parse the *target* sentence.² This results in an *alignment forest*, which is a compact representation of all the derivations of the sentence pair under the translation model. From this forest, the Viterbi or maximum *a posteriori* word alignment can be generated. This alignment algorithm is explored in depth by Dyer (2010). Note that if the phase 1 forest has been pruned in some way, or the grammar does not derive the sentence pair, the target intersection parse may fail, meaning that an alignment will not be recoverable.

3 Translation hypergraphs

Recent research has proposed a unified representation for the various translation and tagging formalisms that is based on weighted logic programming (Lopez, 2009). In this view, translation (or tagging) deductions have the structure of a *context-free forest*, or directed hypergraph, where edges have a single head and 0 or more tail nodes (Nederhof, 2003). Once a forest has been constructed representing the possible translations, general inference algorithms can be applied.

In `cdec`'s translation hypergraph, a node represents a contiguous sequence of target language words. For SCFG models and sequential tagging models, a node also corresponds to a source span and non-terminal type, but for word-based and phrase-based models, the relationship to the source string (or lattice) may be more complicated. In a phrase-based translation hypergraph, the node will correspond to a source *coverage vector* (Koehn et al., 2003). In word-based models, a single node may derive multiple different source language coverages since word based models impose no requirements on covering all words in the input. Figure 3 illustrates two example hypergraphs, one generated using a SCFG model and other from a phrase-based model.

Edges are associated with exactly one synchronous production in the source and target language, and alternative translation possibilities are expressed as alternative edges. Edges are further annotated with feature values, and are annotated with the source span vector the edge corresponds to. An edge's output label may contain mixtures of terminal symbol yields and positions indicating where a child node's yield should be substituted.

²The parser is smart enough to detect the left-branching grammars generated by lexical translation and tagging models, and use a more efficient intersection algorithm.

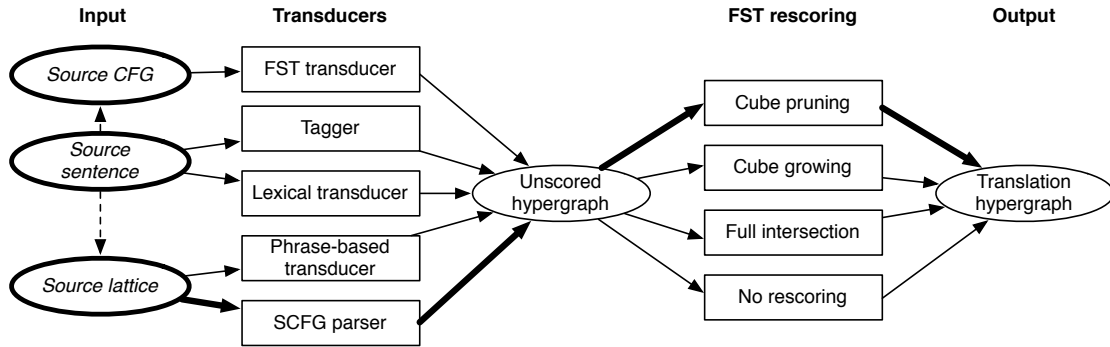


Figure 1: Forest generation workflow (first half of decoding pipeline). The decoder’s configuration specifies what path is taken from the input (one of the bold ovals) to a unified translation hypergraph. The highlighted path is the workflow used in the test reported in §7.

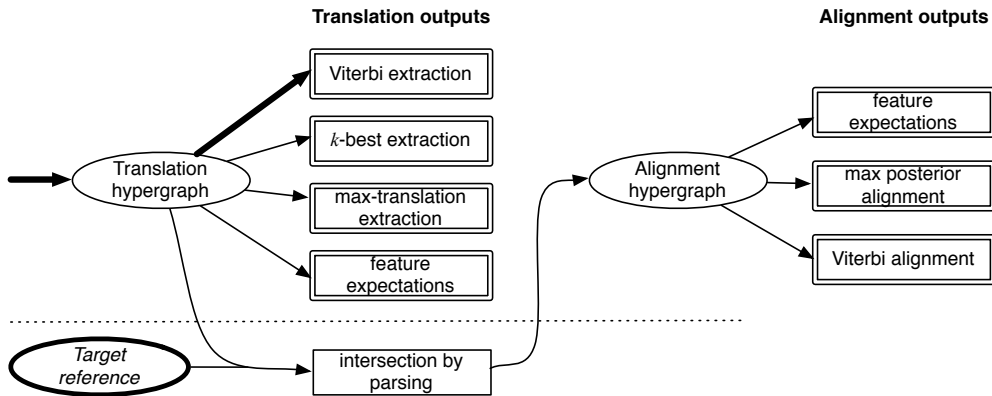


Figure 2: Output generation workflow (second half of decoding pipeline). Possible output types are designated with a double box.

In the case of SCFG grammars, the edges correspond simply to rules in the synchronous grammar. For non-SCFG translation models, there are two kinds of edges. The first have zero tail nodes (i.e., an arity of 0), and correspond to word or phrase translation pairs (with all translation options existing on edges deriving the same head node), or *glue rules* that glue phrases together. For tagging, word-based, and phrase-based models, these are strictly arranged in a monotone, left-branching structure.

4 Rescoring with weighted FSTs

The design of `cdec` separates the creation of a translation forest from its rescoring with a language models or similar models.³ Since the structure of the unified search space is context free (§3), we use the logic for language model rescoring described by Chiang (2007), although any weighted intersection algorithm can be applied. The rescoring

³Other rescoring models that depend on sequential context include distance-based reordering models or Markov features in tagging models.

models need not be explicitly represented as FSTs—the state space can be inferred.

Although intersection using the Chiang algorithm runs in polynomial time and space, the resulting rescored forest may still be too large to represent completely. `cdec` therefore supports three pruning strategies that can be used during intersection: full unpruned intersection (useful for tagging models to incorporate, e.g., Markov features, but not generally practical for translation), cube pruning, and cube growing (Huang and Chiang, 2007).

5 Semiring framework

Semirings are a useful mathematical abstraction for dealing with translation forests since many useful quantities can be computed using a single linear-time algorithm but with different semirings. A semiring is a 5-tuple $(\mathbb{K}, \oplus, \otimes, \bar{0}, \bar{1})$ that indicates the set from which the values will be drawn, \mathbb{K} , a generic addition and multiplication operation, \oplus and \otimes , and their identities $\bar{0}$ and $\bar{1}$. Multiplication and addition must be associative. Multiplication must distribute over addition, and $v \otimes \bar{0}$

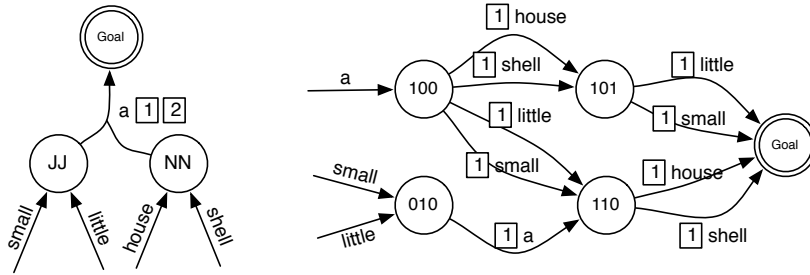


Figure 3: Example unrescored translation hypergraphs generated for the German input *ein* (a) *kleines* (small/little) *Haus* (house/shell) using a SCFG-based model (left) and phrase-based model with a distortion limit of 1 (right).

must equal $\bar{0}$. Values that can be computed using the semirings include the number of derivations, the expected translation length, the entropy of the translation posterior distribution, and the expected values of feature functions (Li and Eisner, 2009).

Since semirings are such a useful abstraction, `cdec` has been designed to facilitate implementation of new semirings. Table 1 shows the C++ representation used for semirings. Note that because of our representation, built-in types like `double`, `int`, and `bool` (together with their default operators) are semirings. Beyond these, the type `prob_t` is provided which stores the logarithm of the value it represents, which helps avoid underflow and overflow problems that may otherwise be encountered. A generic first-order expectation semiring is also provided (Li and Eisner, 2009).

Table 1: Semiring representation. `T` is a C++ type name.

Element	C++ representation
\mathbb{K}	<code>T</code>
\oplus	<code>T::operator+=</code>
\otimes	<code>T::operator*+=</code>
$\bar{0}$	<code>T()</code>
$\bar{1}$	<code>T(1)</code>

Three standard algorithms parameterized with semirings are provided: `INSIDE`, `OUTSIDE`, and `INSIDEOUTSIDE`, and the semiring is specified using C++ generics (templates). Additionally, each algorithm takes a *weight function* that maps from hypergraph edges to a value in \mathbb{K} , making it possible to use many different semirings without altering the underlying hypergraph.

5.1 Viterbi and k -best extraction

Although Viterbi and k -best extraction algorithms are often expressed as `INSIDE` algorithms with

the tropical semiring, `cdec` provides a separate derivation extraction framework that makes use of a `<` operator (Huang and Chiang, 2005). Thus, many of the semiring types define not only the elements shown in Table 1 but `T::operator<` as well. The k -best extraction algorithm is also parameterized by an optional predicate that can filter out derivations at each node, enabling extraction of only derivations that yield different strings as in Huang et al. (2006).

6 Model training

Two training pipelines are provided with `cdec`. The first, called Viterbi envelope semiring training, `VEST`, implements the minimum error rate training (MERT) algorithm, a gradient-free optimization technique capable of maximizing arbitrary loss functions (Och, 2003).

6.1 VEST

Rather than computing an error surface using k -best approximations of the decoder search space, `cdec`'s implementation performs inference over the full hypergraph structure (Kumar et al., 2009). In particular, by defining a semiring whose values are sets of line segments, having an addition operation equivalent to union, and a multiplication operation equivalent to a linear transformation of the line segments, Och's line search can be computed simply using the `INSIDE` algorithm. Since the translation hypergraphs generated by `cdec` may be quite large making inference expensive, the logic for constructing error surfaces is factored according to the MapReduce programming paradigm (Dean and Ghemawat, 2004), enabling parallelization across a cluster of machines. Implementations of the BLEU and TER loss functions are provided (Papineni et al., 2002; Snover et al., 2006).

6.2 Large-scale discriminative training

In addition to the widely used MERT algorithm, `cdec` also provides a training pipeline for discriminatively trained probabilistic translation models (Blunsom et al., 2008; Blunsom and Osborne, 2008). In these models, the translation model is trained to maximize conditional log likelihood of the training data under a specified grammar. Since log likelihood is differentiable with respect to the feature weights in an exponential model, it is possible to use gradient-based optimization techniques to train the system, enabling the parameterization of the model using millions of sparse features. While this training approach was originally proposed for SCFG-based translation models, it can be used to train *any* model type in `cdec`. When used with sequential tagging models, this pipeline is identical to traditional sequential CRF training (Sha and Pereira, 2003).

Both the objective (conditional log likelihood) and its gradient have the form of a difference in two quantities: each has one term that is computed over the *translation* hypergraph which is subtracted from the result of the same computation over the *alignment* hypergraph (refer to Figures 1 and 2). The conditional log likelihood is the difference in the log partition of the translation and alignment hypergraph, and is computed using the INSIDE algorithm. The gradient with respect to a particular feature is the difference in this feature’s expected value in the translation and alignment hypergraphs, and can be computed using either INSIDEOUTSIDE or the expectation semiring and INSIDE. Since a translation forest is generated as an intermediate step in generating an alignment forest (§2) this computation is straightforward.

Since gradient-based optimization techniques may require thousands of evaluations to converge, the batch training pipeline is split into map and reduce components, facilitating distribution over very large clusters. Briefly, the `cdec` is run as the map function, and sentence pairs are mapped over. The reduce function aggregates the results and performs the optimization using standard algorithms, including LBFGS (Liu et al., 1989), RPROP (Riedmiller and Braun, 1993), and stochastic gradient descent.

7 Experiments

Table 2 compares the performance of `cdec`, Hiero, and Joshua 1.3 (running with 1 or 8 threads) decoding using a hierarchical phrase-based trans-

lation grammar and identical pruning settings.⁴ Figure 4 shows the `cdec` configuration and weights file used for this test.

The workstation used has two 2GHz quad-core Intel Xenon processors, 32GB RAM, is running Linux kernel version 2.6.18 and gcc version 4.1.2. All decoders use SRI’s language model toolkit, version 1.5.9 (Stolcke, 2002). Joshua was run on the Sun HotSpot JVM, version 1.6.0_12. A hierarchical phrase-based translation grammar was extracted for the NIST MT03 Chinese-English translation using a suffix array rule extractor (Lopez, 2007). A non-terminal span limit of 15 was used, and all decoders were configured to use cube pruning with a limit of 30 candidates at each node and no further pruning. All decoders produced a BLEU score between 31.4 and 31.6 (small differences are accounted for by different tie-breaking behavior and OOV handling).

Table 2: Memory usage and average per-sentence running time, in seconds, for decoding a Chinese-English test set.

Decoder	Lang.	Time (s)	Memory
<code>cdec</code>	C++	0.37	1.0Gb
Joshua (1×)	Java	0.98	1.5Gb
Joshua (8×)	Java	0.35	2.5Gb
Hiero	Python	4.04	1.1Gb

```
formalism=scfg
grammar=grammar.mt03.scfg.gz
add_pass_through_rules=true
scfg_max_span_limit=15
feature_function=LanguageModel \
    en.3gram.pruned.lm.gz -o 3
feature_function=WordPenalty
intersection_strategy=cube_pruning
cubepruning_pop_limit=30
```

```
LanguageModel 1.12
WordPenalty -4.26
PhraseModel0 0.963
PhraseModel1 0.654
PhraseModel2 0.773
PassThroughRule -20
```

Figure 4: Configuration file (above) and feature weights file (below) used for the decoding test described in §7.

⁴<http://sourceforge.net/projects/joshua/>

8 Future work

cdec continues to be under active development. We are taking advantage of its modular design to study alternative algorithms for language model integration. Further training pipelines are under development, including minimum risk training using a linearly decomposable approximation of BLEU (Li and Eisner, 2009), and MIRA training (Chiang et al., 2009). All of these will be made publicly available as the projects progress. We are also improving support for parallel training using Hadoop (an open-source implementation of MapReduce).

Acknowledgements

This work was partially supported by the GALE program of the Defense Advanced Research Projects Agency, Contract No. HR0011-06-2-001. Any opinions, findings, conclusions or recommendations expressed in this paper are those of the authors and do not necessarily reflect the views of the sponsors. Further support was provided the EuroMatrix project funded by the European Commission (7th Framework Programme). Discussions with Philipp Koehn, Chris Callison-Burch, Zhifei Li, Lane Schwarz, and Jimmy Lin were likewise crucial to the successful execution of this project.

References

- P. Blunsom and M. Osborne. 2008. Probabilistic inference for machine translation. In *Proc. of EMNLP*.
- P. Blunsom, T. Cohn, and M. Osborne. 2008. A discriminative latent variable model for statistical machine translation. In *Proc. of ACL-HLT*.
- P. F. Brown, V. J. Della Pietra, S. A. Della Pietra, and R. L. Mercer. 1993. The mathematics of statistical machine translation: parameter estimation. *Computational Linguistics*, 19(2):263–311.
- D. Chiang, K. Knight, and W. Wang. 2009. 11,001 new features for statistical machine translation. In *Proc. of NAACL*, pages 218–226.
- D. Chiang. 2007. Hierarchical phrase-based translation. *Comp. Ling.*, 33(2):201–228.
- J. Dean and S. Ghemawat. 2004. MapReduce: Simplified data processing on large clusters. In *Proc. of the 6th Symposium on Operating System Design and Implementation (OSDI 2004)*, pages 137–150.
- C. Dyer and P. Resnik. 2010. Context-free reordering, finite-state translation. In *Proc. of HLT-NAACL*.
- C. Dyer, S. Muresan, and P. Resnik. 2008. Generalizing word lattice translation. In *Proc. of HLT-ACL*.
- C. Dyer. 2010. Two monolingual parses are better than one (synchronous parse). In *Proc. of HLT-NAACL*.
- L. Huang and D. Chiang. 2005. Better k -best parsing. In *In Proc. of IWPT*, pages 53–64.
- L. Huang and D. Chiang. 2007. Forest rescoring: Faster decoding with integrated language models. In *Proc. ACL*.
- L. Huang, K. Knight, and A. Joshi. 2006. A syntax-directed translator with extended domain of locality. In *Proc. of AMTA*.
- P. Koehn, F. J. Och, and D. Marcu. 2003. Statistical phrase-based translation. In *Proc. of HLT/NAACL*, pages 48–54.
- P. Koehn, H. Hoang, A. B. Mayne, C. Callison-Burch, M. Federico, N. Bertoldi, B. Cowan, W. Shen, C. Moran, R. Zens, C. Dyer, O. Bojar, A. Constantin, and E. Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In *Proc. of ACL, Demonstration Session*, pages 177–180, June.
- S. Kumar, W. Macherey, C. Dyer, and F. Och. 2009. Efficient minimum error rate training and minimum Bayes-risk decoding for translation hypergraphs and lattices. In *Proc. of ACL*, pages 163–171.
- Z. Li and J. Eisner. 2009. First- and second-order expectation semirings with applications to minimum-risk training on translation forests. In *Proc. of EMNLP*, pages 40–51.
- Z. Li, C. Callison-Burch, C. Dyer, J. Ganitkevitch, S. Khudanpur, L. Schwartz, W. N. G. Thornton, J. Weese, and O. F. Zaidan. 2009. Joshua: an open source toolkit for parsing-based machine translation. In *Proc. of the Fourth Workshop on Stat. Machine Translation*, pages 135–139.
- D. C. Liu, J. Nocedal, D. C. Liu, and J. Nocedal. 1989. On the limited memory BFGS method for large scale optimization. *Mathematical Programming B*, 45(3):503–528.
- A. Lopez. 2007. Hierarchical phrase-based translation with suffix arrays. In *Proc. of EMNLP*, pages 976–985.
- A. Lopez. 2008. Statistical machine translation. *ACM Computing Surveys*, 40(3), Aug.
- A. Lopez. 2009. Translation as weighted deduction. In *Proc. of EACL*, pages 532–540.
- M.-J. Nederhof. 2003. Weighted deductive parsing and Knuth’s algorithm. *Comp. Ling.*, 29(1):135–143, Mar.
- F. Och. 2003. Minimum error rate training in statistical machine translation. In *Proc. of ACL*, pages 160–167.
- K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In *Proc. of ACL*, pages 311–318.
- M. Riedmiller and H. Braun. 1993. A direct adaptive method for faster backpropagation learning: The RPROP algorithm. In *Proc. of the IEEE international conference on neural networks*, pages 586–591.
- F. Sha and F. Pereira. 2003. Shallow parsing with conditional random fields. In *Proc. of NAACL*, pages 134–141.
- M. Snover, B. Dorr, R. Schwartz, L. Micciulla, and J. Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proc. AMTA*.
- A. Stolcke. 2002. SRILM – an extensible language modeling toolkit. In *Intl. Conf. on Spoken Language Processing*.

BEETLE II: a system for tutoring and computational linguistics experimentation

Myroslava O. Dzikovska and **Johanna D. Moore**

School of Informatics, University of Edinburgh, Edinburgh, United Kingdom
{m.dzikovska, j.moore}@ed.ac.uk

Natalie Steinhauser and **Gwendolyn Campbell**

Naval Air Warfare Center Training Systems Division, Orlando, FL, USA
{gwendolyn.campbell, natalie.steihauser}@navy.mil

Elaine Farrow

Heriot-Watt University
Edinburgh, United Kingdom
e.farrow@hw.ac.uk

Charles B. Callaway

University of Haifa
Mount Carmel, Haifa, Israel
ccallawa@gmail.com

Abstract

We present BEETLE II, a tutorial dialogue system designed to accept unrestricted language input and support experimentation with different tutorial planning and dialogue strategies. Our first system evaluation used two different tutorial policies and demonstrated that the system can be successfully used to study the impact of different approaches to tutoring. In the future, the system can also be used to experiment with a variety of natural language interpretation and generation techniques.

1 Introduction

Over the last decade there has been a lot of interest in developing tutorial dialogue systems that understand student explanations (Jordan et al., 2006; Graesser et al., 1999; Alevan et al., 2001; Buckley and Wolska, 2007; Nielsen et al., 2008; VanLehn et al., 2007), because high percentages of self-explanation and student contentful talk are known to be correlated with better learning in human-human tutoring (Chi et al., 1994; Litman et al., 2009; Purandare and Litman, 2008; Steinhauser et al., 2007). However, most existing systems use pre-authored tutor responses for addressing student errors. The advantage of this approach is that tutors can devise remediation dialogues that are highly tailored to specific misconceptions many students share, providing step-by-step scaffolding and potentially suggesting additional problems. The disadvantage is a lack of adaptivity and generality: students often get the same remediation for the same error regardless of their past performance or dialogue context, as it is infeasible to

author a different remediation dialogue for every possible dialogue state. It also becomes more difficult to experiment with different tutorial policies within the system due to the inherent complexities in applying tutoring strategies consistently across a large number of individual hand-authored remediations.

The BEETLE II system architecture is designed to overcome these limitations (Callaway et al., 2007). It uses a deep parser and generator, together with a domain reasoner and a diagnoser, to produce detailed analyses of student utterances and generate feedback automatically. This allows the system to consistently apply the same tutorial policy across a range of questions. To some extent, this comes at the expense of being able to address individual student misconceptions. However, the system's modular setup and extensibility make it a suitable testbed for both computational linguistics algorithms and more general questions about theories of learning.

A distinguishing feature of the system is that it is based on an introductory electricity and electronics course developed by experienced instructional designers. The course was first created for use in a human-human tutoring study, without taking into account possible limitations of computer tutoring. The exercises were then transferred into a computer system with only minor adjustments (e.g., breaking down compound questions into individual questions). This resulted in a realistic tutoring setup, which presents interesting challenges to language processing components, involving a wide variety of language phenomena.

We demonstrate a version of the system that has undergone a successful user evaluation in

2009. The evaluation results indicate that additional improvements to remediation strategies, and especially to strategies dealing with interpretation problems, are necessary for effective tutoring. At the same time, the successful large-scale evaluation shows that BEETLE II can be used as a platform for future experimentation.

The rest of this paper discusses the BEETLE II system architecture (Section 2), system evaluation (Section 3), and the range of computational linguistics problems that can be investigated using BEETLE II (Section 4).

2 System Architecture

The BEETLE II system delivers basic electricity and electronics tutoring to students with no prior knowledge of the subject. A screenshot of the system is shown in Figure 1. The student interface includes an area to display reading material, a circuit simulator, and a dialogue history window. All interactions with the system are typed. Students read pre-authored curriculum slides and carry out exercises which involve experimenting with the circuit simulator and explaining the observed behavior. The system also asks some high-level questions, such as “What is voltage?”.

The system architecture is shown in Figure 2. The system uses a standard interpretation pipeline, with domain-independent parsing and generation components supported by domain specific reasoners for decision making. The architecture is discussed in detail in the rest of this section.

2.1 Interpretation Components

We use the TRIPS dialogue parser (Allen et al., 2007) to parse the utterances. The parser provides a domain-independent semantic representation including high-level word senses and semantic role labels. The contextual interpreter then uses a reference resolution approach similar to Byron (2002), and an ontology mapping mechanism (Dzikovska et al., 2008a) to produce a domain-specific semantic representation of the student’s output. Utterance content is represented as a set of extracted objects and relations between them. Negation is supported, together with a heuristic scoping algorithm. The interpreter also performs basic ellipsis resolution. For example, it can determine that in the answer to the question “Which bulbs will be on and which bulbs will be off in this diagram?”, “off” can be taken to mean “all bulbs in the di-

agram will be off.” The resulting output is then passed on to the domain reasoning and diagnosis components.

2.2 Domain Reasoning and Diagnosis

The system uses a knowledge base implemented in the KM representation language (Clark and Porter, 1999; Dzikovska et al., 2006) to represent the state of the world. At present, the knowledge base represents 14 object types and supports the curriculum containing over 200 questions and 40 different circuits.

Student explanations are checked on two levels, verifying *factual* and *explanation* correctness. For example, for a question “Why is bulb A lit?”, if the student says “it is in a closed path”, the system checks two things: a) is the bulb indeed in a closed path? and b) is being in a closed path a reasonable explanation for the bulb being lit? Different remediation strategies need to be used depending on whether the student made a factual error (i.e., they misread the diagram and the bulb is not in a closed path) or produced an incorrect explanation (i.e., the bulb is indeed in a closed path, but they failed to mention that a battery needs to be in the same closed path for the bulb to light).

The knowledge base is used to check the factual correctness of the answers first, and then a diagnoser checks the explanation correctness. The diagnoser, based on Dzikovska et al. (2008b), outputs a diagnosis which consists of lists of correct, contradictory and non-mentioned objects and relations from the student’s answer. At present, the system uses a heuristic matching algorithm to classify relations into the appropriate category, though in the future we may consider a classifier similar to Nielsen et al. (2008).

2.3 Tutorial Planner

The tutorial planner implements a set of generic tutoring strategies, as well as a policy to choose an appropriate strategy at each point of the interaction. It is designed so that different policies can be defined for the system. The currently implemented strategies are: acknowledging the correct part of the answer; suggesting a slide to read with background material; prompting for missing parts of the answer; hinting (low- and high- specificity); and giving away the answer. Two or more strategies can be used together if necessary.

The hint selection mechanism generates hints automatically. For a low specificity hint it selects

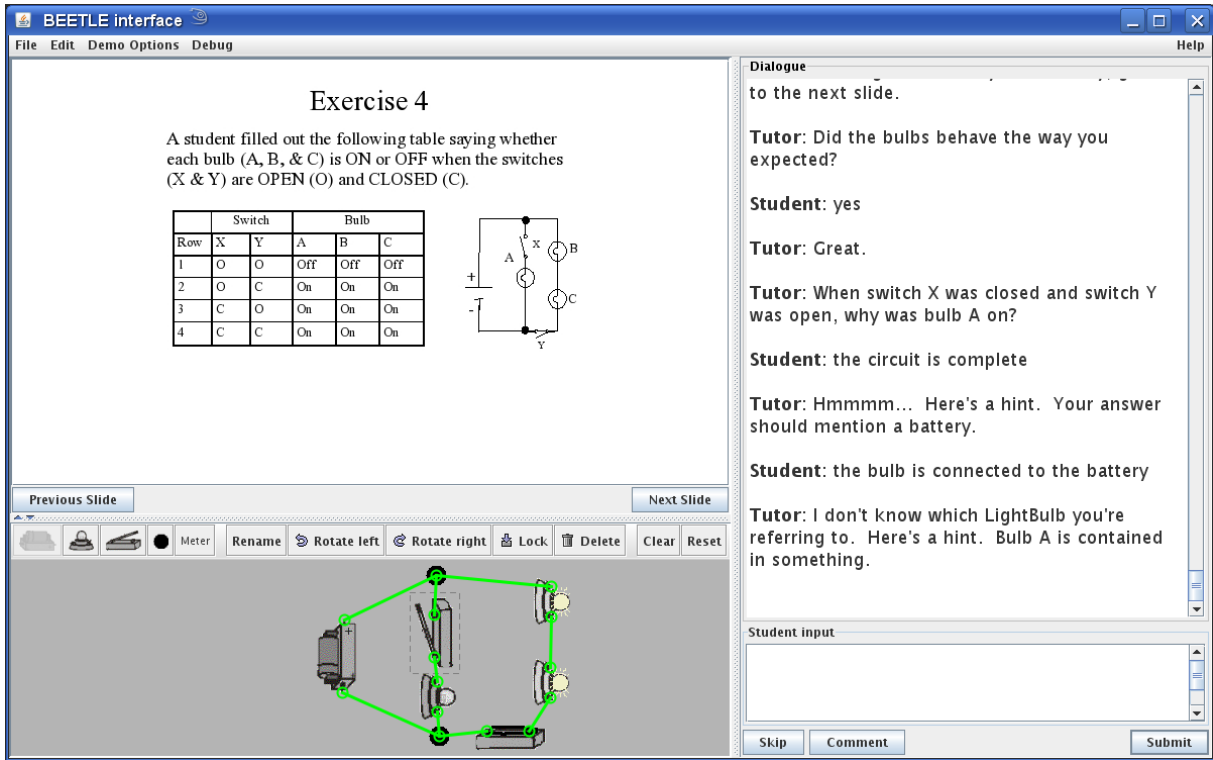


Figure 1: Screenshot of the BEETLE II system

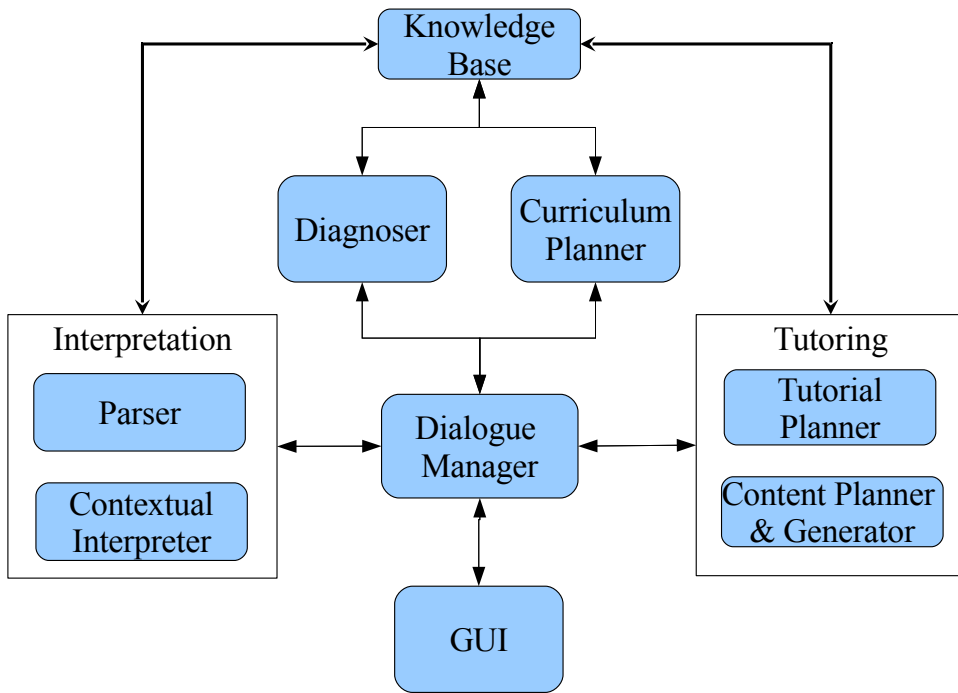


Figure 2: System architecture diagram

an as-yet unmentioned object and hints at it, for example, “Here’s a hint: Your answer should mention a battery.” For high-specificity, it attempts to hint at a two-place relation, for example, “Here’s a hint: the battery is connected to something.”

The tutorial policy makes a high-level decision as to which strategy to use (for example, “acknowledge the correct part and give a high specificity hint”) based on the answer analysis and dialogue context. At present, the system takes into consideration the number of incorrect answers received in response to the current question and the number of uninterpretable answers.¹

In addition to a remediation policy, the tutorial planner implements an error recovery policy (Dzikovska et al., 2009). Since the system accepts unrestricted input, interpretation errors are unavoidable. Our recovery policy is modeled on the TargetedHelp (Hockey et al., 2003) policy used in task-oriented dialogue. If the system cannot find an interpretation for an utterance, it attempts to produce a message that describes the problem but without giving away the answer, for example, “I’m sorry, I’m having a problem understanding. I don’t know the word *power*.” The help message is accompanied with a hint at the appropriate level, also depending on the number of previous incorrect and non-interpretable answers.

2.4 Generation

The strategy decision made by the tutorial planner, together with relevant semantic content from the student’s answer (e.g., part of the answer to confirm), is passed to content planning and generation. The system uses a domain-specific content planner to produce input to the surface realizer based on the strategy decision, and a FUF/SURGE (Elhadad and Robin, 1992) generation system to produce the appropriate text. Templates are used to generate some stock phrases such as “When you are ready, go on to the next slide.”

2.5 Dialogue Management

Interaction between components is coordinated by the dialogue manager which uses the information-state approach (Larsson and Traum, 2000). The dialogue state is represented by a cumulative answer analysis which tracks, over multiple turns, the correct, incorrect, and not-yet-mentioned parts

¹Other factors such as student confidence could be considered as well (Callaway et al., 2007).

of the answer. Once the complete answer has been accumulated, the system accepts it and moves on. Tutor hints can contribute parts of the answer to the cumulative state as well, allowing the system to jointly construct the solution with the student.

3 Evaluation

The first experimental evaluation involving 81 participants (undergraduates recruited from a South-eastern University in the USA) was completed in 2009. Participants had little or no prior knowledge of the domain. Each participant took a pre-test, worked through a lesson with the system, took a post-test, and completed a user satisfaction survey. Each session lasted approximately 4 hours.

We implemented two different tutoring policies in the system for this evaluation. The baseline policy used an “accept and bottom out” strategy for all student answers, regardless of their content. The students were always given the correct answer, but the system made no attempt at remediation, and never indicated whether the student was understood (or correct). In comparison, the full adaptive policy attempted to select an appropriate strategy based on student answer analysis and dialogue context, as described above. Students were randomly assigned to a condition. On average, each student session contained 232 student language turns in the fully adaptive condition ($SD = 25.6$) and 156 turns in the baseline condition ($SD = 2.02$). Additional time was taken by reading and interacting with the simulation environment.

Out of 81 students, 76 successfully completed the evaluation (data from 5 participants had to be discarded due to system crashes). There was no significant difference in learning gain between the two conditions (baseline: $M = 0.63$, $SD = 0.26$, full: $M = 0.57$, $SD = 0.23$, t -test $p = 0.32$), and students gave higher satisfaction scores to the baseline condition (baseline: $M = 3.32$, $SD = 0.65$, full: $M = 2.56$, $SD = 0.65$, t -test, $p < 0.01$). We are currently analyzing the data to investigate the reasons for this outcome. However, the evaluation demonstrated that the system can be successfully deployed in experiments with naive users and therefore can be used to study the impact of different natural language processing algorithms on learning gain and user satisfaction.

An example dialogue with the system, taken from the evaluation corpus, is shown in Figure 3.

Tutor: Why was bulb A on when switch Y was open and switch Z was closed?

Student: because it had a closed path

Tutor: Right. There is a closed path. Keep going.

Student: there was a closed path not effected by the switch

Tutor: I'm sorry, I'm having a problem understanding. I don't understand the word 'effect' here. Here's a hint. Your answer should mention a battery.

Student: it has a closed path with the battery

Tutor: You're on the right track. Here's the answer. Bulb A is still contained in a closed path with the battery.

Figure 3: Example interaction with the system from our corpus

It shows three key system properties: after the student's first turn, the system rephrases its understanding of the correct part of the student answer and prompts the student to supply the missing information. In the second turn, the student utterance could not be interpreted and the system responds with a targeted help message and a hint about the object that needs to be mentioned. Finally, in the last turn the system combines the information from the tutor's hint and the student's answers and restates the complete answer since the current answer was completed over multiple turns.

4 Conclusions and Future Work

The BEETLE II system we present was built to serve as a platform for research in computational linguistics and tutoring, and can be used for task-based evaluation of algorithms developed for other domains. We are currently developing an annotation scheme for the data we collected to identify student paraphrases of correct answers. The annotated data will be used to evaluate the accuracy of existing paraphrasing and textual entailment approaches and to investigate how to combine such algorithms with the current deep linguistic analysis to improve system robustness. We also plan to annotate the data we collected for evidence of misunderstandings, i.e., situations where the system arrived at an incorrect interpretation of a student utterance and took action on it. Such annotation can provide useful input for statistical learning algorithms to detect and recover from misun-

derstandings.

In dialogue management and generation, the key issue we are planning to investigate is that of linguistic alignment. The analysis of the data we have collected indicates that student satisfaction may be affected if the system rephrases student answers using different words (for example, using better terminology) but doesn't explicitly explain the reason why different terminology is needed (Dzikovska et al., 2010). Results from other systems show that measures of semantic coherence between a student and a system were positively associated with higher learning gain (Ward and Litman, 2006). Using a deep generator to automatically generate system feedback gives us a level of control over the output and will allow us to devise experiments to study those issues in more detail.

From the point of view of tutoring research, we are planning to use the system to answer questions about the effectiveness of different approaches to tutoring, and the differences between human-human and human-computer tutoring. Previous comparisons of human-human and human-computer dialogue were limited to systems that asked short-answer questions (Litman et al., 2006; Rosé and Torrey, 2005). Having a system that allows more unrestricted language input will provide a more balanced comparison. We are also planning experiments that will allow us to evaluate the effectiveness of individual strategies implemented in the system by comparing system versions using different tutoring policies.

Acknowledgments

This work has been supported in part by US Office of Naval Research grants N000140810043 and N0001410WX20278. We thank Katherine Harrison and Leanne Taylor for their help running the evaluation.

References

- V. Aleven, O. Popescu, and K. R. Koedinger. 2001. Towards tutorial dialog to support self-explanation: Adding natural language understanding to a cognitive tutor. In *Proceedings of the 10th International Conference on Artificial Intelligence in Education (AIED '01)*.
- James Allen, Myroslava Dzikovska, Mehdi Manshadi, and Mary Swift. 2007. Deep linguistic processing for spoken dialogue systems. In *Proceedings of the ACL-07 Workshop on Deep Linguistic Processing*.

- Mark Buckley and Magdalena Wolska. 2007. Towards modelling and using common ground in tutorial dialogue. In *Proceedings of DECALOG, the 2007 Workshop on the Semantics and Pragmatics of Dialogue*, pages 41–48.
- Donna K. Byron. 2002. *Resolving Pronominal Reference to Abstract Entities*. Ph.D. thesis, University of Rochester.
- Charles B. Callaway, Myroslava Dzikovska, Elaine Farrow, Manuel Marques-Pita, Colin Matheson, and Johanna D. Moore. 2007. The Beetle and BeeD-iff tutoring systems. In *Proceedings of SLaTE'07 (Speech and Language Technology in Education)*.
- Michelene T. H. Chi, Nicholas de Leeuw, Mei-Hung Chiu, and Christian LaVancher. 1994. Eliciting self-explanations improves understanding. *Cognitive Science*, 18(3):439–477.
- Peter Clark and Bruce Porter, 1999. *KM (1.4): Users Manual*. <http://www.cs.utexas.edu/users/mfkb/km>.
- Myroslava O. Dzikovska, Charles B. Callaway, and Elaine Farrow. 2006. Interpretation and generation in a knowledge-based tutorial system. In *Proceedings of EACL-06 workshop on knowledge and reasoning for language processing*, Trento, Italy, April.
- Myroslava O. Dzikovska, James F. Allen, and Mary D. Swift. 2008a. Linking semantic and knowledge representations in a multi-domain dialogue system. *Journal of Logic and Computation*, 18(3):405–430.
- Myroslava O. Dzikovska, Gwendolyn E. Campbell, Charles B. Callaway, Natalie B. Steinhauer, Elaine Farrow, Johanna D. Moore, Leslie A. Butler, and Colin Matheson. 2008b. Diagnosing natural language answers to support adaptive tutoring. In *Proceedings 21st International FLAIRS Conference*, Coconut Grove, Florida, May.
- Myroslava O. Dzikovska, Charles B. Callaway, Elaine Farrow, Johanna D. Moore, Natalie B. Steinhauer, and Gwendolyn C. Campbell. 2009. Dealing with interpretation errors in tutorial dialogue. In *Proceedings of SIGDIAL-09*, London, UK, Sep.
- Myroslava O. Dzikovska, Johanna D. Moore, Natalie Steinhauer, and Gwendolyn Campbell. 2010. The impact of interpretation problems on tutorial dialogue. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics (ACL-2010)*.
- Michael Elhadad and Jacques Robin. 1992. Controlling content realization with functional unification grammars. In R. Dale, E. Hovy, D. Rösner, and O. Stock, editors, *Proceedings of the Sixth International Workshop on Natural Language Generation*, pages 89–104, Berlin, April. Springer-Verlag.
- A. C. Graesser, P. Wiemer-Hastings, P. Wiemer-Hastings, and R. Kreuz. 1999. Autotutor: A simulation of a human tutor. *Cognitive Systems Research*, 1:35–51.
- Beth Ann Hockey, Oliver Lemon, Ellen Campana, Laura Hiatt, Gregory Aist, James Hieronymus, Alexander Gruenstein, and John Dowding. 2003. Targeted help for spoken dialogue systems: intelligent feedback improves naive users' performance. In *Proceedings of the tenth conference on European chapter of the Association for Computational Linguistics*, pages 147–154, Morristown, NJ, USA.
- Pamela Jordan, Maxim Makatchev, Umarani Pappuswamy, Kurt VanLehn, and Patricia Albacete. 2006. A natural language tutorial dialogue system for physics. In *Proceedings of the 19th International FLAIRS conference*.
- Staffan Larsson and David Traum. 2000. Information state and dialogue management in the TRINDI Dialogue Move Engine Toolkit. *Natural Language Engineering*, 6(3-4):323–340.
- Diane Litman, Carolyn P. Rosé, Kate Forbes-Riley, Kurt VanLehn, Dumisizwe Bhembe, and Scott Silliman. 2006. Spoken versus typed human and computer dialogue tutoring. *International Journal of Artificial Intelligence in Education*, 16:145–170.
- Diane Litman, Johanna Moore, Myroslava Dzikovska, and Elaine Farrow. 2009. Generalizing tutorial dialogue results. In *Proceedings of 14th International Conference on Artificial Intelligence in Education (AIED)*, Brighton, UK, July.
- Rodney D. Nielsen, Wayne Ward, and James H. Martin. 2008. Learning to assess low-level conceptual understanding. In *Proceedings 21st International FLAIRS Conference*, Coconut Grove, Florida, May.
- Amruta Purandare and Diane Litman. 2008. Content-learning correlations in spoken tutoring dialogs at word, turn and discourse levels. In *Proceedings 21st International FLAIRS Conference*, Coconut Grove, Florida, May.
- C.P. Rosé and C. Torrey. 2005. Interactivity versus expectation: Eliciting learning oriented behavior with tutorial dialogue systems. In *Proceedings of Interact'05*.
- N. B. Steinhauer, L. A. Butler, and G. E. Campbell. 2007. Simulated tutors in immersive learning environments: Empirically-derived design principles. In *Proceedings of the 2007 Interservice/Industry Training, Simulation and Education Conference*, Orlando, FL.
- Kurt VanLehn, Pamela Jordan, and Diane Litman. 2007. Developing pedagogically effective tutorial dialogue tactics: Experiments and a testbed. In *Proceedings of SLaTE Workshop on Speech and Language Technology in Education*, Farmington, PA, October.
- Arthur Ward and Diane Litman. 2006. Cohesion and learning in a tutorial spoken dialog system. In *Proceedings of 19th International FLAIRS (Florida Artificial Intelligence Research Society) Conference*, Melbourne Beach, FL.

GernEdiT: A Graphical Tool for GermaNet Development

Verena Henrich

University of Tübingen
Tübingen, Germany.

verena.henrich@uni-
tuebingen.de

Erhard Hinrichs

University of Tübingen
Tübingen, Germany.

erhard.hinrichs@uni-
tuebingen.de

Abstract

GernEdiT (short for: GermaNet Editing Tool) offers a graphical interface for the lexicographers and developers of GermaNet to access and modify the underlying GermaNet resource. GermaNet is a lexical-semantic wordnet that is modeled after the Princeton WordNet for English. The traditional lexicographic development of GermaNet was error prone and time-consuming, mainly due to a complex underlying data format and no opportunity of automatic consistency checks. GernEdiT replaces the earlier development by a more user-friendly tool, which facilitates automatic checking of internal consistency and correctness of the linguistic resource. This paper presents all these core functionalities of GernEdiT along with details about its usage and usability.

1 Introduction

The main purpose of the GermaNet Editing Tool GernEdiT tool is to support lexicographers in accessing, modifying, and extending the GermaNet data (Kunze and Lemnitzer, 2002; Henrich and Hinrichs, 2010) in an easy and adaptive way and to aid in the navigation through the GermaNet word class hierarchies, so as to find the appropriate place in the hierarchy for new synsets (short for: synonymy set) and lexical units. GernEdiT replaces the traditional GermaNet development based on lexicographer files (Fellbaum, 1998) by a more user-friendly visual tool that supports versioning and collaborative annotation by several lexicographers working in parallel.

Furthermore, GernEdiT facilitates internal consistency of the GermaNet data such as appropriate linking of lexical units with synsets, connectedness of the synset graph, and automatic

closure among relations and their inverse counterparts.

All these functionalities along with the main aspects of GernEdiT's usage and usability are presented in this paper.

2 The Structure of GermaNet

GermaNet is a lexical-semantic wordnet that is modeled after the Princeton WordNet for English (Fellbaum, 1998). It covers the three word categories of adjectives, nouns, and verbs and partitions the lexical space into a set of concepts that are interlinked by semantic relations. A semantic concept is modeled by a *synset*. A synset is a set of words (called *lexical units*) where all the words are taken to have (almost) the same meaning. Thus a synset is a set-representation of the semantic relation of synonymy, which means that it consists of a list of lexical units.

There are two types of semantic relations in GermaNet: *conceptual* and *lexical relations*. Conceptual relations hold between two semantic concepts, i.e. synsets. They include relations such as hyperonymy, part-whole relations, entailment, or causation. GermaNet is hierarchically structured in terms of the hyperonymy relation. Lexical relations hold between two individual lexical units. Antonymy, a pair of opposites, is an example of a lexical relation.

3 The GermaNet Editing Tool

The GermaNet Editing Tool GernEdiT provides a graphical user interface, implemented as a Java Swing application, which primarily allows maintaining the GermaNet data in a user-friendly way. The editor represents an interface to a relational database, where all GermaNet data is stored from now on.

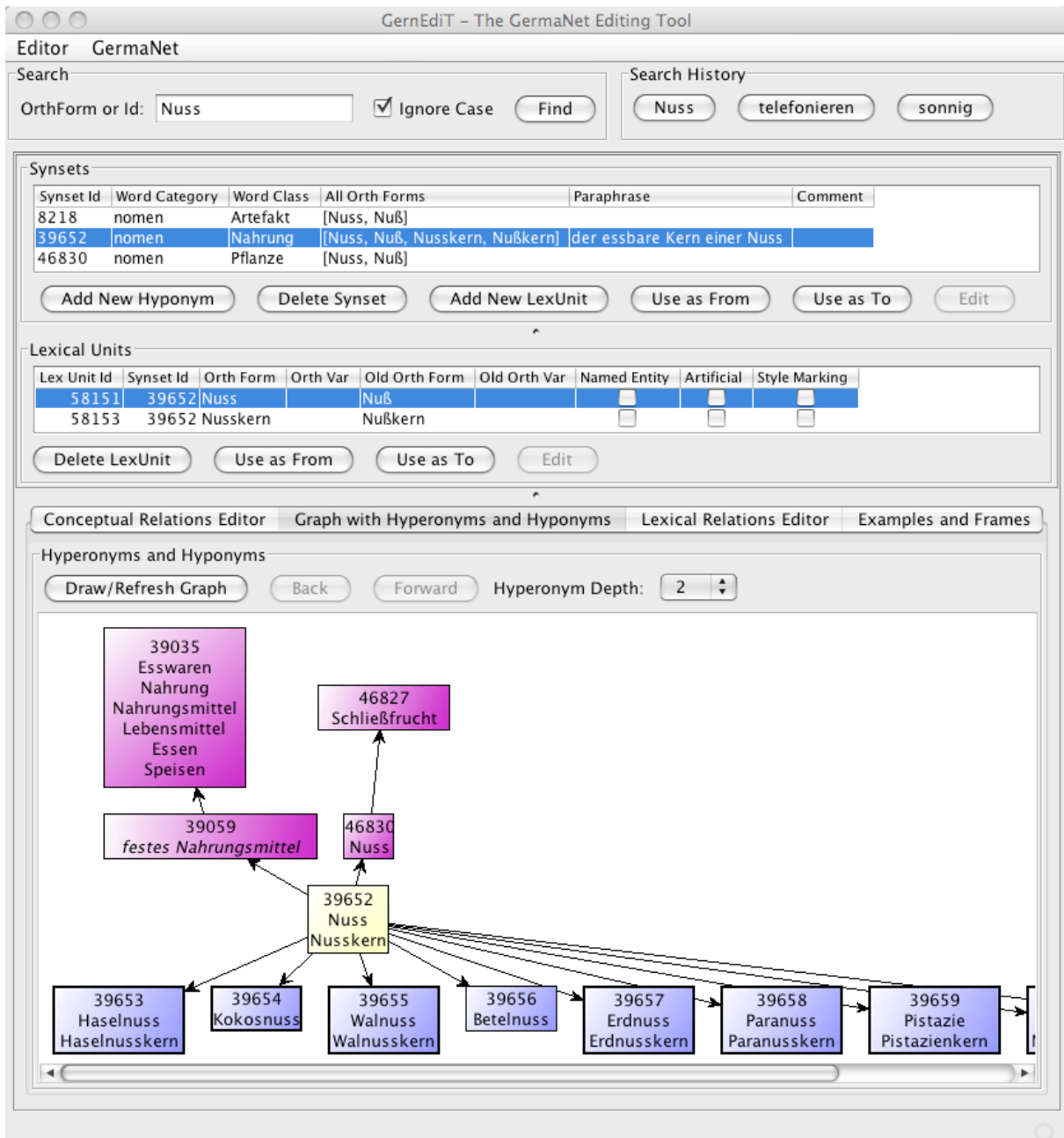


Figure 1. The main view of GernEdiT.

3.1 Motivation

The traditional lexicographic development of GermaNet was error prone and time-consuming, mainly due to a complex underlying data format and no opportunity of automatic consistency checks. This is exactly why GernEdiT was developed: It supports lexicographers who need to access, modify, and extend GermaNet data by providing these functions through simple button-clicks, searches, and form editing. There are several ways to search data and browse through the GermaNet graph. These functionalities allow lexicographers, among other things, to find the

appropriate place in the hierarchy for the insertion of new synsets and lexical units. Last but not least, GernEdiT facilitates internal consistency and correctness of the linguistic resource and supports versioning and collaborative annotation of GermaNet by several lexicographers working in parallel.

3.2 The Main User Interface

Figure 1 illustrates the main user panel of GernEdiT. It shows a *Search* panel above, two panels for *Synsets* and *Lexical Units* in the middle, and four tabs below: a *Conceptual Relations Editor*, a *Graph with Hyperonyms and Hyponyms*, a *Lexi-*

Lex Unit Id	Orth Form	Synset Id	Synonyms	Frames	Examples
84970	ausweiten	60447	[ausweiten]	NN.AN,NN.AN	"Wir werden unsere Bemühungen auf weitere Zielgruppen auswe..."
76139	abhalten	53955	[abhalten]	NN.AN.PP	"Wir versuchen, die Jugendlichen vom Rauchen abzuhalten."
84392	abspalten	60021	[abspalten]	NN.AN	"Wir müssen diese Eiweißmoleküle abspalten."
75212	aufwarten	53276	[aufwarten]	NN.DN	"Wichtige Politiker warteten dem König auf."
82681	aufleuchten	58722	[aufscheinen, aufleuchten]	NN	"Wenn das Warnlicht aufleuchtet, sollten Sie die Maschine sofort..."
74744	ausrasten	52956	[ausrasten, durchdrehen,...]	NN	"Wegen der Nachricht rastete er aus."
76593	abtöten	54272	[deitalisieren, abtöten]	NN.AN	"Vor der Wurzelbehandlung musste die Zahnärztin den Nerv abt..."
84660	ausrotten	60223	[ausrotten]	NN.AN	"Viele Tierarten sind durch Jagd und Industrie schon ausgerottet..."
74652	auseinanderdriften	52885	[auseinanderdriften]	NN	"Unsere Standpunkte drifteten zusehends auseinander."
76264	auswerten	54043	[auswerten]	NN.AN	"Sie werteten die neueste Marktanalyse aus."
75667	auftreten	53610	[auftreten, agieren]	NN.BM,NN.BR,NN.BL...	"Sie tritt selbstbewusst und entschieden auf.", "Sie tritt selbstbew..."
75876	auftreten	53772	[auftreten]	NN.PP.BM	"Sie tritt in der Öffentlichkeit sicher auf."
75415	abtreten	53424	[abtreten]	NN.PP	"Sie trat als Schauspielerin von der Bühne ab."
75515	abstreiten	53495	[ablehnen, abstreiten]	NN.AZ,NN.DS,NN.AN	"Sie stritt seine Anschuldigungen ab.", "Sie stritt seine Anschuldigu..."
85540	aufsplitten	60882	[splitten, aufsplitten]	NN.AN.PP,NN.AN,N...	"Sie splitteten das Vermögen in vier gleich große Teile auf.", "Sie..."
81894	antreten	58152	[antreten]	NN.BU	"Sie sind pünktlich zum Dienst angetreten."
79224	ausschelten	56189	[ausschimpfen, auszanke...	NN.AN	"Sie schalt ihren Sohn aus, nachdem er heimlich abgehauen war."
84227	ausschalten	59890	[abschalten, ausschalten,...]	NN.AN,NN.AN	"Sie schaltete den Staubsauger ab.", "Er schaltete das Radio aus."
84221	anschalten	59887	[anstellen, anschalten, ein...	NN.AN	"Sie schaltete das Licht an."
78914	abraten	55945	[abraten]	NN.AZ,NN.DN,NN.DN.P...	"Sie riet ihm ab, nochmals dort anzurufen.", "Sie riet ihm ab, noch..."
77956	anarbeiten	55249	[anarbeiten]	NN.PP	"Sie muss gegen die neuen Richtlinien anarbeiten."
77111	abschalten	54617	[ausspannen, ausruhen, a...	NN	"Sie kann gut abschalten."
78039	auseinanderhalten	55311	[differenzieren, auseinan...	NN.AN	"Sie kann die Zwillingsschwestern nicht auseinanderhalten."
80008	aufheften	56784	[aufheften]	NN.AN.Pp	"Sie heftete das Tuch auf einer einfachen Plastikoberfläche auf."
75029	aufarbeiten	53141	[aufarbeiten]	NN.AN	"Sie hat viel Stoff aufzuarbeiten."
79575	anraten	56454	[anempfehlen, anraten]	NN.AN.DN	"Sie hatte ihm mehr Einfühlungsvermögen angeraten."
83768	abhärten	59546	[abhärten]	NN.AN.Pp,NN.AR.Pp...	"Sie hat ihren Körper durch gesunde Ernährung abgehärtet.", "Sie..."
78140	abtöten	55388	[abtöten, abwürgen]	NN.AN	"Sie hat ihre Gefühle für ihn abgetötet."

Figure 2: Filtered list of lexical units.

cal Relations Editor, and an Examples and Frames tab.

In Figure 1, a search for synsets consisting of lexical units with the word *Nuss* (German noun for: *nut*) has been executed. Accordingly, the *Synsets* panel displays the three resulting synsets that match the search item. The *Synset Id* is the unique database ID that unambiguously identifies a synset, and which can also be used to search for exactly that synset. *Word Category* specifies whether a synset is an adjective (*adj*), a noun (*nomen*), or a verb (*verben*), whereas *Word Class* classifies the synsets into semantic fields. The word class of the selected synset in Figure 1 is *Nahrung* (German noun for: *food*). The *Paraphrase* column contains a description of a synset, e.g., for the selected synset the paraphrase is: *der essbare Kern einer Nuss* (German phrase for: *the edible kernel of a nut*). The column *All Orth Forms* simply lists all orthographical variants of all its lexical units.

Which lexical units are listed in the *Lexical Units* panel depends on the selected synset in the *Synsets* panel. Here, *Lex Unit Id* and *Synset Id* again reflect the corresponding unique database IDs. *Orth Form* (short for: *orthographic form*) represents the correct spelling of a word according to the rules of the spelling reform *Neue Deutsche Rechtschreibung* (Rat für deutsche Rechtschreibung, 2006), a recently adopted spelling reform. In our example, the main orthographic form is *Nuss*. *Orth Var* may contain an

alternative spelling that is admissible according to the *Neue Deutsche Rechtschreibung*.¹ *Old Orth Form* represents the main orthographic form prior to the *Neue Deutsche Rechtschreibung*. This means that *Nuß* was the correct spelling instead of *Nuss* before the German spelling reform. *Old Orth Var* contains any accepted variant prior to the *Neue Deutsche Rechtschreibung*. The *Old Orth Var* field is filled only if it is no longer allowed in the new orthography.

The Boolean values *Named Entity*, *Artificial*, and *Style Marking* express further properties of a lexical unit, whether the lexical unit is a named entity, an artificial concept node, or a stylistic variant.

For both the lexical units and the synsets, there are two buttons *Use as From* and *Use as To*, which help to add new relations (see the explanation of Figure 3 in section 3.6 below which explains the creation of new relations).

3.3 Search Functionalities

It is possible to search for words or synset database IDs via the search panel (see Figure 1 at the top). The check box *Ignore Case* offers the possibility of searching without distinguishing between upper and lower case.

¹ An example of this kind is the German word *Delfin* (German noun for: *dolphin*). Apart from the main form *Delfin*, there is an orthographic variant *Delphin*.

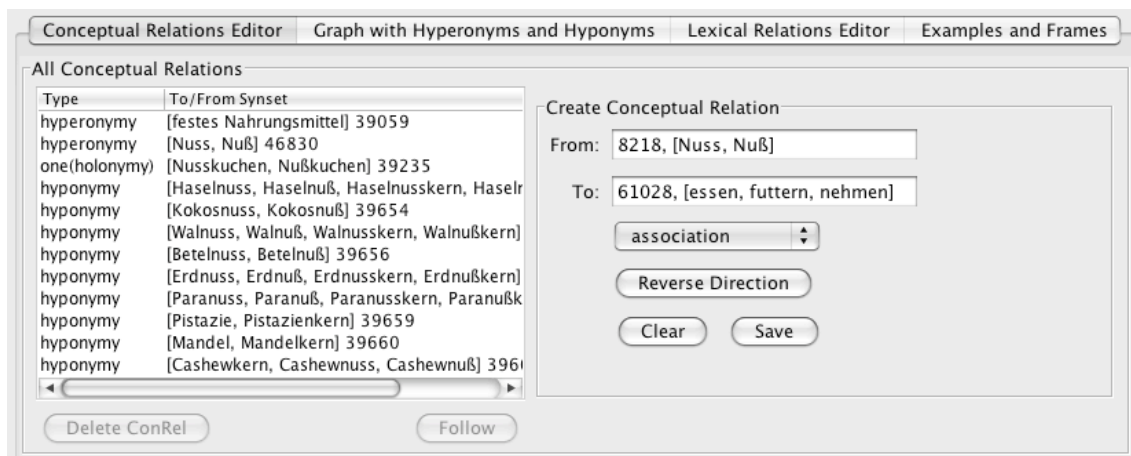


Figure 3. Conceptual Relations Editor tab.

Via the file menu, lists of all synsets or lexical units with their properties can be accessed. To these lists, very detailed filters can be applied: e.g., filtering the lexical units or synsets by parts of their orthographical forms. Figure 2 shows a list of lexical units to which a detailed filter has been applied: verbs have been chosen (see the chosen tab) whose orthographical forms start with an *a-* (see *starts with* check box and corresponding text field) and end with the suffix *-ten* (see *ends with* check box and corresponding text field). Only verbs that have a frame that contains *NN* are chosen (see *Frame contains* check box and corresponding text field). Furthermore, the resulting filtered list is sorted in descending order by their examples (see the little triangle in the *Examples* header of the result table). The number in the brackets behind the word category in the tab title indicates the count of the filtered lexical units (in this example 193 verbs pass the filter).

3.4 Visualization of the Graph Hierarchy

There is the possibility to display a graph with all hyperonyms and hyponyms of a selected synset. This is shown in the bottom half of Figure 1 in the tab *Graph with Hyperonyms and Hyponyms*. The graph in Figure 1 visualizes a part of the hierarchical structure of GermaNet centered around the synset containing *Nuss* and displays the hyperonyms and hyponyms of this synset up to a certain parameterized depth (in this case depth 2 has been chosen). The *Hyperonym Depth* chooser allows unfolding the graph to the top up to the preselected depth. As it is not possible to visualize the whole GermaNet contents at once, the graph can be seen as a window to GermaNet.

A click on any synset node within the graph, navigates to that synset. This functionality supports lexicographers especially in finding the appropriate place in the hierarchy for the insertion of new synsets.

3.5 Modifications of Existing Items

If the lexicographers' task is to modify existing synsets or lexical units, this is done by selecting a synset or lexical unit displayed in the *Synsets* and the *Lexical Units* panels shown in Figure 1. The properties of such selected items can be edited by a click in the corresponding table cell. For example by clicking in the cell *Orth Form* the spelling of a lexical unit can be corrected in case of an earlier typo was made.

If lexicographers want to edit examples, frames, conceptual, or lexical relations this is done by choosing the appropriate tab indicated at the bottom of Figure 1. By clicking one of these tabs, the corresponding panel appears below these tabs. In Figure 1 the panel for *Graph with Hyperonyms and Hyponyms* is displayed.

It is possible to edit the examples and frames associated with a lexical unit via the *Examples and Frames* tab. Frames specify the syntactic valence of a lexical unit. Each frame can have an associated example that indicates a possible usage of the lexical unit for that particular frame. The tab *Examples and Frames* is thus particularly geared towards the editing of verb entries. By clicking on the tab all examples and frames of a lexical unit are listed and can then be modified by choosing the appropriate editing buttons. For more information about these editing functions see Henrich and Hinrichs (2010).

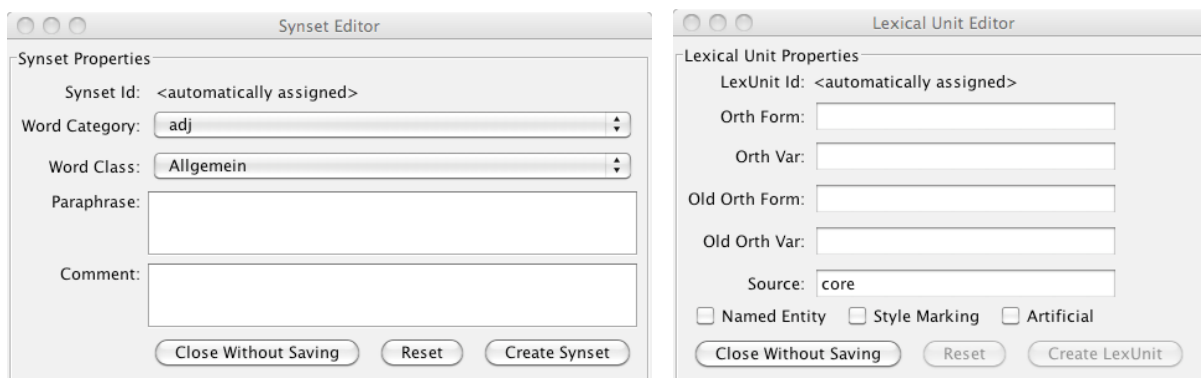


Figure 4. Synset Editor (left). Lexical Units Editor (right).

3.6 Editing of Relations

If lexicographers want to add new conceptual or lexical relations to a synset or a lexical unit this is done by clicking on the *Conceptual Relations Editor* or the *Lexical Relations Editor* shown in Figure 1.

Figure 3 shows the panel that appears if the *Conceptual Relations Editor* has been chosen for the synset containing *Nuss*. To create a new relation, the lexicographer needs to use the buttons *Use as From* and *Use as To* shown in Figure 1. This will insert the ID of the selected synsets from the *Synsets* panel in the corresponding *From* or *To* field in Figure 3. The button *Delete ConRel* allows deletion of a conceptual relation, if all consistency checks are passed.

The *Lexical Relations Editor* tab supports editing all lexical relations. It is not displayed separately for reasons of space, but it is analogue to the *Conceptual Relations Editor* tab for editing conceptual relations.

3.7 Adding Synsets and Lexical Units

The buttons *Add New Hyponym* and *Add New LexUnit* in the *Synsets* panel (see Figure 1) can be used to insert a new synset or lexical unit at the selected place in the GermaNet graph, and the buttons *Delete Synset* and *Delete LexUnit* remove the selected entry, respectively.

The *Synset Editor* in Figure 4 (on the left) shows the window which appears after a click on *Add New Hyponym*. When clicking on the button *Create Synset*, the *Lexical Unit Editor* (shown in Figure 4, right) pops up. This workflow forces the parallel creation of a lexical unit while creating a synset.

3.8 Consistency Checks

GernEiT facilitates internal consistency of the GermaNet data. This is achieved by the

workflow-oriented design of the editor. It is not possible to create a synset without creating a lexical unit in parallel (as described in section 3.7). Furthermore, it is not possible to insert a new synset without specifying the place in the GermaNet hierarchy where the new synset should be added. This is achieved by the button *Add New Hyponym* (see Figure 1) which forces the user to identify the appropriate hyperonym for the new synset to be added. Furthermore, it is not possible to insert a lexical unit without specifying the corresponding synset. On deletion of a synset, all corresponding data such as conceptual relations, lexical units with their lexical relations, frames, and examples, are deleted automatically.

Consistency checks also take effect for the table cell editing in the *Synsets* and *Lexical Units* panels of the main user interface (see Figure 1), e.g., the main orthographic form of a lexical unit may never be empty.

All buttons in GernEiT are enabled only if the corresponding functionalities meet the consistency requirements, e.g., if a synset consists only of one lexical unit, it is not possible to delete that lexical unit and thus the button *Delete LexUnit* is disabled. Also, if the deletion of a synset or a relation would violate the complete connectedness of the GermaNet graph, it is not possible to delete that synset.

3.9 Further Functionalities

There are further functionalities available through the file menu. Besides retrieving the up-to-date statistics of GermaNet, an editing history makes it possible to list all modifications on the GermaNet data, with the information about who made the change and how the modified item looked before.

GernEiT supports various export functionalities. For example, it is possible to export all GermaNet contents into XML files, which are used as an exchange format of GermaNet, or to

export a list of all verbs with their corresponding frames and examples.

4 Tool Evaluation

In order to assess the usefulness of GernEdiT, we conducted in depth interviews with the GermaNet lexicographers and with the senior researcher who oversees all lexicographic development. At the time of the interview all of these researchers had worked with the tool for about eight months. The present section summarizes the feedback about GernEdiT that was obtained in this way.

The initial learning curve for getting familiar with GernEdiT is considerably lower compared to the learning curve required for the traditional development based on lexicographer files. Moreover, the GermaNet development with GernEdiT is both more efficient and accurate compared to the traditional development along the following dimensions:

1. The menu-driven and graphics-based navigation through the GermaNet graph is much easier compared to finding the correct entry point in the purely text-based format of lexicographer files.
2. Lexicographers no longer need to learn the complex specification syntax of the lexicographer files. Thereby, syntax errors in the specification language – a frequent source of errors prior to development with GernEdiT – are entirely eliminated.
3. GernEdiT facilitates automatic checking of internal consistency and correctness of the GermaNet data such as appropriate linking of lexical units with synsets, connectedness of the synset graph, and automatic closure among relations and their inverse counterparts.
4. It is now even possible to perform further queries, which were not possible before, e.g., listing all hyponyms of a synset.
5. Especially for the senior researcher who is responsible for coordinating the GermaNet lexicographers, it is now much easier to trace back changes and to verify who was responsible for them.
6. The collaborative annotation by several lexicographers working in parallel is now easily possible and does not cause any management overhead as before.

In sum, the lexicographers of GermaNet gave very positive feedback about the use of GernEdiT and also made smaller suggestions for improving its user-friendliness further. This underscores the utility of GernEdiT from a practical point of view.

5 Conclusion and Future Work

In this paper we have described the functionality of GernEdiT. The extremely positive feedback of the GermaNet lexicographers underscores the practical benefits gained by using the GernEdiT tool in practice.

At the moment, GernEdiT is customized for maintaining the GermaNet data. In future work, we plan to adapt the tool so that it can be used with wordnets for other languages as well. This would mean that the wordnet data for a given language would have to be stored in a relational database and that the tool itself can handle the language specific data structures of the wordnet in question.

Acknowledgements

We would like to thank all GermaNet lexicographers for their willingness to experiment with GernEdiT and to be interviewed about their experiences with the tool.

Special thanks go to Reinhild Barkey for her valuable input on both the features and user-friendliness of GernEdiT and to Alexander Kislev for his contributions to the underlying database format.

References

- Claudia Kunze and Lothar Lemnitzer. 2002. *GermaNet – representation, visualization, application*. Proceedings of LREC 2002, Main Conference, Vol V. pp. 1485-1491, 2002.
- Christiane Fellbaum (ed.). 1998. *WordNet – An Electronic Lexical Database*. Cambridge, MA: MIT Press.
- Verena Henrich and Erhard Hinrichs. 2010. *GernEdiT – The GermaNet Editing Tool*. Proceedings of LREC 2010, Main Conference, Valletta, Malta.
- Rat für deutsche Rechtschreibung (eds.) (2006). *Deutsche Rechtschreibung – Regeln und Wörterverzeichnis: Amtliche Regelung*. Gunter Narr Verlag Tübingen.

WebLicht: Web-based LRT services for German

Erhard Hinrichs, Marie Hinrichs, Thomas Zastrow
Seminar für Sprachwissenschaft, University of Tübingen
firstname.lastname@uni-tuebingen.de

Abstract

This software demonstration presents WebLicht (short for: *Web-Based Linguistic Chaining Tool*), a web-based service environment for the integration and use of language resources and tools (LRT). WebLicht is being developed as part of the D-SPIN project¹. WebLicht is implemented as a web application so that there is no need for users to install any software on their own computers or to concern themselves with the technical details involved in building tool chains. The integrated web services are part of a prototypical infrastructure that was developed to facilitate chaining of LRT services. WebLicht allows the integration and use of distributed web services with standardized APIs. The nature of these open and standardized APIs makes it possible to access the web services from nearly any programming language, shell script or workflow engine (UIMA, Gate etc.) Additionally, an application for integration of additional services is available, allowing anyone to contribute his own web service.

1 Introduction

Currently, WebLicht offers LRT services that were developed independently at the Institut für Informatik, Abteilung Automatische Sprachverarbeitung at the University of Leipzig (tokenizer, lemmatizer, co-occurrence extraction, and frequency analyzer), at the Institut für Maschinelle Sprachverarbeitung at the University of Stuttgart (tokenizer, tagger/lemmatizer, German morphological analyser SMOR, constituent and dependency parsers), at the Berlin Brandenburgische Akademie der Wissenschaften (conversion of plain text to D-Spin format, tokenizer, taggers, NE recog-

¹ D-SPIN stands for **D**eutsche **S**prachressourcen **I**nfrastruktur; the D-SPIN project is partly financed by the BMBF; it is a national German complement to the EU-project CLARIN. See the URLs <http://www.d-spin.org> and <http://www.clarin.eu> for details

nizer) and at the Seminar für Sprachwissenschaft/Computerlinguistik at the University of Tübingen (conversion of plain text to D-Spin format, GermaNet, Open Thesaurus synonym service, and Treebank browser). They cover a wide range of linguistic applications, like tokenization, co-occurrence extraction, POS Tagging, lexical and semantic analysis, and several languages (currently German, English, Italian, French, Romanian, Spanish and Finnish). For some of these tasks, more than one web service is available. As a first external partner, the University of Helsinki in Finland contributed a set of web services to create morphological annotated text corpora in the Finnish language. With the help of the webbased user interface, these individual web services can be combined into a chain of linguistic applications.

2 Service Oriented Architecture

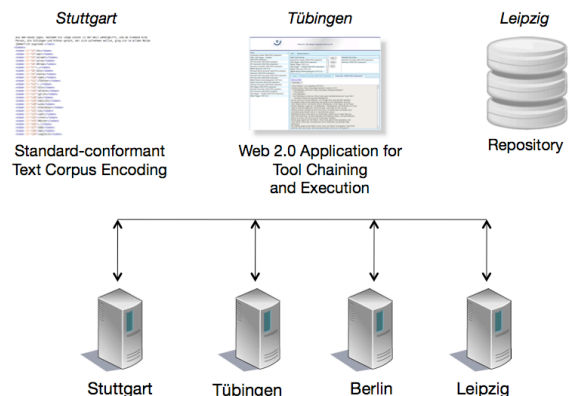


Figure 1: The Overall Structure of WebLicht

WebLicht is a so-called Service Oriented Architecture (Binildas et al., 2008), which means that distributed and independent services (Tanenbaum et al, 2002) are combined together to a chain of LRT tools. A centralized database, the repository, stores technical and content-related metadata about each service. With the help of

this repository, the chaining mechanism as described in section 3 is implemented. The WebLicht user interface encapsulates this chaining mechanism in an AJAX driven web application. Since web applications can be invoked from any browser, downloading and installation of individual tools on the user's local computer is avoided. But using WebLicht web services is not restricted to the use of the integrated user interface. It is also possible to access the web services from nearly any programming language, shell script or workflow engine (UIMA, Gate etc.). Figure 1 depicts the overall structure of WebLicht.

An important part of Service Oriented Architectures is ensuring interoperability between the underlying services. Interoperability of web services, as they are implemented in WebLicht, refers to the seamless flow of data between them. To be interoperable, these web services must first agree on protocols defining the interaction between the services (WSDL/SOAP, REST, XML-RPC). They must also use a shared and standardized data exchange format, which is preferably based on widely accepted formats already in use (UTF-8, XML). WebLicht uses the RESTstyle API and its own XML-based data exchange format (Text Corpus Format, TCF).

3 The Service Repository

Every tool included in WebLicht is registered in a central repository, located in Leipzig. Also realized as a web service, it offers metadata and processing information about each registered tool. For example, the metadata includes information about the creator, name and the adress of the service. The input and output specifications of each web service are required in order to determine which processing chains are possible. Combining the metadata and the processing information, the repository is able to offer functions for the chain building process.

Wrappers: TCF, 0.3 / TCF, 0.3	
Inputs	Outputs
lemmas postags -tagset: stts	sem_lex_rels -source: GermaNet

Table 1: Input and Output Specifications of Tübingen's Semantic Annotator

A specialized tool for registering new web services in the repository is available.

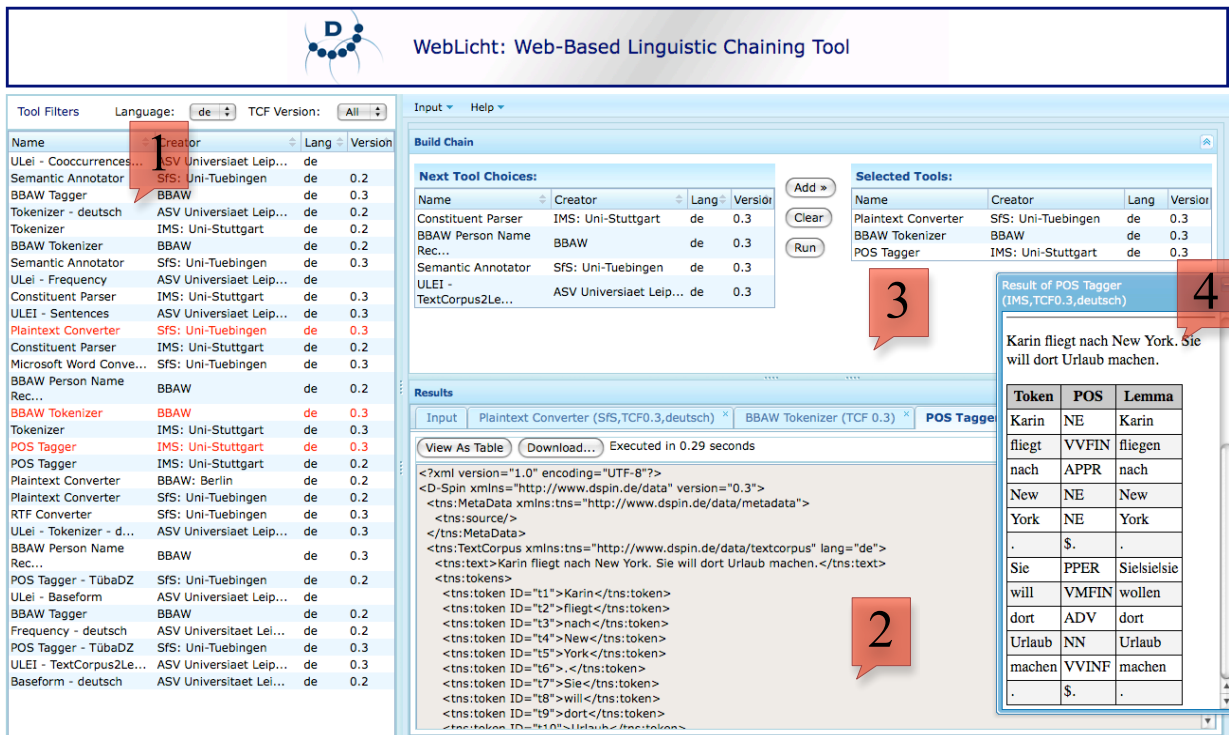


Figure 2: A Screenshot of the WebLicht Webinterface

4 The WebLicht User Interface

Figure 2 shows a screenshot of the WebLicht web interface, developed and hosted in Tübingen. Area 1 shows a list of all WebLicht web services along with a subset of metadata (author, URL, description etc.). This list is extracted on-the-fly from a centralized repository located in Leipzig. This means that after registration in the repository, a web service is immediately available for inclusion in a processing chain.

The *Language Filter* selection box allows the selection of any language for which tools are available in WebLicht (currently, German, English, Italian, French, Romanian, Spanish or Finnish). The majority of the presently integrated web services operates on German input. The platform, however, is language-independent and supports LRT resources for any language.

Plain text input to the service chain can be specified in one of three ways: a) entered by the user in the *Input tab*, b) file upload from the user's local harddrive or c) selecting one of the sample texts offered by WebLicht (Area 2). Various format converters can be used to convert uploaded files into the data exchange format (TCF) used by WebLicht. Input file formats accepted by WebLicht currently include plain text, Microsoft Word, RTF and PDF.

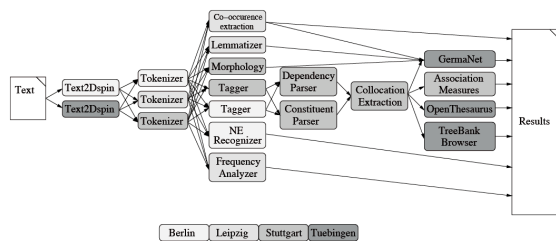


Figure 3: A Choice of Alternative Services

In Area 3, one can assemble the service tool chain and execute it on the input text. The *Selected Tools* list displays all web services that have already been entered into the web service chain. The list under *Next Tool Choices* then offers the set of tools that can be entered as next into the chain. This list is generated by inspecting the metadata of the tools which are already in the chain. The chaining mechanism ensures that this list only contains tools, that are a valid next step in the chain. For example, a Part-of-Speech

Tagger can only be added to a chain after a tokenizer has been added. The metadata of each tool contains information about the annotations which are required in the input data and which annotations are added by that tool.

As Figure 3 shows, the user sometimes has a choice of alternative tools - in the example at hand a wide variety of services are offered as candidates. Figure 3 shows a subset of web service workflows currently available in WebLicht. Notice that these workflows can combine tools from various institutions and are not restricted to predefined combinations of tools. This allows users to compare the results of several tool chains and find the best solution for their individual use case.

The final result of running the tool chain as well as each individual step can be visualized in a *Table View* (implemented as a separate frame, Area 4), or downloaded to the user's local harddrive in WebLicht's own data exchange format TCF.

5 The TCF Format

The D-SPIN *Text Corpus Format* TCF (Heid et al, 2010) is used by WebLicht as an internal data

```
<?xml version="1.0" encoding="UTF-8" ?>
<D-Spin xmlns="http://www.dspin.de/data" version="0.3">
  <tns:MetaData xmlns:tns="http://www.dspin.de/data/metadata">
    <tns:source/>
  </tns:MetaData>
  <tns:TextCorpus xmlns:tns="http://www.dspin.de/data/textcorpus" lang="en">
    <tns:text>Bob went to the zoo.</tns:text>
    <tns:tokens>
      <tns:token ID="t0">Bob</tns:token>
      <tns:token ID="t1">went</tns:token>
      <tns:token ID="t2">to</tns:token>
      <tns:token ID="t3">the</tns:token>
      <tns:token ID="t4">zoo</tns:token>
      <tns:token ID="t5">.</tns:token>
    </tns:tokens>
    <tns:POSTags tagset="PennTB">
      <tns:tag tokID="t0">NP</tns:tag>
      <tns:tag tokID="t1">VBD</tns:tag>
      <tns:tag tokID="t2">TO</tns:tag>
      <tns:tag tokID="t3">DT</tns:tag>
      <tns:tag tokID="t4">NN</tns:tag>
      <tns:tag tokID="t5">.</tns:tag>
    </tns:POSTags>
    <tns:lemmas>
      <tns:lemma tokID="t0">Bob</tns:lemma>
      <tns:lemma tokID="t1">go</tns:lemma>
      <tns:lemma tokID="t2">to</tns:lemma>
      <tns:lemma tokID="t3">the</tns:lemma>
      <tns:lemma tokID="t4">zoo</tns:lemma>
      <tns:lemma tokID="t5">.</tns:lemma>
    </tns:lemmas>
  </tns:TextCorpus>
</D-Spin>
```

Figure 4: A Short Example of a TCF Document, Containing the Plain Text, Tokens and POS Tags and Lemmas

exchange format. The TCF format allows the combination of the different linguistic annotations produced by the tool chain. It supports incremental enrichment of linguistic annotations at different levels of analysis in a common XML-based format (see Figure 4).

The Text Corpus Format was designed to efficiently enable the seamless flow of data between the individual services of a Service Oriented Architecture.

Figure 4 shows a data sample in the D-SPIN Text Corpus Format. Lexical tokens are identified via token IDs which serve as unique identifiers in different annotation layers. From an organizational point-of-view, tokens can be seen as the central, atomic elements in TCF to which other annotation layers refer. For example, the POS annotations refer to the token IDs in the token annotation layer via the attribute *tokID*. The annotation layers are rendered in a stand-off annotation format. TCF stores all linguistic annotation layers in one single file. That means that during the chaining process, the file grows (see Figure 5). Each tool is permitted to add an arbitrary number of layers, but it is not allowed to change or delete any existing layer.

Within the D-SPIN project, several other XML based data formats were developed beside the TCF format (for example, an encoding for lexicon based data). In order to avoid any confusion of element names between these different formats, namespaces for the different contextual scopes within each format have been introduced. At the end of the chaining process, converter services will convert the textcorpora from the

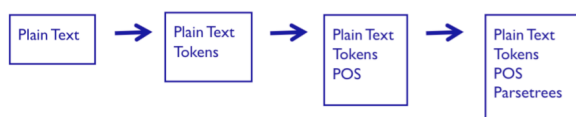


Figure 5: Annotation Layers are Added to the TCF Document by Each Service

TCF format into other common and standardized data formats, for example MAF/SynAF or TEI.

6 Implementation Details

The web services are available in RESTstyle and use the TCF data format for input and output. The concrete implementation can use any combination of programming language and server environment.

The repository is a relational database, offering its content also as RESTstyle web services.

The user interface is a Rich Internet Application (RIA), using an AJAX driven toolkit. It incorporates the Java EE 5 technology and can be deployed in any Java application server.

7 How to Participate in WebLicht

Since WebLicht follows the paradigm of a Service Oriented Architecture, it is easily extendable by adding new services. In order to participate in WebLicht by donating additional tools, one must implement the tool as a RESTful web service using the TCF data format. You can find further information including a tutorial on the D-SPIN homepage².

8 Further Work

The WebLicht platform in its current form moves the functionality of LRT tools from the users desktop computer into the net (Gray et al, 2005). At this point, the user must download the results of the chaining process and deal with them on his local machine again. In the future, an online workspace has to be implemented so that annotated textcorpora created with WebLicht can also be stored in and retrieved from the net. For that purpose, an integration of the eSciDoc research environment³ into Weblicht is planned. The eSciDoc infrastructure enables sustainable and reliable long-term preservation of primary research and analysis data.

To make the use of WebLicht more convenient to the end user, there will be predefined processing chains. These will consist of the most commonly used processing chains and will relieve the user of having to define the chains manually. In the last year, WebLicht has proven to be a realizable and useful service environment for the humanities. In its current state, WebLicht is still a prototype: due to the restrictions of the underlying hardware, WebLicht cannot yet be made available to the general public.

9 Scope of the Software Demonstration

This demonstration will present the core functionalities of WebLicht as well as related modules and applications. The process of building language-specific processing tool chains will be shown. WebLicht's capability of offering only appropriate tools at each step in the chain-building process will be demonstrated.

² <http://weblicht.sfs.uni-tuebingen.de/englisch/weblichttutorial.shtml>

³ For further information about the eSciDoc platform, see <https://www.escidoc.org/>

The selected tool chain can be applied to any arbitrary uploaded text. The resulting annotated text corpus can be downloaded or visualized using an integrated software module.

All these functions will be shown live using just a webbrowser during the software demonstration. Demo Preview and Hardware Requirements

The call for papers asks submitters of software demonstrations to provide pointers to demo previews and to provide technical details about hardware requirements for the actual demo at the conference.

The WebLicht web application is currently password protected. Access can be granted by requesting an account (weblicht@d-spin.org).

If the software demonstration is accepted, internet access is necessary at the conference, but no special hardware is required. The authors will bring a laptop of their own and if necessary also a beamer.

Acknowledgments

WebLicht is the product of a combined effort within the D-SPIN projects (www.d-spin.org). Currently, partners include: Seminar für Sprachwissenschaft/Computerlinguistik, Universität Tübingen, Abteilung für Automatische Sprachverarbeitung, Universität Leipzig, Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart and Berlin Brandenburgische Akademie der Wissenschaften.

References

- Binildas, C.A., Malhar Barai et.al. (2008). *Service Oriented Architectures with Java*. PACKT Publishing, Birmingham – Mumbai
- Gray, J., Liu, D., Nieto-Santisteban, M., Szalay, A., DeWitt, D., Heber, G. (2005). Scientific Data Management in the Coming Decade. Technical Report MSR-TR-2005-10, Microsoft Research.
- Heid, U., Schmid, H., Eckart, K., Hinrichs, E. (2010). A Corpus Representation Format for Linguistic Web Services: the D_SPIN Text Corpus Format and its Relationship with ISO Standards. In Proceedings of LREC 2010, Malta.
- Tanenbaum, A., van Steen, M. (2002). *Distributed Systems*, Prentice Hall, Upper Saddle River, NJ, 1st Edition.

The S-Space Package: An Open Source Package for Word Space Models

David Jurgens

University of California, Los Angeles,
4732 Boelter Hall
Los Angeles, CA 90095
jurgens@cs.ucla.edu

Keith Stevens

University of California, Los Angeles,
4732 Boelter Hall
Los Angeles, CA 90095
kstevens@cs.ucla.edu

Abstract

We present the S-Space Package, an open source framework for developing and evaluating word space algorithms. The package implements well-known word space algorithms, such as LSA, and provides a comprehensive set of matrix utilities and data structures for extending new or existing models. The package also includes word space benchmarks for evaluation. Both algorithms and libraries are designed for high concurrency and scalability. We demonstrate the efficiency of the reference implementations and also provide their results on six benchmarks.

1 Introduction

Word similarity is an essential part of understanding natural language. Similarity enables meaningful comparisons, entailments, and is a bridge to building and extending rich ontologies for evaluating word semantics. Word space algorithms have been proposed as an automated approach for developing meaningfully comparable semantic representations based on word distributions in text.

Many of the well known algorithms, such as Latent Semantic Analysis (Landauer and Dumais, 1997) and Hyperspace Analogue to Language (Burgess and Lund, 1997), have been shown to approximate human judgements of word similarity in addition to providing computational models for other psychological and linguistic phenomena. More recent approaches have extended this approach to model phenomena such as child language acquisition (Baroni et al., 2007) or semantic priming (Jones et al., 2006). In addition, these models have provided insight in fields outside of linguistics, such as information retrieval, natural language processing and cognitive psychology. For a recent survey of word space approaches and applications, see (Turney and Pantel, 2010).

The parallel development of word space models in different fields has often resulted in duplicated work. The pace of development presents a need for a reliable method for accurate comparisons between new and existing approaches. Furthermore, given the frequent similarity of approaches, we argue that the research community would greatly benefit from a common library and evaluation utilities for word spaces. Therefore, we introduce the **S-Space Package**, an open source framework with four main contributions:

1. reference implementations of frequently cited algorithms
2. a comprehensive, highly concurrent library of tools for building new models
3. an evaluation framework for testing models on standard benchmarks, e.g. the TOEFL Synonym Test (Landauer et al., 1998)
4. a standardized interface for interacting with all word space models, which facilitates word space based applications.

The package is written in Java and defines a standardized Java interface for word space algorithms. While other word space frameworks exist, e.g. (Widdows and Ferraro, 2008), the focus of this framework is to ease the development of new algorithms and the comparison against existing models. Compared to existing frameworks, the S-Space Package supports a much wider variety of algorithms and provides significantly more reusable developer utilities for word spaces, such as tokenizing and filtering, sparse vectors and matrices, specialized data structures, and seamless integration with external programs for dimensionality reduction and clustering. We hope that the release of this framework will greatly facilitate other researchers in their efforts to develop and validate new word space models. The toolkit is available at <http://code.google.com/p/airhead-research/>, which includes a wiki

containing detailed information on the algorithms, code documentation and mailing list archives.

2 Word Space Models

Word space models are based on the contextual distribution in which a word occurs. This approach has a long history in linguistics, starting with Firth (1957) and Harris (1968), the latter of whom defined this approach as the Distributional Hypothesis: for two words, their similarity in meaning is predicted by the similarity of the distributions of their co-occurring words. Later models have expanded the notion of co-occurrence but retain the premise that distributional similarity can be used to extract meaningful relationships between words.

Word space algorithms consist of the same core algorithmic steps: word features are extracted from a corpus and the distribution of these features is used as a basis for semantic similarity. Figure 1 illustrates the shared algorithmic structure of all the approaches, which is divided into four components: corpus processing, context selection, feature extraction and global vector space operations.

Corpus processing normalizes the input to create a more uniform set of features on which the algorithm can work. Corpus processing techniques frequently include stemming and filtering of stop words or low-frequency words. For web-gathered corpora, these steps also include removal of non linguistic tokens, such as html markup, or restricting documents to a single language.

Context selection determines which tokens in a document may be considered for features. Common approaches use a lexical distance, syntactic relation, or document co-occurrence to define the context. The various decisions for selecting the context accounts for many differences between otherwise similar approaches.

Feature extraction determines the dimensions of the vector space by selecting which tokens in the context will count as features. Features are commonly word co-occurrences, but more advanced models may perform a statistical analysis to select only those features that best distinguish word meanings. Other models approximate the full set of features to enable better scalability.

Global vector space operations are applied to the entire space once the initial word features have been computed. Common operations include altering feature weights and dimensionality reduc-

Document-Based Models
LSA (Landauer and Dumais, 1997)
ESA (Gabrilovich and Markovitch, 2007)
Vector Space Model (Salton et al., 1975)
Co-occurrence Models
HAL (Burgess and Lund, 1997)
COALS (Rohde et al., 2009)
Approximation Models
Random Indexing (Sahlgren et al., 2008)
Reflective Random Indexing (Cohen et al., 2009)
TRI (Jurgens and Stevens, 2009)
BEAGLE (Jones et al., 2006)
Incremental Semantic Analysis (Baroni et al., 2007)
Word Sense Induction Models
Purandare and Pedersen (Purandare and Pedersen, 2004)
HERMIT (Jurgens and Stevens, 2010)

Table 1: Algorithms in the S-Space Package

tion. These operations are designed to improve word similarity by changing the feature space itself.

3 The S-Space Framework

The S-Space framework is designed to be extensible, simple to use, and scalable. We achieve these goals through the use of Java interfaces, reusable word space related data structures, and support for multi-threading. Each word space algorithm is designed to run as a stand alone program and also to be used as a library class.

3.1 Reference Algorithms

The package provides reference implementations for twelve word space algorithms, which are listed in Table 1. Each algorithm is implemented in its own Java package, and all commonalities have been factored out into reusable library classes. The algorithms implement the same Java interface, which provides a consistent abstraction of the four processing stages.

We divide the algorithms into four categories based on their structural similarity: document-based, co-occurrence, approximation, and Word Sense Induction (WSI) models. Document-based models divide a corpus into discrete documents and construct the vector space from word frequencies in the documents. The documents are defined independently of the words that appear in them. Co-occurrence models build the vector space using the distribution of co-occurring words in a context, which is typically defined as a region around a word or paths rooted in a parse tree. The third category of models approximate

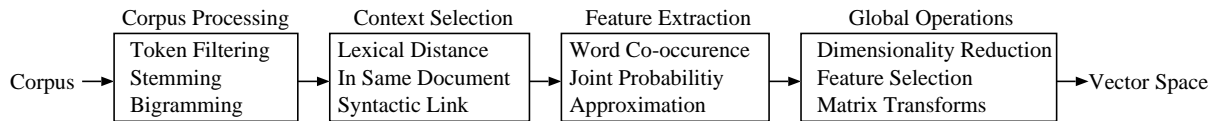


Figure 1: A high-level depiction of common algorithmic steps that convert a corpus into a word space

co-occurrence data rather than model it explicitly in order to achieve better scalability for larger data sets. WSI models also use co-occurrence but also attempt to discover distinct word senses while building the vector space. For example, these algorithms might represent “earth” with two vectors based on its meanings “planet” and “dirt.”

3.2 Data Structures and Utilities

The S-Space Package provides efficient implementations for matrices, vectors, and specialized data structures such as multi-maps and tries. Implementations are modeled after the `java.util` library and offer concurrent implementations when multi-threading is required. In addition, the libraries provide support for converting between multiple matrix formats, enabling interaction with external matrix-based programs. The package also provides support for parsing different corpora formats, such as XML or email threads.

3.3 Global Operation Utilities

Many algorithms incorporate dimensionality reduction to smooth their feature data, e.g. (Landauer and Dumais, 1997; Rohde et al., 2009), or to improve efficiency, e.g. (Sahlgren et al., 2008; Jones et al., 2006). The S-Space Package supports two common techniques: the Singular Value Decomposition (SVD) and randomized projections. All matrix data structures are designed to seamlessly integrate with six SVD implementations for maximum portability, including `SVDLIBJ`¹, a Java port of `SVDLIBC`², a scalable sparse SVD library. The package also provides a comprehensive library for randomized projections, which project high-dimensional feature data into a lower dimensional space. The library supports both integer-based projections (Kanerva et al., 2000) and Gaussian-based (Jones et al., 2006).

The package supports common matrix transformations that have been applied to word spaces: point wise mutual information (Dekang,

1998), term frequency-inverse document frequency (Salton and Buckley, 1988), and log entropy (Landauer and Dumais, 1997).

3.4 Measurements

The choice of similarity function for the vector space is the least standardized across approaches. Typically the function is empirically chosen based on a performance benchmark and different functions have been shown to provide application specific benefits (Weeds et al., 2004). To facilitate exploration of the similarity function parameter space, the S-Space Package provides support for multiple similarity functions: cosine similarity, Euclidean distance, KL divergence, Jaccard Index, Pearson product-moment correlation, Spearman’s rank correlation, and Lin Similarity (Dekang, 1998)

3.5 Clustering

Clustering serves as a tool for building and refining word spaces. WSI algorithms, e.g. (Purandare and Pedersen, 2004), use clustering to discover the different meanings of a word in a corpus. The S-Space Package provides bindings for using the CLUTO clustering package³. In addition, the package provides Java implementations of Hierarchical Agglomerative Clustering, Spectral Clustering (Kannan et al., 2004), and the Gap Statistic (Tibshirani et al., 2000).

4 Benchmarks

Word space benchmarks assess the semantic content of the space through analyzing the geometric properties of the space itself. Currently used benchmarks assess the semantics by inspecting the representational similarity of word pairs. Two types of benchmarks are commonly used: word choice tests and association tests. The S-Space Package supports six tests, and has an easily extensible model for adding new tests.

¹<http://bender.unibe.ch/svn/codemap/Archive/svdlbj/>

²<http://tedlab.mit.edu/~dr/SVDLIBC/>

³<http://glaros.dtc.umn.edu/gkhome/views/cluto>

Algorithm	Corpus	Word Choice			Word Association		
		TOEFL	ESL	RDWP	R-G	WordSim353	Deese
BEAGLE	TASA	46.03	35.56	46.99	0.431	0.342	0.235
COALS	TASA	65.33	60.42	93.02	0.572	0.478	0.388
HAL	TASA	44.00	20.83	50.00	0.173	0.180	0.318
HAL	Wiki	50.00	31.11	43.44	0.261	0.195	0.042
ISA	TASA	41.33	18.75	33.72	0.245	0.150	0.286
LSA	TASA	56.00 ^a	50.00	45.83	0.652	0.519	0.349
LSA	Wiki	60.76	54.17	59.20	0.681	0.614	0.206
P&P	TASA	34.67	20.83	31.39	0.088	-0.036	0.216
RI	TASA	42.67	27.08	34.88	0.224	0.201	0.211
RI	Wiki	68.35	31.25	40.80	0.226	0.315	0.090
RI + Perm. ^b	TASA	52.00	33.33	31.39	0.137	0.260	0.268
RRI	TASA	36.00	22.92	34.88	0.088	0.138	0.109
VSM	TASA	61.33	52.08	84.88	0.496	0.396	0.200

^a Landauer et al. (1997) report a score of 64.4 for this test, while Rohde et al. (2009) report a score of 53.4.

^b + Perm indicates that permutations were used with Random Indexing, as described in (Sahlgren et al., 2008)

Table 2: A comparison of the implemented algorithms on common evaluation benchmarks

4.1 Word Choice

Word choice tests provide a target word and a list of options, one of which has the desired relation to the target. Word space models solve these tests by selecting the option whose representation is most similar. Three word choice benchmarks that measure synonymy are supported.

The first test is the widely-reported Test of English as a Foreign Language (TOEFL) synonym test from (Landauer et al., 1998), which consists of 80 multiple-choice questions with four options. The second test comes from the English as a Second Language (ESL) exam and consists of 50 question with four choices (Turney, 2001). The third consists of 200 questions from the Canadian Reader’s Digest Word Power (RDWP) (Jarmasz and Szpakowicz, 2003), which unlike the previous two tests, allows the target and options to be multi-word phrases.

4.2 Word Association

Word association tests measure the semantic relatedness of two words by comparing word space similarity with human judgements. Frequently, these tests measure synonymy; however, other types of word relations such as antonymy (“hot” and “cold”) or functional relatedness (“doctor” and “hospital”) are also possible. The S-Space Package supports three association tests.

The first test uses data gathered by Rubenstein

and Goodneough (1965). To measure word similarity, word similarity scores of 51 human reviewers were gathered a set of 65 noun pairs, scored on a scale of 0 to 4. The ratings are then correlated with word space similarity scores.

Finkelstein et al. (2002) test for relatedness. 353 word pairs were rated by either 13 or 16 subjects on a 0 to 10 scale for how related the words are. This test is notably more challenging for word space models because human ratings are not tied to a specific semantic relation.

The third benchmark considers the antonym association. Deese (1964) introduced 39 antonym pairs that Greffenstette (1992) used to assess whether a word space modeled the antonymy relationship. We quantify this relationship by measuring the similarity rank of each word in an antonym pair, w_1, w_2 , i.e. w_2 is the k^{th} most-similar word to w_1 in the vector space. The antonym score is calculated as $\frac{2}{rank_{w_1}(w_2) + rank_{w_2}(w_1)}$. The score ranges from $[0, 1]$, where 1 indicates that the most similar neighbors in the space are antonyms. We report the mean score for all 39 antonyms.

5 Algorithm Analysis

The content of a word space is fundamentally dependent upon the corpus used to construct it. Moreover, algorithms which use operations such as the SVD have a limit to the corpora sizes they

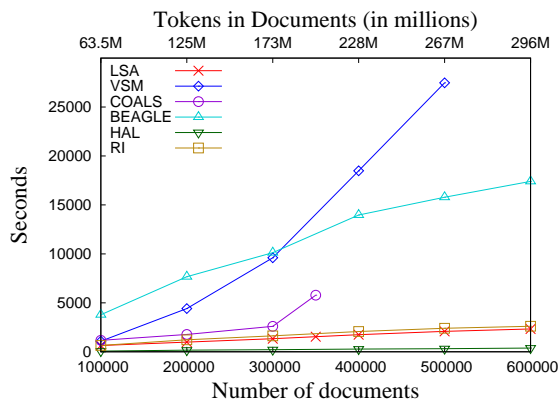


Figure 2: Processing time across different corpus sizes for a word space with the 100,000 most frequent words

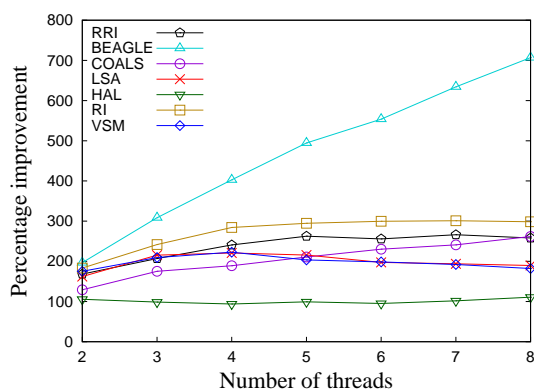


Figure 3: Run time improvement as a factor of increasing the number of threads

can process. We therefore highlight the differences in performance using two corpora. TASA is a collection of 44,486 topical essays introduced in (Landauer and Dumais, 1997). The second corpus is built from a Nov. 11, 2009 Wikipedia snapshot, and filtered to contain only articles with more than 1000 words. The resulting corpus consists of 387,082 documents and 917 million tokens.

Table 2 reports the scores of reference algorithms on the six benchmarks using cosine similarity. The variation in scoring illustrates that different algorithms are more effective at capturing certain semantic relations. We note that scores are likely to change for different parameter configurations of the same algorithm, e.g. token filtering or changing the number of dimensions.

As a second analysis, we report the efficiency of reference implementations by varying the corpus size and number of threads. Figure 2 reports the total amount of time each algorithm needs for processing increasingly larger segments of a web-

gathered corpus when using 8 threads. In all cases, only the top 100,000 words were counted as features. Figure 3 reports run time improvements due to multi-threading on the TASA corpus.

Algorithm efficiency is determined by three factors: contention on global statistics, contention on disk I/O, and memory limitations. Multi-threading benefits increase proportionally to the amount of work done per context. Memory limitations account for the largest efficiency constraint, especially as the corpus size and number of features grow. Several algorithms lack data points for larger corpora and show a sharp increase in running time in Figure 2, reflecting the point at which the models no longer fit into 8GB of memory.

6 Future Work and Conclusion

We have described a framework for developing and evaluating word space algorithms. Many well known algorithms are already provided as part of the framework as reference implementations for researchers in distributional semantics. We have shown that the provided algorithms and libraries scale appropriately. Last, we motivate further research by illustrating the significant performance differences of the algorithms on six benchmarks.

Future work will be focused on providing support for syntactic features, including dependency parsing as described by (Padó and Lapata, 2007), reference implementations of algorithms that use this information, non-linear dimensionality reduction techniques, and more advanced clustering algorithms.

References

- Marco Baroni, Alessandro Lenci, and Luca Onnis. 2007. Isa meets lara: A fully incremental word space model for cognitively plausible simulations of semantic learning. In *Proceedings of the 45th Meeting of the Association for Computational Linguistics*.
- Curt Burgess and Kevin Lund. 1997. Modeling parsing constraints with high-dimensional context space. *Language and Cognitive Processes*, 12:177210.
- Trevor Cohen, Roger Schvaneveldt, and Dominic Widows. 2009. Reflective random indexing and indirect inference: A scalable method for discovery of implicit connections. *Journal of Biomedical Informatics*, 43.
- J. Deese. 1964. The associative structure of some common english adjectives. *Journal of Verbal Learning and Verbal Behavior*, 3(5):347–357.

- Lin Dekang. 1998. Automatic retrieval and clustering of similar words. In *Proceedings of the Joint Annual Meeting of the Association for Computational Linguistics and International Conference on Computational Linguistics*, pages 768–774.
- L. Finkelstein, E. Gabrilovich, Y. Matias, E. Z. S. Rivlin, G. Wolfman, and E. Ruppin. 2002. Placing search in context: The concept revisited. *ACM Transactions of Information Systems*, 20(1):116–131.
- J. R. Firth, 1957. *A synopsis of linguistic theory 1930-1955*. Oxford: Philological Society. Reprinted in F. R. Palmer (Ed.), (1968). *Selected papers of J. R. Firth 1952-1959*, London: Longman.
- Evgeniy Gabrilovich and Shaul Markovitch. 2007. Computing semantic relatedness using wikipedia-based explicit semantic analysis. In *IJCAI'07: Proceedings of the 20th international joint conference on Artificial intelligence*, pages 1606–1611.
- Gregory Grefenstette. 1992. Finding semantic similarity in raw text: The Deese antonyms. In *Working notes of the AAAI Fall Symposium on Probabilistic Approaches to Natural Language*, pages 61–65. AAAI Press.
- Zellig Harris. 1968. *Mathematical Structures of Language*. Wiley, New York.
- Mario Jarmasz and Stan Szpakowicz. 2003. Roget's thesaurus and semantic similarity. In *Conference on Recent Advances in Natural Language Processing*, pages 212–219.
- Michael N. Jones, Walter Kintsch, and Douglas J. K. Mewhort. 2006. High-dimensional semantic space accounts of priming. *Journal of Memory and Language*, 55:534–552.
- David Jurgens and Keith Stevens. 2009. Event detection in blogs using temporal random indexing. In *Proceedings of RANLP 2009: Events in Emerging Text Types Workshop*.
- David Jurgens and Keith Stevens. 2010. HERMIT: Flexible Clustering for the SemEval-2 WSI Task. In *Proceedings of the 5th International Workshop on Semantic Evaluations (SemEval-2010)*. Association of Computational Linguistics.
- P. Kanerva, J. Kristoferson, and A. Holst. 2000. Random indexing of text samples for latent semantic analysis. In L. R. Gleitman and A. K. Josh, editors, *Proceedings of the 22nd Annual Conference of the Cognitive Science Society*, page 1036.
- Ravi Kannan, Santosh Vempala, and Adrian Vetta. 2004. On clusterings: Good, bad and spectral. *Journal of the ACM*, 51(3):497–515.
- Thomas K. Landauer and Susan T. Dumais. 1997. A solution to Plato's problem: The Latent Semantic Analysis theory of the acquisition, induction, and representation of knowledge. *Psychological Review*, 104:211–240.
- T. K. Landauer, P. W. Foltz, and D. Laham. 1998. Introduction to Latent Semantic Analysis. *Discourse Processes*, (25):259–284.
- Sebastian Padó and Mirella Lapata. 2007. Dependency-Based Construction of Semantic Space Models. *Computational Linguistics*, 33(2):161–199.
- Amruta Purandare and Ted Pedersen. 2004. Word sense discrimination by clustering contexts in vector and similarity spaces. In *HLT-NAACL 2004 Workshop: Eighth Conference on Computational Natural Language Learning (CoNLL-2004)*, pages 41–48. Association for Computational Linguistics.
- Douglas L. T. Rohde, Laura M. Gonnerman, and David C. Plaut. 2009. An improved model of semantic similarity based on lexical co-occurrence. *Cognitive Science*. submitted.
- H. Rubenstein and J. B. Goodenough. 1965. Contextual correlates of synonymy. *Communications of the ACM*, 8:627–633.
- M. Sahlgren, A. Holst, and P. Kanerva. 2008. Permutations as a means to encode order in word space. In *Proceedings of the 30th Annual Meeting of the Cognitive Science Society (CogSci'08)*.
- G. Salton and C. Buckley. 1988. Term-weighting approaches in automatic text retrieval. *Information Processing & Management*, 24:513–523.
- G. Salton, A. Wong, and C. S. Yang. 1975. A vector space model for automatic indexing. *Communications of the ACM*, 18(11):613–620.
- Robert Tibshirani, Guenther Walther, and Trevor Hastie. 2000. Estimating the number of clusters in a dataset via the gap statistic. *Journal Royal Statistics Society B*, 63:411–423.
- Peter D. Turney and Patrick Pantel. 2010. From Frequency to Meaning: Vector Space Models of Semantics. *Journal of Artificial Intelligence Research*, 37:141–188.
- Peter D. Turney. 2001. Mining the Web for synonyms: PMI-IR versus LSA on TOEFL. In *Proceedings of the Twelfth European Conference on Machine Learning (ECML-2001)*, pages 491–502.
- Julie Weeds, David Weir, and Diana McCarty. 2004. Characterising measures of lexical distributional similarity. In *Proceedings of the 20th International Conference on Computational Linguistics COLING'04*, pages 1015–1021.
- Dominic Widdows and Kathleen Ferraro. 2008. Semantic vectors: a scalable open source package and online technology management application. In *Proceedings of the Sixth International Language Resources and Evaluation (LREC'08)*.

Talking NPCs in a Virtual Game World

Tina Klüwer, Peter Adolphs, Feiyu Xu, Hans Uszkoreit, Xiwen Cheng

Deutsches Forschungszentrum für Künstliche Intelligenz (DFKI)

Projektbüro Berlin

Alt-Moabit 91c

10559 Berlin

Germany

{tina.kluewer,peter.adolphs,feiyu,uszkoreit,xiwen.cheng}@dfki.de

Abstract

This paper describes the *KomParse* system, a natural-language dialog system in the three-dimensional virtual world *Twinity*. In order to fulfill the various communication demands between non-player characters (NPCs) and users in such an online virtual world, the system realizes a flexible and hybrid approach combining knowledge-intensive domain-specific question answering, task-specific and domain-specific dialog with robust chatbot-like chitchat.

1 Introduction

In recent years multi-user online games in virtual worlds such as *Second Life* or *World of Warcraft* have attracted large user communities. Such virtual online game worlds provide new social and economic platforms for people to work and interact in. Furthermore, virtual worlds open new perspectives for research in the social, behavioral, and economic sciences, as well as in human-centered computer science (Bainbridge, 2007). Depending on the game type, non-player characters (NPCs) are often essential for supporting the game plot, for making the artificial world more vivid and ultimately for making it more immersive. In addition, NPCs are useful to populate new worlds by carrying out jobs the user-led characters come in touch with. The range of functions to be filled by NPCs is currently still strongly restricted by their limited capabilities in autonomous acting and communication. This shortcoming creates a strong need for progress in areas such as AI and NLP, especially their planning and dialog systems.

The *KomParse* system, described in this paper, provides NPCs for a virtual online world named *Twinity*, a product of the Berlin startup company

*Metaversum*¹. The *KomParse* NPCs offer various services through conversation with game users using question-answering and dialog functionality. The utilization of Semantic Web technology with RDF-encoded generic and domain-specific ontologies furthermore enables semantic search and inference.

This paper is organized as follows: Section 2 presents the NPC modelling and explains the application scenarios. Section 3 details the knowledge representation and semantic inference in our system. Section 4 explains the system architecture and its key components. Section 5 describes the *KomParse* dialog system. Section 7 gives a conclusion and closes off with our future work.

2 Application Scenario and NPC Modelling

The online game *Twinity* extends the *Second Life* idea by mirroring an urban part of the real world. At the time of this writing, the simulated section of reality already contains 3D models of the cities of Berlin, Singapore and London and it keeps growing. Users can log into the virtual world, where they can meet other users and communicate with them using the integrated chat function or talk to each other via Voice-over-IP. They can style their virtual appearance, can rent or buy their own flats and decorate them as to their preferences and tastes.

Out of many types of NPCs useful for this application such as pedestrians, city guides and personnel in stores, restaurants and bars, we start with two specific characters: a female furniture sales agent and a male bartender. The furniture seller is designed for helping users furnish their virtual apartments. Users can buy pieces of furniture and room decoration from the NPC by describing their demands and wishes in a text chat. During the di-

¹<http://www.metaversum.com/>



Figure 1: The furniture sales NPC selling a sofa

along with the NPC, the preferred objects are then selected and directly put into a location in the apartment, which can be further refined with the user interfaces that *Twinity* provides.

In the second scenario, the bartender sells virtual drinks. He can talk about cocktails with users, but moreover, he can also entertain his guests by providing trivia-type information about popular celebrities and various relations among them.

We chose these two characters not only because of their value for the *Twinity* application but also for our research goals. They differ in many interesting aspects. First of all, the furniture sales agent is controlled by a complex task model including ontology-driven and data-driven components to guide the conversation. This agent also possesses a much more fine-grained action model, which allows several different actions to cover the potential conversation situations for the selling task. The bartender agent on the other hand is designed not to fulfill one strict task because his clients do not follow a specific goal except ordering drinks. Our bartender has the role of a conversation companion and is able to entertain clients with his broad knowledge. Thus, he is allowed to access to several knowledge bases and is able to handle questions (and later conversations) about a much larger domain called the “gossip domain” which enables conversation about pop stars, movie actors and other celebrities as well as the relations between these people. In order to achieve a high robustness, we integrate a chatbot into the bartender agent to catch chitchat utterances we cannot handle.



Figure 2: Our bartender NPC in his bar in *Twinity*

3 Knowledge Representation and Semantic Inference

Semantic Web technology is employed for modelling the knowledge of the NPCs. The Resource Description Format (RDF) serves as the base for the actual encoding. An RDF statement is a binary relation instance between two individuals, that is a triple of a predicate and two arguments, called the subject and the object, and written as *subj pred obj* (e.g. *f:Sofa_Alatea f:hasMainColour f:Burgundy*).

All objects and properties the NPC can talk about are modelled in this way. Therefore the knowledge base has to reflect the physical properties of the virtual objects in *Twinity* as faithfully as possible. For instance, specific pieces of furniture are described by their main color, material or style, whereas cocktails are characterized by their ingredients, color, consistence and taste. Furthermore, references to the 3D models of the objects are stored in order to create, find and remove such objects in the virtual world.

The concepts and individuals of the particular domain are structured and organized in domain-specific ontologies. These ontologies are modelled in the Web Ontology Language (OWL). OWL allows us to define concept hierarchies, relations between concepts, domains and ranges of these relations, as well as specific relation instances between instances of a concept. Our ontologies are defined by the freely available ontology editor Protégé 4.0². The advantage of using an ontology for structuring the domain knowledge is

²<http://protege.stanford.edu/>, as accessed 27 Oct 2009

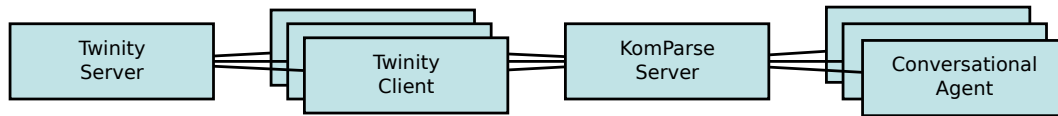


Figure 3: Overall System Architecture – Server/Client Architecture for NPC Control

the modular non-redundant encoding. When combined with a reasoner, only a few statements about an individual have to be asserted explicitly, while the rest can be inferred from the ontology. We employ several ontologies, among which the following are relevant for modelling the specific domains of our NPCs:

- An extensive furniture ontology, created by our project partner ZAS Berlin, defining kinds of furniture, room parts, colors and styles as well as the specific instances of furniture in *Twinity*. This knowledge base contains 95,258 triples, 123 furniture classes, 20 color classes, 243 color instances and various classes defining styles and similar concepts.
- A cocktail ontology, defining 13 cocktail classes with ingredients and tastes in 21,880 triples.
- A biographical ontology, the “gossip ontology”, defining biographical and career-specific concepts for people. This ontology is accompanied by a huge database of celebrities, which has been automatically acquired from the Web and covers nearly 600,000 persons and relations between these people like family relationships, marriages and professional relations. (Adolphs et al., 2010)

The furniture ontology is the only knowledge base for the furniture sales agent, whereas the bartender NPC has access to both the cocktail as well as the gossip knowledge base.

We use *SwiftOwlim*³ for storing and querying the data. *SwiftOwlim* is a “triple store”, a kind of database which is specifically built for storing and querying RDF data. It provides a forward-chaining inference engine which evaluates the domain definitions when loading the knowledge repository, and makes implicit knowledge explicit by asserting triples that must also hold true according to the ontology. Once the reasoner is finished, the triple store can be queried directly using the RDF query language SPARQL.

³<http://www.ontotext.com/owlim/>

4 Overall System Architecture

Figure 3 shows the overall system architecture. *Twinity* is a server/client application, in which the server hosts the virtual world and coordinates the user interactions. In order to use *Twinity*, users have to download the *Twinity* client. The client allows the user to control the physical representation of the user’s character in the virtual world, also called the “avatar”. Thus the client is responsible for displaying the graphics, calculating the effects of physical interactions, handling the user’s input and synchronizing the 3D data and user actions with the *Twinity* server.

Each NPC comprises two major parts: whereas its avatar is the physical appearance of the NPC in the virtual world, the “conversational agent” provides the actual control logic which controls the avatar autonomously. It is in particular able to hold a conversation with *Twinity* users in that it reacts to a user’s presence, interprets user’s utterances in dialog context and generates adequate responses.

The *KomParse* server is a multi-client, multi-threaded server written in Java that hosts the conversational agents for the NPCs (section 5). The NPC’s avatar, on the other hand, is realized by a modified *Twinity* client. We utilize the Python interface provided by the *Twinity* client to call our own plugin which opens a bidirectional socket connection to the *KomParse* server. The plugin is started together with the *Twinity* client and serves as a mediator between the *Twinity* server and the *KomParse* server from then on (fig. 3). It sends all in-game events relevant to our system to the server and translates the commands sent by the server into *Twinity*-specific actions.

The integration architecture allows us to be maximally independent of the specific game platform. Rather than using the particular programming language and development environment of the platform for realizing the conversational agent or reimplementing a whole client/server protocol for connecting the avatar to the corresponding agent, we use an interface tailored to the specific needs of our system. Thus the *KomParse* system

can be naturally extended to other platforms since only the avatar interfaces have to be adapted.

The integration architecture also has the advantage that the necessary services can be easily distributed in a networked multi-platform environment. The *Twinity* clients require a Microsoft Windows machine with a 3D graphics card supporting DirectX 9.0c or higher, 1 GB RAM and a CPU core per instance. The *KomParse* server requires roughly 1 GB RAM. The triple store is run as a separate server process and is accessed by an XML-RPC interface. Roughly 1.2 GB RAM are required for loading our current knowledge base.

5 Conversational Agent: KomParse Dialog System

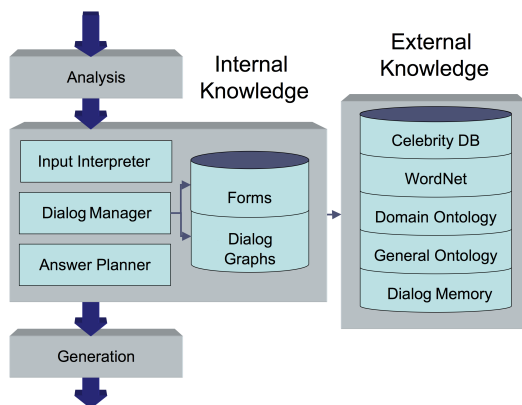


Figure 4: Dialog System: Conversational Agent

The *KomParse* dialog system, the main functionality of the conversational agent, consists of the following three major components: input analyzer, dialog manager and output generator (fig.4).

The *input analyzer* is responsible for the linguistic analysis of the user’s textual input including preprocessing such as string cleaning, part-of-speech tagging, named entity recognition, parsing and semantic interpretation. It yields a semantic representation which is the input for the dialog manager.

The *dialog manager* takes the result of the input analyzer and delivers an interpretation based on the dialog context and the available knowledge. It also controls the task conversation chain and handles user requests. The dialog manager determines the next system action based on the interpreted parameters.

The *output generator* realizes the action defined by the dialog manager with its multimodal gener-

ation competence. The generated results can be verbal, gestural or a combination of both.

As mentioned above, our dialog system has to deal with two different scenarios. While the focal point of the bartender agent lies in the question answering functionality, the furniture sales agent is driven by a complex dialog task model based on a dialog graph. Thus, the bartender agent relies mainly on question answering technology, in that it needs to understand questions and extract the right answer from our knowledge bases, whereas the sales agent has to accommodate various dialog situations with respect to a sales scenario. It therefore has to understand the dialog acts intended by the user and trigger the corresponding reactions, such as presenting an object, memorizing user preferences, negotiating further sales goals, etc.

The task model for sales conversations is inspired by a corpus resulting from the annotation of a Wizard-of-Oz experiment in the furniture sales agent scenario carried out by our project partner at ZAS (Bertomeu and Benz, 2009). In these experiments, 18 users spent one hour each on furnishing a virtual living room in a *Twinity* apartment by talking to a human wizard controlling the virtual sales agent. The final corpus consists of 18 dialogs containing 3,171 turns with 4,313 utterances and 23,015 alpha-numerical strings (words). The following example shows a typical part of such a conversation:

USR.1: And do we have a little side table for the TV?
NPC.1: I could offer you another small table or a sideboard.
USR.2: Then I'll take a sideboard that's similar to my shelf.
NPC.2: Let me check if we have something like that.

Table 1: Example Conversation from the Wizard-of-Oz Experiment

The flow of the task-based conversation is controlled by a data-driven finite-state model, which is the backbone of the dialog manager. During a sales conversation, objects and features of objects mentioned by the NPC and the user are extracted from the knowledge bases and added into the underspecified graph nodes and edges at runtime. This strategy keeps the finite-state graph as small as possible. Discussed objects and their features are stored in a frame-based sub-component named "form". The form contains entries which correspond to ontological concepts in the furni-

ture ontology. During conversation, these entries will be specified with the values of the properties of the discussed objects. This frame-based approach increases the flexibility of the dialog manager (McTear, 2002) and is particularly useful for a task-driven dialog system. As long as the negotiated object is not yet fully specified, the form represents the underspecified object description according to the ontology concept. Every time the user states a new preference or request, the form is enriched with additional features until the set of objects is small enough to be presented to the user for final selection. Thus the actual flow of dialog according to the task model does not have to be expressed by the graph but can be derived on demand from the knowledge and handled by the form which in turn activates the appropriate dialog subgraphs. This combination of graph-based dialog models and form-based task modelling effectively accounts for the interaction of sequential dialog strategies and the non-sequential nature of complex dialog goals.

Given a resolved semantic representation, the dialog manager triggers either a semantic search in the knowledge bases to deliver factual answers as needed in a gossip conversation or a further dialog response for example providing choices for the user in a sales domain. The semantic search is needed in both domains. In case that the semantic representation can neither be resolved in the task domain nor in the gossip domain, it gets passed to the embedded A.L.I.C.E. chatbot that uses its own understanding and generation components (Wallace and Bush, 2001).

5.1 Semantic Representation

The input understanding of the system is implemented as one single understanding pipeline. The understanding pipeline delivers a semantic representation which is the basis for the decision of the dialog manager which action to perform next.

This semantic representation can be extracted from the user input by our understanding component via a robust hybrid approach: either via a number of surface patterns containing regular expressions or via patterns reflecting the syntactic analysis of a dependency parser (de Marneffe and Manning, 2008).

The representation's structure is inspired by our knowledge representation design described in section 3 as well as by predicate logic. The core of the

representation is a predicate-argument structure limited to two arguments including message type and the whole syntactic information found by the analysis pipeline. The field "Message Type" can have one of the following values: wh-question, yes/no-question, declarative. Predicates can often be instantiated with the lemmatized matrix verb of the successfully analysed piece of the input. If the input contains a wh-question, the questioned fact is marked as an unfilled argument slot. The general structure can be simplified described as:

```
<PREDICATE, ARG1, ARG2, [message-type]>
```

The following examples show the structure used for different input:

- "Who is the boyfriend of Madonna?"
<hasBoyfriend, Madonna, ?, [wh]>
- "I want to buy a sofa."
<buy, I, "a sofa", [declarative]>

5.2 Information Extraction

Both scenarios make use of state-of-the-art information extraction approaches to extract the important pieces from the user input. While the bartender depends on relation extraction to detect the fact or relation questioned by the user (Xu et al., 2007), the sales agent uses information extraction methods to recognize user wishes and demands. As a result, the questioned fact or the demanded object feature equals the ontology structure containing the knowledge needed to handle the user input. The input "Do you have any red couches?" for example needs to get processed by the system in such a way that the information regarding the sofa with red color is extracted.

This is done by the system in a data-driven way. The input analysis first tries to find a demanded object in the input via asking the ontology: Every object which can be discussed in the scenario is encoded in the sales agents knowledge base. This can be seen as a Named Entity Recognition step. In case of success, the system tries to detect one of the possible relations of the object found in the input. This is achieved by querying the ontology about what kind of relations the identified object can satisfy. Possible relations are encoded in the class description of the given object. As a result the system can detect a relation "hasColour" for the found object "sofa" and the color value "red". The found information gets inserted into the form which gets more and more similar or if possible equal to a search query via RDF.

Input	Extracted Information	Knowledge Base
Do you have any red couches?	#hasMainColour (#Sofa, #Reds)	{f:Sofa_Alatea f:hasMainColour f:Reds}, {f:Sofa_Franky f:hasMainColour f:Reds}, {...}
I would like a leather one	#hasMaterial(#Sofa, #Leather)	{f:Sofa_Alatea f:hasMaterial f:leather}, {f:Sofa_Franky f:hasMaterial f:leather}, {...}
Who is the boyfriend of Madonna?	#hasBoyfriend(#Madonna, ?)	{f:Madonna f:hasBoyfriend f:GuyRichie}, {f:Madonna f:hasBoyfriend f:VanillaIce}, {...}

Figure 5: Comparison of Input, Extracted Information and Knowledge Base

6 Conclusion and Future Work

The *KomParse* system demonstrates an attractive application area for dialog systems that bears great future potential. Natural language dialog with NPCs is an important factor in making virtual worlds more interesting, interactive and immersive. Virtual worlds with conversing characters will also find many additional applications in education, marketing, and entertainment.

KomParse is an ambitious and nevertheless pragmatic attempt to bring NLP into the world of virtual games. We develop a new strategy to integrate task models and domain ontologies into dialog models. This strategy is useful for task-driven NPCs such as furniture sellers. With the chatty bartender, a combination of task-specific dialog and domain-specific question answering enables a smart wide-domain off-task conversation. Since the online game employs bubble-chat as a mode of communication in addition to Voice-over-IP, we are able to test our dialog system in a real-time application without being hindered by imperfect speech recognition.

The system presented here is still work in progress. The next goals will include various evaluation steps. On the one hand we will focus on single components like hybrid parsing of input utterances and dialog interpretation in terms of precision and recall. On the other hand an evaluation of the two different scenarios regarding the usability are planned in experiments with end users. Moreover we will integrate some opinion mining and sentiment analysis functionality which can be helpful to better detect and understand the user's preferences in the furniture sales agents scenario.

Acknowledgements

The project *KomParse* is funded by the ProFIT programme of the Federal State of Berlin, co-

funded by the EFRE programme of the European Union. The research presented here is additionally supported through a grant to the project TAKE, funded by the German Ministry for Education and Research (BMBF, FKZ: 01IW08003). Many thanks go to our project partners at the Centre for General Linguistics (ZAS) in Berlin as well as to the supporting company Metaversum.

References

- Peter Adolphs, Xiwen Cheng, Tina Klüwer, Hans Uszkoreit, and Feiyu Xu. 2010. Question answering biographic information and social network powered by the semantic web. In *Proceedings of LREC 2010*, Valletta, Malta.
- William Sims Bainbridge. 2007. The scientific research potential of virtual worlds. *Science*, 317.
- Nuria Bertomeu and Anton Benz. 2009. Annotation of joint projects and information states in human-npc dialogues. In *Proceedings of the First International Conference on Corpus Linguistics (CILC-09)*, Murcia, Spain.
- Marie C. de Marneffe and Christopher D. Manning. 2008. The Stanford typed dependencies representation. In *Coling 2008: Proceedings of the workshop on Cross-Framework and Cross-Domain Parser Evaluation*, Manchester, UK.
- Michael F. McTear. 2002. Spoken dialogue technology: enabling the conversational user interface. *ACM Comput. Surv.*, 34(1).
- Richard Wallace and Noel Bush. 2001. Artificial intelligence markup language (aiml) version 1.0.1 (2001). Unpublished A.L.I.C.E. AI Foundation Working Draft (rev 006).
- Feiyu Xu, Hans Uszkoreit, and Hong Li. 2007. A seed-driven bottom-up machine learning framework for extracting relations of various complexity. In *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, pages 584–591, Prague, Czech Republic, June. Association for Computational Linguistics.

An Open-Source Package for Recognizing Textual Entailment

Milen Kouylekov and Matteo Negri
FBK - Fondazione Bruno Kessler
Via Sommarive 18, 38100 Povo (TN), Italy
[kouylekov, negri]@fbk.eu

Abstract

This paper presents a general-purpose open source package for recognizing Textual Entailment. The system implements a collection of algorithms, providing a configurable framework to quickly set up a working environment to experiment with the RTE task. Fast prototyping of new solutions is also allowed by the possibility to extend its modular architecture. We present the tool as a useful resource to approach the Textual Entailment problem, as an instrument for didactic purposes, and as an opportunity to create a collaborative environment to promote research in the field.

1 Introduction

Textual Entailment (TE) has been proposed as a unifying generic framework for modeling language variability and semantic inference in different Natural Language Processing (NLP) tasks. The Recognizing Textual Entailment (RTE) task (Dagan and Glickman, 2007) consists in deciding, given two text fragments (respectively called *Text* - *T*, and *Hypothesis* - *H*), whether the meaning of *H* can be inferred from the meaning of *T*, as in:

T: "Yahoo acquired Overture"
H: "Yahoo owns Overture"

The RTE problem is relevant for many different areas of text processing research, since it represents the core of the semantic-oriented inferences involved in a variety of practical NLP applications including Question Answering, Information Retrieval, Information Extraction, Document Summarization, and Machine Translation. However, in spite of the great potential of integrating RTE into complex NLP architectures, little has been done to actually move from the controlled scenario pro-

posed by the RTE evaluation campaigns¹ to more practical applications. On one side, current RTE technology might not be mature enough to provide reliable components for such integration. Due to the intrinsic complexity of the problem, in fact, state of the art results still show large room for improvement. On the other side, the lack of available tools makes experimentation with the task, and the fast prototyping of new solutions, particularly difficult. To the best of our knowledge, the broad literature describing RTE systems is not accompanied with a corresponding effort on making these systems open-source, or at least freely available. We believe that RTE research would significantly benefit from such availability, since it would allow to quickly set up a working environment for experiments, encourage participation of newcomers, and eventually promote state of the art advances.

The main contribution of this paper is to present the latest release of EDITS (Edit Distance Textual Entailment Suite), a freely available, open source software package for recognizing Textual Entailment. The system has been designed following three basic requirements:

Modularity. System architecture is such that the overall processing task is broken up into major modules. Modules can be composed through a configuration file, and extended as plug-ins according to individual requirements. System's workflow, the behavior of the basic components, and their IO formats are described in a comprehensive documentation available upon download.

Flexibility. The system is general-purpose, and suited for any TE corpus provided in a simple XML format. In addition, both language dependent and language independent configurations are allowed by algorithms that manipulate different representations of the input data.

¹TAC RTE Challenge: <http://www.nist.gov/tac>
EVALITA TE task: <http://evalita.itc.it>

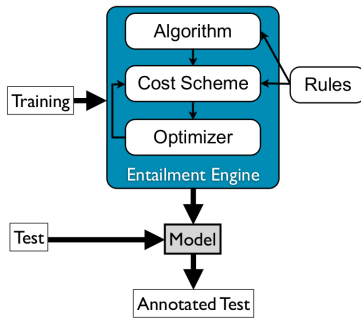


Figure 1: Entailment Engine, main components and workflow

Adaptability. Modules can be tuned over training data to optimize performance along several dimensions (*e.g.* overall Accuracy, Precision/Recall trade-off on YES and NO entailment judgements). In addition, an optimization component based on genetic algorithms is available to automatically set parameters starting from a basic configuration.

EDITS is open source, and available under GNU Lesser General Public Licence (LGPL). The tool is implemented in Java, it runs on Unix-based Operating Systems, and has been tested on MAC OSX, Linux, and Sun Solaris. The latest release of the package can be downloaded from <http://edits.fbk.eu>.

2 System Overview

The EDITS package allows to:

- Create an *Entailment Engine* (Figure 1) by defining its basic components (*i.e.* algorithms, cost schemes, rules, and optimizers);
- Train such *Entailment Engine* over an annotated RTE corpus (containing T-H pairs annotated in terms of entailment) to learn a *Model*;
- Use the *Entailment Engine* and the *Model* to assign an entailment judgement and a confidence score to each pair of an un-annotated test corpus.

EDITS implements a distance-based framework which assumes that the probability of an entailment relation between a given T-H pair is inversely proportional to the *distance* between T and H (*i.e.* the higher the distance, the lower is the probability of entailment). Within this framework the system implements and harmonizes different approaches to distance computation, providing both *edit distance* algorithms, and *similarity* algorithms (see

Section 3.1). Each algorithm returns a normalized distance score (a number between 0 and 1). At a training stage, distance scores calculated over annotated T-H pairs are used to estimate a threshold that best separates positive from negative examples. The threshold, which is stored in a *Model*, is used at a test stage to assign an entailment judgement and a confidence score to each test pair.

In the creation of a distance Entailment Engine, algorithms are combined with *cost schemes* (see Section 3.2) that can be *optimized* to determine their behaviour (see Section 3.3), and optional external knowledge represented as *rules* (see Section 3.4). Besides the definition of a single Entailment Engine, a unique feature of EDITS is that it allows for the combination of multiple Entailment Engines in different ways (see Section 4.4).

Pre-defined basic components are already provided with EDITS, allowing to create a variety of entailment engines. Fast prototyping of new solutions is also allowed by the possibility to extend the modular architecture of the system with new algorithms, cost schemes, rules, or plug-ins to new language processing components.

3 Basic Components

This section overviews the main components of a distance Entailment Engine, namely: *i)* algorithms, *iii)* cost schemes, *iii)* the cost optimizer, and *iv)* entailment/contradiction rules.

3.1 Algorithms

Algorithms are used to compute a distance score between T-H pairs.

EDITS provides a set of predefined algorithms, including edit distance algorithms, and similarity algorithms adapted to the proposed distance framework. The choice of the available algorithms is motivated by their large use documented in RTE literature².

Edit distance algorithms cast the RTE task as the problem of mapping the whole content of H into the content of T. Mappings are performed as sequences of editing operations (*i.e.* insertion, deletion, substitution of text portions) needed to transform T into H, where each edit operation has a cost associated with it. The distance algorithms available in the current release of the system are:

²Detailed descriptions of all the systems participating in the TAC RTE Challenge are available at <http://www.nist.gov/tac/publications>

- Token Edit Distance: a token-based version of the Levenshtein distance algorithm, with edit operations defined over sequences of tokens of T and H;
- Tree Edit Distance: an implementation of the algorithm described in (Zhang and Shasha, 1990), with edit operations defined over single nodes of a syntactic representation of T and H.

Similarity algorithms are adapted to the EDITS distance framework by transforming measures of the lexical/semantic similarity between T and H into distance measures. These algorithms are also adapted to use the three edit operations to support overlap calculation, and define term weights. For instance, substitutable terms in T and H can be treated as equal, and non-overlapping terms can be weighted proportionally to their insertion/deletion costs. Five similarity algorithms are available, namely:

- Word Overlap: computes an overall (distance) score as the proportion of common words in T and H;
- Jaro-Winkler distance: a similarity algorithm between strings, adapted to similarity on words;
- Cosine Similarity: a common vector-based similarity measure;
- Longest Common Subsequence: searches the longest possible sequence of words appearing both in T and H in the same order, normalizing its length by the length of H;
- Jaccard Coefficient: confronts the intersection of words in T and H to their union.

3.2 Cost Schemes

Cost schemes are used to define the cost of each edit operation.

Cost schemes are defined as XML files that explicitly associate a cost (a positive real number) to each edit operation applied to elements of T and H. Elements, referred to as A and B, can be of different types, depending on the algorithm used. For instance, Tree Edit Distance will manipulate *nodes* in a dependency tree representation, whereas Token Edit Distance and similarity algorithms will manipulate *words*. Figure 2 shows an example of

```
<scheme>
  <insertion><cost>10</cost></insertion>
  <deletion><cost>10</cost></deletion>
  <substitution>
    <condition>(equals A B)</condition>
    <cost>0</cost>
  </substitution>
  <substitution>
    <condition>(not (equals A B))</condition>
    <cost>20</cost>
  </substitution>
</scheme>
```

Figure 2: Example of XML Cost Scheme

cost scheme, where edit operation costs are defined as follows:

Insertion(B)=10 - inserting an element B from H to T, no matter what B is, always costs 10;

Deletion(A)=10 - deleting an element A from T, no matter what A is, always costs 10;

substitution(A,B)=0 if A=B - substituting A with B costs 0 if A and B are equal;

substitution(A,B)=20 if A≠B - substituting A with B costs 20 if A and B are different.

In the distance-based framework adopted by EDITS, the interaction between algorithms and cost schemes plays a central role. Given a T-H pair, in fact, the distance score returned by an algorithm directly depends on the cost of the operations applied to transform T into H (edit distance algorithms), or on the cost of mapping words in H with words in T (similarity algorithms). Such interaction determines the overall behaviour of an Entailment Engine, since distance scores returned by the same algorithm with different cost schemes can be considerably different. This allows users to define (and optimize, as explained in Section 3.3) the cost schemes that best suit the RTE data they want to model³.

EDITS provides two predefined cost schemes:

- Simple Cost Scheme - the one shown in Figure 2, setting fixed costs for each edit operation.
- IDF Cost Scheme - insertion and deletion costs for a word w are set to the inverse document frequency of w ($IDF(w)$). The substitution cost is set to 0 if a word $w1$ from T and a word $w2$ from H are the same, and $IDF(w1)+IDF(w2)$ otherwise.

³For instance, when dealing with T-H pairs composed by texts that are much longer than the hypotheses (as in the RTE5 Campaign), setting low deletion costs avoids penalization to short Hs fully contained in the Ts.

In the creation of new cost schemes, users can express edit operation costs, and conditions over the A and B elements, using a meta-language based on a lisp-like syntax (*e.g.* (+ (IDF A) (IDF B)), (not (equals A B))). The system also provides functions to access data stored in hash files. For example, the IDF Cost Scheme accesses the IDF values of the most frequent 100K English words (calculated on the Brown Corpus) stored in a file distributed with the system. Users can create new hash files to collect statistics about words in other languages, or other information to be used inside the cost scheme.

3.3 Cost Optimizer

A cost optimizer is used to adapt cost schemes (either those provided with the system, or new ones defined by the user) to specific datasets.

The optimizer is based on cost adaptation through genetic algorithms, as proposed in (Mehdad, 2009). To this aim, cost schemes can be parametrized by externalizing as parameters the edit operations costs. The optimizer iterates over training data using different values of these parameters until an optimal set is found (*i.e.* the one that best performs on the training set).

3.4 Rules

Rules are used to provide the Entailment Engine with knowledge (*e.g.* lexical, syntactic, semantic) about the probability of entailment or contradiction between elements of T and H. Rules are invoked by cost schemes to influence the cost of substitutions between elements of T and H. Typically, the cost of the substitution between two elements A and B is inversely proportional to the probability that A entails B.

Rules are stored in XML files called Rule Repositories, with the format shown in Figure 3. Each rule consists of three parts: *i)* a left-hand side, *ii)* a right-hand side, *iii)* a probability that the left-hand side entails (or contradicts) the right-hand side.

EDITS provides three predefined sets of lexical entailment rules acquired from lexical resources widely used in RTE: WordNet⁴, Lin’s word similarity dictionaries⁵, and VerbOcean⁶.

⁴<http://wordnet.princeton.edu>

⁵<http://webdocs.cs.ualberta.ca/~lindek/downloads.htm>

⁶<http://demo.patrickpantel.com/Content/verbocan>

```
<rule entailment="ENTAILMENT">
  <t>acquire</t>
  <h>own</h>
  <probability>0.95</probability>
</rule>
<rule entailment="CONTRADICTION">
  <t>beautiful</t>
  <h>ugly</h>
  <probability>0.88</probability>
</rule>
```

Figure 3: Example of XML Rule Repository

4 Using the System

This section provides basic information about the use of EDITS, which can be run with commands in a Unix Shell. A complete guide to all the parameters of the main script is available as HTML documentation downloadable with the package.

4.1 Input

The input of the system is an entailment corpus represented in the EDITS Text Annotation Format (ETAF), a simple XML internal annotation format. ETAF is used to represent both the input T-H pairs, and the entailment and contradiction rules. ETAF allows to represent texts at two different levels: *i)* as sequences of tokens with their associated morpho-syntactic properties, or *ii)* as syntactic trees with structural relations among nodes.

Plug-ins for several widely used annotation tools (including TreeTagger, Stanford Parser, and OpenNLP) can be downloaded from the system’s website. Users can also extend EDITS by implementing plug-ins to convert the output of other annotation tools in ETAF.

Publicly available RTE corpora (RTE 1-3, and EVALITA 2009), annotated in ETAF at both the annotation levels, are delivered together with the system to be used as first experimental datasets.

4.2 Configuration

The creation of an Entailment Engine is done by defining its basic components (algorithms, cost schemes, optimizer, and rules) through an XML configuration file. The configuration file is divided in modules, each having a set of options. The following XML fragment represents a simple example of configuration file:

```
<module alias="distance">
  <module alias="tree"/>
  <module alias="xml">
    <option name="scheme-file"
```

```

        value="IDF_Scheme.xml" />
    </module>
    <module alias="pso" />
</module>

```

This configuration defines a distance Entailment Engine that combines Tree Edit Distance as a core distance algorithm, and the predefined IDF Cost Scheme that will be optimized on training data with the Particle Swarm Optimization algorithm (“*pso*”) as in (Mehdad, 2009). Adding external knowledge to an entailment engine can be done by extending the configuration file with a reference to a rules file (e.g. “*rules.xml*”) as follows:

```

<module alias="rules">
  <option name="rules-file"
    value="rules.xml" />
</module>

```

4.3 Training and Test

Given a configuration file and an RTE corpus annotated in ETAF, the user can run the **training** procedure to learn a model. At this stage, EDITS allows to tune performance along several dimensions (e.g. overall Accuracy, Precision/Recall trade-off on YES and/or NO entailment judgments). By default the system maximizes the overall accuracy (distinction between YES and NO pairs). The output of the training phase is a *model*: a zip file that contains the learned threshold, the configuration file, the cost scheme, and the entailment/contradiction rules used to calculate the threshold. The explicit availability of all this information in the model allows users to share, replicate and modify experiments⁷.

Given a model and an un-annotated RTE corpus as input, the **test** procedure produces a file containing for each pair: *i*) the decision of the system (YES, NO), *ii*) the confidence of the decision, *iii*) the entailment score, *iv*) the sequence of edit operations made to calculate the entailment score.

4.4 Combining Engines

A relevant feature of EDITS is the possibility to combine multiple Entailment Engines into a single one. This can be done by grouping their definitions as sub-modules in the configuration file. EDITS allows users to define customized combination strategies, or to use two predefined combination modalities provided with the package,

⁷Our policy is to publish online the models we use for participation in the RTE Challenges. We encourage other users of EDITS to do the same, thus creating a collaborative environment, allow new users to quickly modify working configurations, and replicate results.

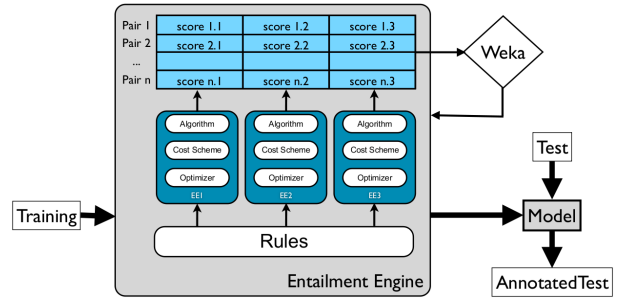


Figure 4: Combined Entailment Engines

namely: *i*) Linear Combination, and *ii*) Classifier Combination. The two modalities combine in different ways the entailment scores produced by multiple independent engines, and return a final decision for each T-H pair.

Linear Combination returns an overall entailment score as the weighted sum of the entailment scores returned by each engine:

$$score_{combination} = \sum_{i=0}^n score_i * weight_i \quad (1)$$

In this formula, $weight_i$ is an ad-hoc weight parameter for each entailment engine. Optimal weight parameters can be determined using the same optimization strategy used to optimize the cost schemes, as described in Section 3.3.

Classifier Combination is similar to the approach proposed in (Malakasiotis and Androutsopoulos, 2007), and is based on using the entailment scores returned by each engine as features to train a classifier (see Figure 4). To this aim, EDITS provides a plug-in that uses the Weka⁸ machine learning workbench as a core. By default the plug-in uses an SVM classifier, but other Weka algorithms can be specified as options in the configuration file.

The following configuration file describes a combination of two engines (i.e. one based on Tree Edit Distance, the other based on Cosine Similarity), used to train a classifier with Weka⁹.

```

<module alias="weka">
  <module alias="distance">
    <module alias="tree" />
    <module alias="xml">
      <option name="scheme-file"
        value="IDF_Scheme.xml" />
    </module>
  </module>
</module>

```

⁸<http://www.cs.waikato.ac.nz/ml/weka>

⁹A linear combination can be easily obtained by changing the alias of the highest-level module (“weka”) into “linear”.

```

<module alias="distance">
  <module alias="cosine"/>
  <module alias="IDF_Scheme.xml"/>
</module>
</module>

```

5 Experiments with EDITS

To give an idea of the potentialities of the EDITS package in terms of flexibility and adaptability, this section reports some results achieved in RTE-related tasks by previous versions of the tool. The system has been tested in different scenarios, ranging from the evaluation of standalone systems within task-specific RTE Challenges, to their integration in more complex architectures.

As regards the RTE Challenges, in the last years EDITS has been used to participate both in the PASCAL/TAC RTE Campaigns for the English language (Mehdad et al., 2009), and in the EVALITA RTE task for Italian (Cabrio et al., 2009). In the last RTE-5 Campaign the result achieved in the traditional “2-way Main task” (60.17% Accuracy) roughly corresponds to the performance of the average participating systems (60.36%). In the “Search” task (which consists in finding all the sentences that entail a given H in a given set of documents about a topic) the same configuration achieved an F1 of 33.44%, ranking 3rd out of eight participants (average score 29.17% F1). In the EVALITA 2009 RTE task, EDITS ranked first with an overall 71.0% Accuracy. To promote the use of EDITS and ease experimentation, the complete models used to produce each submitted run can be downloaded with the system. An improved model obtained with the current release of EDITS, and trained over RTE-5 data (61.83% Accuracy on the “2-way Main task” test set), is also available upon download.

As regards application-oriented integrations, EDITS has been successfully used as a core component in a Restricted-Domain Question Answering system within the EU-Funded QALL-ME Project¹⁰. Within this project, an entailment-based approach to Relation Extraction has been defined as the task of checking for the existence of entailment relations between an input question (the *text* in RTE parlance), and a set of textual realizations of domain-specific binary relations (the *hypotheses* in RTE parlance). In recognizing 14 relations relevant in the CINEMA domain present in a collection of spoken English requests, the system

¹⁰<http://qallme.fbk.eu>

achieved an F1 of 72.9%, allowing to return correct answers to 83% of 400 test questions (Negri and Kouylekov, 2009).

6 Conclusion

We have presented the first open source package for recognizing Textual Entailment. The system offers a modular, flexible, and adaptable working environment to experiment with the task. In addition, the availability of pre-defined system configurations, tested in the past Evaluation Campaigns, represents a first contribution to set up a collaborative environment, and promote advances in RTE research. Current activities are focusing on the development of a Graphical User Interface, to further simplify the use of the system.

Acknowledgments

The research leading to these results has received funding from the European Community’s Seventh Framework Programme (FP7/2007-2013) under Grant Agreement n. 248531 (CoSyne project).

References

- Prodromos Malakasiotis and Ion Androutsopoulos 2007. *Learning Textual Entailment using SVMs and String Similarity Measures*. Proc. of the ACL ’07 Workshop on Textual Entailment and Paraphrasing.
- Ido Dagan and Oren Glickman 2004. *Probabilistic Textual Entailment: Generic Applied Modeling of Language Variability*. Proc. of the PASCAL Workshop on Learning Methods for Text Understanding and Mining.
- Kaizhong Zhang and Dennis Shasha 1990. *Fast Algorithm for the Unit Cost Editing Distance Between Trees*. Journal of Algorithms. vol.11.
- Yashar Mehdad 2009. *Automatic Cost Estimation for Tree Edit Distance Using Particle Swarm Optimization*. Proc. of ACL-IJCNLP 2009.
- Matteo Negri and Milen Kouylekov 2009. *Question Answering over Structured Data: an Entailment-Based Approach to Question Analysis*. Proc. of RANLP-2009.
- Elena Cabrio, Yashar Mehdad, Matteo Negri, Milen Kouylekov, and Bernardo Magnini 2009. *Recognizing Textual Entailment for Italian EDITS @ EVALITA 2009* Proc. of EVALITA 2009.
- Yashar Mehdad, Matteo Negri, Elena Cabrio, Milen Kouylekov, and Bernardo Magnini 2009. *Recognizing Textual Entailment for English EDITS @ TAC 2009* To appear in Proceedings of TAC 2009.

Personalising speech-to-speech translation in the EMIME project

Mikko Kurimo^{1†}, William Byrne⁶, John Dines³, Philip N. Garner³, Matthew Gibson⁶, Yong Guan⁵, Teemu Hirsimäki¹, Reima Karhila¹, Simon King², Hui Liang³, Keiichiro Oura⁴, Lakshmi Saheer³, Matt Shannon⁶, Sayaka Shiota⁴, Jilei Tian⁵, Keiichi Tokuda⁴, Mirjam Wester², Yi-Jian Wu⁴, Junichi Yamagishi²

¹ Aalto University, Finland, ² University of Edinburgh, UK, ³ Idiap Research Institute, Switzerland, ⁴ Nagoya Institute of Technology, Japan, ⁵ Nokia Research Center Beijing, China, ⁶ University of Cambridge, UK

†Corresponding author: Mikko.Kurimo@tkk.fi

Abstract

In the EMIME project we have studied unsupervised cross-lingual speaker adaptation. We have employed an HMM statistical framework for both speech recognition and synthesis which provides transformation mechanisms to adapt the synthesized voice in TTS (text-to-speech) using the recognized voice in ASR (automatic speech recognition). An important application for this research is personalised speech-to-speech translation that will use the voice of the speaker in the input language to utter the translated sentences in the output language. In mobile environments this enhances the users' interaction across language barriers by making the output speech sound more like the original speaker's way of speaking, even if she or he could not speak the output language.

1 Introduction

A mobile real-time speech-to-speech translation (S2ST) device is one of the grand challenges in natural language processing (NLP). It involves several important NLP research areas: automatic speech recognition (ASR), statistical machine translation (SMT) and speech synthesis, also known as text-to-speech (TTS). In recent years significant advance have also been made in relevant technological devices: the size of powerful computers has decreased to fit in a mobile phone and fast WiFi and 3G networks have spread widely to connect them to even more powerful computation servers. Several hand-held S2ST applications and devices have already become available, for ex-

ample by IBM, Google or Jibbig¹, but there are still serious limitations in vocabulary and language selection and performance.

When an S2ST device is used in practical human interaction across a language barrier, one feature that is often missed is the personalization of the output voice. Whoever speaks to the device in what ever manner, the output voice always sounds the same. Producing high-quality synthesis voices is expensive and even if the system had many output voices, it is hard to select one that would sound like the input voice. There are many features in the output voice that could raise the interaction experience to a much more natural level, for example, emotions, speaking rate, loudness and the speaker identity.

After the recent development in hidden Markov model (HMM) based TTS, it has become possible to adapt the output voice using model transformations that can be estimated from a small number of speech samples. These techniques, for instance the maximum likelihood linear regression (MLLR), are adopted from HMM-based ASR where they are very powerful in fast adaptation of speaker and recording environment characteristics (Gales, 1998). Using hierarchical regression trees, the TTS and ASR models can further be coupled in a way that enables unsupervised TTS adaptation (King et al., 2008). In unsupervised adaptation samples are annotated by applying ASR. By eliminating the need for human intervention it becomes possible to perform voice adaptation for TTS in almost real-time.

The target in the EMIME project² is to study unsupervised cross-lingual speaker adaptation for S2ST systems. The first results of the project have

¹<http://www.jibbig.com>

²<http://emime.org>

been, for example, to bridge the gap between the ASR and TTS (Dines et al., 2009), to improve the baseline ASR (Hirsimäki et al., 2009) and SMT (de Gispert et al., 2009) systems for morphologically rich languages, and to develop robust TTS (Yamagishi et al., 2010). The next step has been preliminary experiments in intra-lingual and cross-lingual speaker adaptation (Wu et al., 2008). For cross-lingual adaptation several new methods have been proposed for mapping the HMM states, adaptation data and model transformations (Wu et al., 2009).

In this presentation we can demonstrate the various new results in ASR, SMT and TTS. Even though the project is still ongoing, we have an initial version of mobile S2ST system and cross-lingual speaker adaptation to show.

2 Baseline ASR, TTS and SMT systems

The baseline ASR systems in the project are developed using the HTK toolkit (Young et al., 2001) for Finnish, English, Mandarin and Japanese. The systems can also utilize various real-time decoders such as Julius (Kawahara et al., 2000), Juicer at IDIAP and the TTK decoder (Hirsimäki et al., 2006). The main structure of the baseline systems for each of the four languages is similar and fairly standard and in line with most other state-of-the-art large vocabulary ASR systems. Some special flavors have been added, such as the morphological analysis for Finnish (Hirsimäki et al., 2009). For speaker adaptation, the MLLR transformation based on hierarchical regression classes is included for all languages.

The baseline TTS systems in the project utilize the HTS toolkit (Yamagishi et al., 2009) which is built on top of the HTK framework. The HMM-based TTS systems have been developed for Finnish, English, Mandarin and Japanese. The systems include an average voice model for each language trained over hundreds of speakers taken from standard ASR corpora, such as Speecon (Iskra et al., 2002). Using speaker adaptation transforms, thousands of new voices have been created (Yamagishi et al., 2010) and new voices can be added using a small number of either supervised or unsupervised speech samples. Cross-lingual adaptation is possible by creating a mapping between the HMM states in the input and the output language (Wu et al., 2009).

Because the resources of the EMIME project

have been focused on ASR, TTS and speaker adaptation, we aim at relying on existing solutions for SMT as far as possible. New methods have been studied concerning the morphologically rich languages (de Gispert et al., 2009), but for the S2ST system we are currently using Google translate³.

3 Demonstrations to show

3.1 Monolingual systems

In robust speech synthesis, a computer can learn to speak in the desired way after processing only a relatively small amount of training speech. The training speech can even be a normal quality recording outside the studio environment, where the target speaker is speaking to a standard microphone and the speech is not annotated. This differs dramatically from conventional TTS, where building a new voice requires an hour or more careful repetition of specially selected prompts recorded in an anechoic chamber with high quality equipment.

Robust TTS has recently become possible using the statistical HMM framework for both ASR and TTS. This framework enables the use of efficient speaker adaptation transformations developed for ASR to be used also for the TTS models. Using large corpora collected for ASR, we can train average voice models for both ASR and TTS. The training data may include a small amount of speech with poor coverage of phonetic contexts from each single speaker, but by summing the material over hundreds of speakers, we can obtain sufficient models for an average speaker. Only a small amount of adaptation data is then required to create transformations for tuning the average voice closer to the target voice.

In addition to the supervised adaptation using annotated speech, it is also possible to employ ASR to create annotations. This unsupervised adaptation enables the system to use a much broader selection of sources, for example, recorded samples from the internet, to learn a new voice.

The following systems will demonstrate the results of monolingual adaptation:

1. In *EMIME Voice cloning in Finnish and English* the goal is that the users can clone their own voice. The user will dictate for about

³<http://translate.google.com>

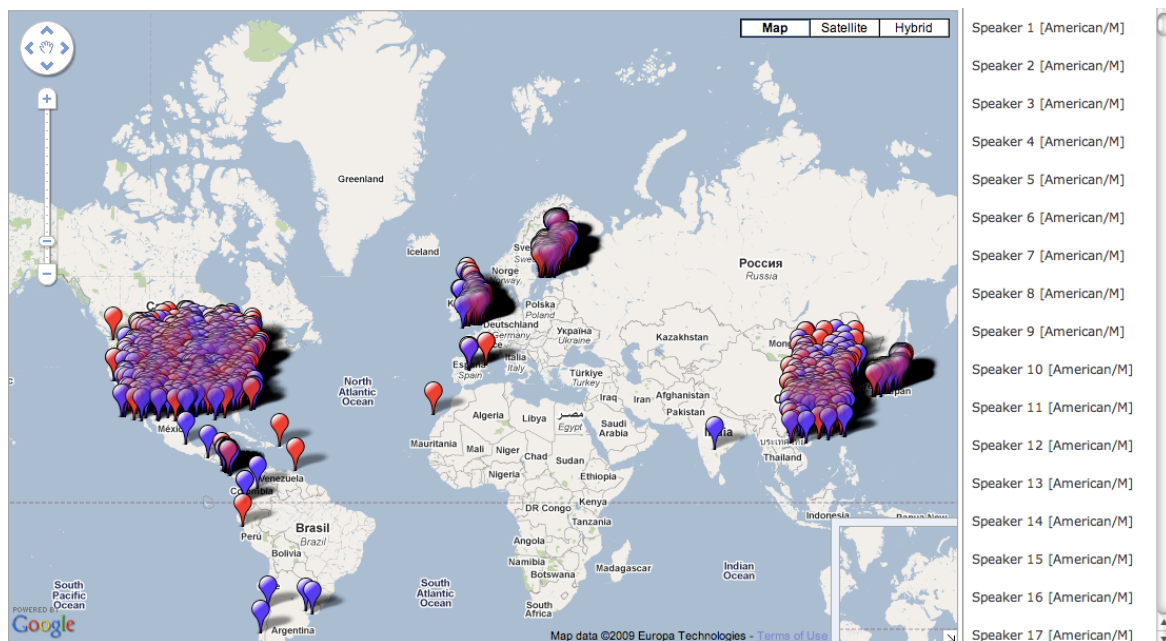


Figure 1: Geographical representation of HTS voices trained on ASR corpora for EMIME projects. Blue markers show male speakers and red markers show female speakers. Available online via <http://www.emime.org/learn/speech-synthesis/listen/Examples-for-D2.1>

10 minutes and then after half an hour of processing time, the TTS system has transformed the average model towards the user's voice and can speak with this voice. The cloned voices may become especially valuable, for example, if a person's voice is later damaged in an accident or by a disease.

2. In *EMIME Thousand voices map* the goal is to browse the world's largest collection of synthetic voices by using a world map interface (Yamagishi et al., 2010). The user can zoom in the world map and select any voice, which are organized according to the place of living of the adapted speaker, to utter the given sentence. This interactive geographical representation is shown in Figure 1. Each marker corresponds to an individual speaker. Blue markers show male speakers and red markers show female speakers. Some markers are in arbitrary locations (in the correct country) because precise location information is not available for all speakers. This geographical representation, which includes an interactive TTS demonstration of many of the voices, is available from the URL provided. Clicking on a marker will play synthetic speech from that speaker⁴. As well as

⁴Currently the interactive mode supports English and Spanish only. For other languages this only provides pre-

being a convenient interface to compare the many voices, the interactive map is an attractive and easy-to-understand demonstration of the technology being developed in EMIME.

3. The models developed in the HMM framework can be demonstrated also in adaptation of an ASR system for *large-vocabulary continuous speech recognition*. By utilizing morpheme-based language models instead of word-based models the Finnish ASR system is able to cover practically an unlimited vocabulary (Hirsimäki et al., 2006). This is necessary for morphologically rich languages where, due to inflection, derivation and composition, there exists so many different word forms that word based language modeling becomes impractical.

3.2 Cross-lingual systems

In the EMIME project the goal is to learn cross-lingual speaker adaptation. Here the output language ASR or TTS system is adapted from speech samples in the input language. The results so far are encouraging, especially for TTS: Even though the cross-lingual adaptation may somewhat degrade the synthesis quality, the adapted speech now sounds more like the target speaker. Several recent evaluations of the cross-lingual speaker

synthesised examples, but we plan to add an interactive type-in text-to-speech feature in the near future.

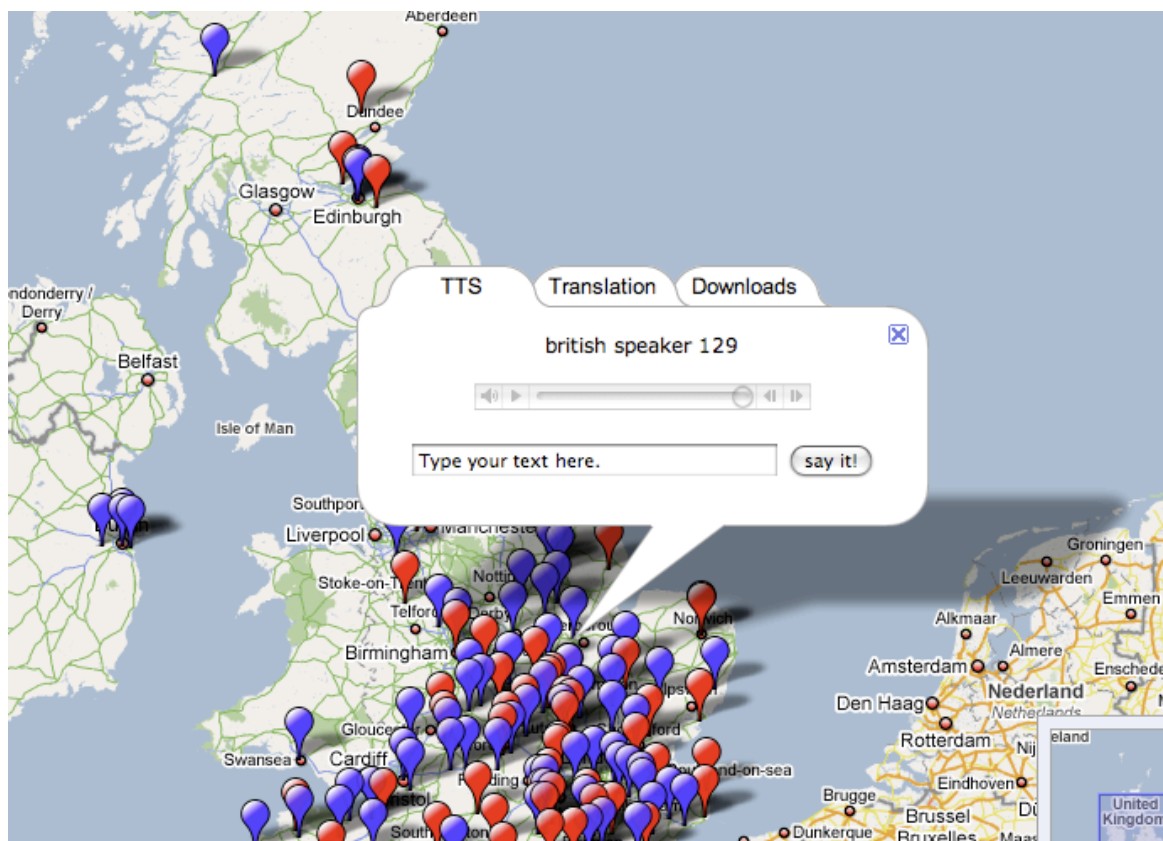


Figure 2: All English HTS voices can be used as online TTS on the geographical map.

adaptation methods can be found in (Gibson et al., 2010; Oura et al., 2010; Liang et al., 2010; Oura et al., 2009).

The following systems have been created to demonstrate cross-lingual adaptation:

1. In *EMIME Cross-lingual Finnish/English and Mandarin/English TTS adaptation* the input language sentences dictated by the user will be used to learn the characteristics of her or his voice. The adapted cross-lingual model will be used to speak output language (English) sentences in the user's voice. The user does not need to be bilingual and only reads sentences in their native language.
2. In *EMIME Real-time speech-to-speech mobile translation demo* two users will interact using a pair of mobile N97 devices (see Figure 3). The system will recognize the phrase the other user is speaking in his native language and translate and speak it in the native language of the other user. After a few sentences the system will have the speaker adaptation transformations ready and can apply them in the synthesized voices to make them sound more like the original speaker instead of a standard voice. The first real-time demo

version is available for the Mandarin/English language pair.

3. *The morpheme-based translation system* for Finnish/English and English/Finnish can be compared to a word based translation for arbitrary sentences. The morpheme-based approach is particularly useful for language pairs where one or both languages are morphologically rich ones where the amount and complexity of different word forms severely limits the performance for word-based translation. The morpheme-based systems can learn translation models for phrases where morphemes are used instead of words (de Gispert et al., 2009). Recent evaluations (Kurimo et al., 2009) have shown that the performance of the unsupervised data-driven morpheme segmentation can rival the conventional rule-based ones. This is very useful if hand-crafted morphological analyzers are not available or their coverage is not sufficient for all languages.

Acknowledgments

The research leading to these results was partly funded from the European Community's Seventh

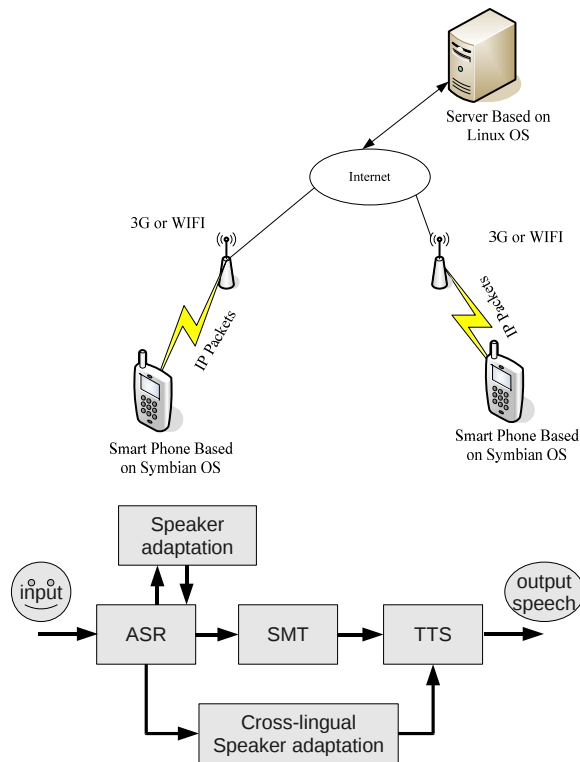


Figure 3: EMIME Real-time speech-to-speech mobile translation demo Framework Programme (FP7/2007-2013) under grant agreement 213845 (the EMIME project).

References

- A. de Gispert, S. Virpioja, M. Kurimo, and W. Byrne. 2009. Minimum Bayes risk combination of translation hypotheses from alternative morphological decompositions. In *Proc. NAACL-HLT*.
- J. Dines, J. Yamagishi, and S. King. 2009. Measuring the gap between HMM-based ASR and TTS. In *Proc. Interspeech '09*, Brighton, UK.
- M. Gales. 1998. Maximum likelihood linear transformations for HMM-based speech recognition. *Computer Speech and Language*, 12(2):75–98.
- M. Gibson, T. Hirsimäki, R. Karhila, M. Kurimo, and W. Byrne. 2010. Unsupervised cross-lingual speaker adaptation for HMM-based speech synthesis using two-pass decision tree construction. In *Proc. of ICASSP*, page to appear, March.
- T. Hirsimäki, M. Creutz, V. Siivola, M. Kurimo, S. Virpioja, and J. Pytkönen. 2006. Unlimited vocabulary speech recognition with morph language models applied to finnish. *Computer Speech & Language*, 20(4):515–541, October.
- T. Hirsimäki, J. Pytkönen, and M. Kurimo. 2009. Importance of high-order n-gram models in morph-based speech recognition. *IEEE Trans. Audio, Speech, and Language Process.*, 17:724–732.
- D. Iskra, B. Grosskopf, K. Marasek, H. van den Heuvel, F. Diehl, and A. Kiessling. 2002. SPEECON speech databases for consumer devices: Database specification and validation. In *Proc. LREC*, pages 329–333.
- T. Kawahara, A. Lee, T. Kobayashi, K. Takeda, N. Minematsu, S. Sagayama, K. Itou, A. Ito, M. Yamamoto, A. Yamada, T. Utsuro, and K. Shikano. 2000. Free software toolkit for japanese large vocabulary continuous speech recognition. In *Proc. ICSLP-2000*, volume 4, pages 476–479.
- S. King, K. Tokuda, H. Zen, and J. Yamagishi. 2008. Unsupervised adaptation for HMM-based speech synthesis. In *Proc. Interspeech 2008*, pages 1869–1872, September.
- Mikko Kurimo, Sami Virpioja, Ville T. Turunen, Graeme W. Blackwood, and William Byrne. 2009. Overview and results of Morpho Challenge 2009. In *Working Notes for the CLEF 2009 Workshop*, Corfu, Greece, September.
- H. Liang, J. Dines, and L. Saheer. 2010. A comparison of supervised and unsupervised cross-lingual speaker adaptation approaches for HMM-based speech synthesis. In *Proc. of ICASSP*, page to appear, March.
- Keiichiro Oura, Junichi Yamagishi, Simon King, Mirjam Wester, and Keiichi Tokuda. 2009. Unsupervised speaker adaptation for speech-to-speech translation system. In *Proc. SLP (Spoken Language Processing)*, number 356 in 109, pages 13–18.
- K. Oura, K. Tokuda, J. Yamagishi, S. King, and M. Wester. 2010. Unsupervised cross-lingual speaker adaptation for HMM-based speech synthesis. In *Proc. of ICASSP*, page to appear, March.
- Y.-J. Wu, S. King, and K. Tokuda. 2008. Cross-lingual speaker adaptation for HMM-based speech synthesis. In *Proc. of ISCSLP*, pages 1–4, December.
- Y.-J. Wu, Y. Nankaku, and K. Tokuda. 2009. State mapping based method for cross-lingual speaker adaptation in HMM-based speech synthesis. In *Proc. of Interspeech*, pages 528–531, September.
- J. Yamagishi, T. Nose, H. Zen, Z.-H. Ling, T. Toda, K. Tokuda, S. King, and S. Renals. 2009. Robust speaker-adaptive HMM-based text-to-speech synthesis. *IEEE Trans. Audio, Speech and Language Process.*, 17(6):1208–1230. (in press).
- J. Yamagishi, B. Usabaev, S. King, O. Watts, J. Dines, J. Tian, R. Hu, K. Oura, K. Tokuda, R. Karhila, and M. Kurimo. 2010. Thousands of voices for hmm-based speech synthesis. *IEEE Trans. Speech, Audio & Language Process.* (in press).

S. Young, G. Everman, D. Kershaw, G. Moore, J. Odell, D. Ollason, V. Valtchev, and P. Woodland, 2001. *The HTK Book Version 3.1*, December.

Hunting for the Black Swan: Risk Mining from Text

Jochen L. Leidner and Frank Schilder

Thomson Reuters Corporation

Research & Development

610 Opperman Drive, St. Paul, MN 55123 USA

FirstName.LastName@ThomsonReuters.com

Abstract

In the business world, analyzing and dealing with risk permeates all decisions and actions. However, to date, *risk identification*, the first step in the risk management cycle, has always been a manual activity with little to no intelligent software tool support. In addition, although companies are required to list risks to their business in their annual SEC filings in the USA, these descriptions are often very high-level and vague.

In this paper, we introduce *Risk Mining*, which is the task of identifying a set of risks pertaining to a business area or entity. We argue that by combining Web mining and Information Extraction (IE) techniques, risks can be detected automatically before they materialize, thus providing valuable business intelligence.

We describe a system that induces a risk taxonomy with concrete risks (e.g., interest rate changes) at its leaves and more abstract risks (e.g., financial risks) closer to its root node. The taxonomy is induced via a bootstrapping algorithms starting with a few seeds. The risk taxonomy is used by the system as input to a risk monitor that matches risk mentions in financial documents to the abstract risk types, thus bridging a lexical gap. Our system is able to automatically generate company specific “risk maps”, which we demonstrate for a corpus of earnings report conference calls.

1 Introduction

Any given human activity with a particular intended outcome is bound to face a non-zero likelihood of failure. In business, companies are exposed to market risks such as new competitors, disruptive technologies, change in customer attitudes, or a changes in government legislation that can dramatically affect their profitability or threaten their business model or mode of operation. Therefore, any tool to assist in the elicitation of otherwise unforeseen risk factors carries tremendous potential value.

However, it is very hard to identify risks exhaustively, and some types (commonly referred to as the *unknown unknowns*) are especially elusive: if a *known unknown* is the established knowledge that important risk factors are known, but it is unclear whether and when they become realized,

then an *unknown unknown* is the lack of awareness, in practice or in principle, of circumstances that may impact the outcome of a project, for example. Nassim Nicholas Taleb calls these “black swans” (Taleb, 2007).

Companies in the US are required to disclose a list of potential risks in their annual Form 10-K SEC filings in order to warn (potential) investors, and risks are frequently the topic of conference phone calls about a company’s earnings. These risks are often reported in general terms, in particular, because it is quite difficult to pinpoint the *unknown unknown*, i.e. what kind of risk is concretely going to materialize. On the other hand, there is a stream of valuable evidence available on the Web, such as news messages, blog entries, and analysts’ reports talking about companies’ performance and products. Financial analysts and risk officers in large companies have not enjoyed any text analytics support so far, and risk lists devised using questionnaires or interviews are unlikely to be exhaustive due to small sample size, a gap which we aim to address in this paper.

To this end, we propose to use a combination of Web Mining (WM) and Information Extraction (IE) to assist humans interested in risk (with respect to an organization) and to bridge the gap between the general language and concrete risks. We describe our system, which is divided in two main parts: (a) an offline **Risk Miner** that facilitates the *risk identification* step of the risk management process, and an online (b) **Risk Monitor** that supports the *risk monitoring* step (cf. Figure 2). In addition, a **Risk Mapper** can aggregate and visualize the evidence in the form of a *risk map*. Our risk mining algorithm combines Riloff hyponym patterns with recursive Web pattern bootstrapping and a graph representation.

We do not know of any other implemented end-to-end system for computer-assisted risk identification/visualization using text mining technology.

2 Related Work

Financial IE. IE systems have been applied to the financial domain on Message Understanding Contest (MUC) like tasks, ranging from named entity tagging to slot filling in templates (Costantino, 1992).

Automatic Knowledge Acquisition. (Hearst, 1992) pioneered the pattern-based extraction of hyponyms from corpora, which laid the groundwork for subsequent work, and which included extraction of knowledge from the Web (e.g. (Etzioni et al., 2004)). To improve precision was the mission of (Kozareva et al., 2008), which was designed to extract hyponymy, but they did so at the expense of recall, using longer *dual anchored patterns* and a pattern linkage graph. However, their method is by its very nature unable to deal with low-frequency items, and their system does not contain a chunker, so only single term items can be extracted. De Saenger et al. (De Saeger et al., 2008) describe an approach that extracts instances of the “trouble” or “obstacle” relations from the Web in the form of pairs of fillers for these binary relations. Their approach, which is described for the Japanese language, uses support vector machine learning and relies on a Japanese syntactic parser, which permits them to process negation. In contrast, and unlike their method, we pursue a more general, open-ended search process, which does not impose as much a priori knowledge. Also, they create a set of pairs, whereas our approach creates a taxonomy tree as output. Most importantly though, our approach is not driven by frequency, and was instead designed to work especially with rare occurrences in mind to permit “black swan”-type risk discovery.

Correlation of Volatility and Text. (Kogan et al., 2009) study the correlation between share price volatility, a proxy for risk, and a set of trigger words occurring in 60,000 SEC 10-K filings from 1995-2006. Since the disclosure of a company’s risks is mandatory by law, SEC reports provide a rich source. Their trigger words are selected a priori by humans; in contrast, risk mining as exercised in this paper aims to find risk-indicative words and phrases automatically.

Kogan and colleagues attempt to find a regression model using very simple unigram features based on whole documents that predicts volatility, whereas our goal is to automatically extract patterns to be used as alerts.

Speculative Language & NLP. Light et al. (Light et al., 2004) found that sub-string matching of 14 pre-defined string literals outperforms an SVM classifier using bag-of-words features in the task of speculative language detection in medical abstracts. (Goldberg et al., 2009) are concerned with automatic recognition of human wishes, as expressed in human notes for Year’s Eve. They use a bi-partite graph-based approach, where one kind of node (content node) represents things people wish for (“world peace”) and the other kind of node (template nodes) represent templates that extract them (e.g. “I wish for ___”). Wishes can be seen as positive Q , in our formalization.

3 Data

We apply the mined risk extraction patterns to a corpus of financial documents. The data originates from the StreetEvents database and was kindly provided to us by Starmine, a Thomson Reuters company. In particular, we are dealing with 170k earning calls transcripts, a text type that contains monologue (company executives reporting about their company’s performance and general situation) as well as dialogue (in the form of questions and answers at the end of each conference call). Participants typically include select business analysts from investment banks, and the calls are published afterwards for the shareholders’ benefits. Figure 1 shows some example excerpts. We randomly took a sample of $N=6,185$ transcripts to use them in our risk alerting experiments.¹

4 Method

4.1 System

The overall system is divided into two core parts: (a) Risk Mining and (b) Risk Monitoring (cf. Figure 2). For demonstration purposes, we add a (c) Risk Mapper, a visualization component. We describe how a variety of risks can be identified given a normally very high-level description of risks, as one can find in earnings reports, other financial news, or the risk section of 10-K SEC filings. Starting with rather abstract descriptions such as *operational risks* and hyponym-inducing pattern “< RISK > such as *”, we use the Web to mine pages from which we can harvest additional,

¹We could also use this data for risk mining, but did not try this due to the small size of the dataset compared to the Web.

CEO: As announced last evening, during our third quarter, we will take the difficult but necessary step to seize [cease] manufacturing at our nearly 100 year old Pennsylvania House casegood plant in Lewisburg, Pennsylvania as well as the nearby Pennsylvania House dining room chair assembly facility in White Deer. Also, the three Lewisburg area warehouses will be consolidated as we assess the logistical needs of the casegood group's existing warehouse operations at an appropriate time in the future to minimize any disruption of service to our customers. This will result in the loss of 425 jobs or approximately 15% of the casegood group's current employee base.

Analyst: Okay, so your comments – and I guess I don't know – I can figure out, as you correctly helped me through, what dollar contribution at GE. I don't know the net equipment sales number last quarter and this quarter. But it sounded like from your comments that if you exclude these fees, that equipment sales were probably flattish. Is that fair to say?

CEO: We're not breaking out the origination fee from the equipment fee, but I think in total, I would say flattish to slightly up.

Figure 1: Example sentences from the earnings conference call dataset. Top: main part. Bottom: Q&A.

and eventually more concrete, candidates, and relate them to risk types via a transitive chain of binary IS-A relations. Contrary to the related work, we use a base NP chunker and download the full pages returned by the search engine rather than search snippets in order to be able to extract risk phrases rather than just terms, which reduces contextual ambiguity and thus increases overall precision. The taxonomy learning method described in the following subsection determines a risk taxonomy and new risks patterns.

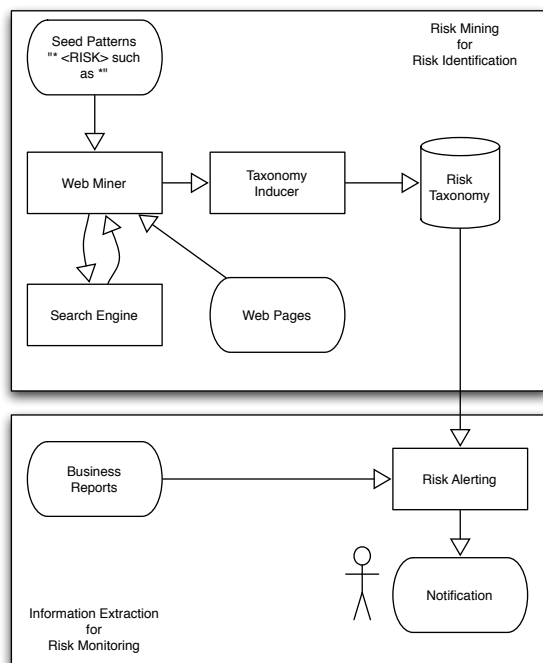


Figure 2: The risk mining and monitoring system architecture

The second part of the system, the Risk Monitor, takes the risks from the risk taxonomy and uses them for monitoring financial text streams such as news, SEC filings, or (in our use case) earnings reports. Using this, an analyst is then able to identify concrete risks in news messages and link them to the high-level risk descriptions. He

or she may want to identify operational risks such as fraud for a particular company, for instance. The risk taxonomy can also derive further risks in this category (e.g., faulty components, brakes) for exploration and drill-down analysis. Thus, news reports about faulty breaks in (e.g. Toyota) or volcano outbreaks (e.g. Iceland) can be directly linked to the risk as stated in earnings reports or security filings.

Our Risk Miner and Risk Monitor are implemented in Perl, with the graph processing of the taxonomy implemented in SWI-Prolog, whereas the Risk Mapper exists in two versions, a static image generator for R^2 and, alternatively, an interactive Web page (DHTML, JavaScript, and using Google's Chart API). We use the Yahoo Web search API.

4.2 Taxonomy induction method

Using frequency to compute confidence in a pattern does not work for risk mining, however, because mention of particular risks might be rare. Instead of frequency based indicators (n-grams, frequency weights), we rely on two types of structural confidence validation, namely (a) previously identified risks and (b) previously acquired structural patterns. Note, however, that we can still use PageRank, a popularity-based graph algorithm, because multiple patterns might be connected to a risk term or phrase, even in the absence of frequency counts for each (i.e., we interpret popularity as having multiple sources of support).

1. Risk Candidate Extraction Step. The first step is used to extract a list of risks based on high precision patterns. However, it has been shown that the use of such patterns (e.g., *such as*) quickly lead to an decrease in precision. Ideally, we want to retrieve specific risks by re-applying the the extract risk descriptions:

²<http://www.r-project.org>

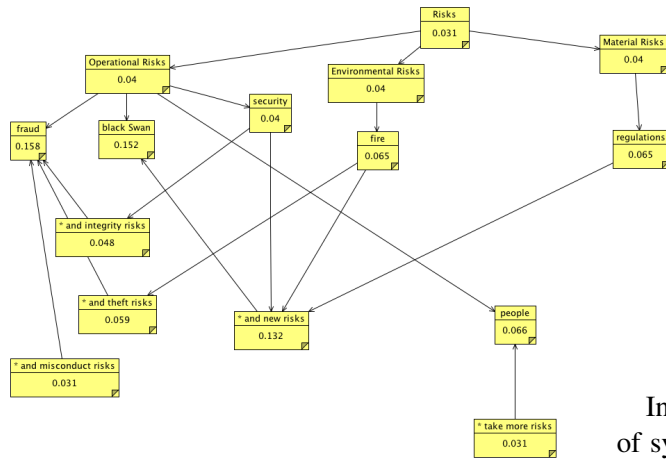


Figure 3: A sample IS-A and Pattern network with sample PageRank scores.

- (a) Take a seed, instantiate " $\langle \text{SEED} \rangle$ " such as "*" pattern with seed, extract candidates:

Input: risks

Method: apply pattern " $\langle \text{SEED} \rangle$ " such as " $\langle \text{INSTANCE} \rangle$ ", where $\langle \text{SEED} \rangle = \text{risks}$

Output: list of instances (e.g., *faulty components*)

- (b) For each candidate from the list of instances, we find a set of additional candidate hyponyms.

Input: faulty components

Method: apply pattern " $\langle \text{SEED} \rangle$ " such as " $\langle \text{INSTANCE} \rangle$ ", where $\langle \text{SEED} \rangle = \text{faulty components}$

Output: list of instances (e.g., *brake*)

2. Risk Validation. Since the Risk Candidate extraction step will also find many false positives, we need to factor in information that validates that the extracted risk is indeed a risk. We do this by constructing a possible pattern containing this new risk.

- (a) Append "*" risks" to the output of 1(b) in order to make sure that the candidate occurs in a risk context.

Input: brake(s)

Pattern: "brake(s) * risk(s)"

Output: a list of patterns (e.g., *minimize such risks, raising the risk*)

- (b) extract new risk pattern by substituting the risk candidate with $\langle \text{RISK} \rangle$; creating a limited number of variations

Input: list of all patterns mined from step 2

(a)

Method: create more pattern variations, such as " $\langle \text{RISK} \rangle$ minimize such risks", "raising the risk of $\langle \text{RISK} \rangle$ " etc.

Output: list of new potential risks (e.g., *deflation*), but also many false positives (e.g., *way*, as in *The best way to minimize such risks*).

In order to benefit from any human observations of system errors in future runs, we also extended the system so as to read in a partial list of pre-defined risks at startup time, which can guide the risk miner; while technically different from active learning, this approach was somewhat inspired by it (but our feedback is more loose).

3. Constructing Risk Graph. We have now reached the point where we constructed a graph with risks and patterns. Risks are connected via IS-A links; risks and patterns are connected via PATTERN links. Note that there are links from risks to patterns and from patterns to risks; some risks back-pointed by a pattern may actually not be a risk (e.g., *people*). However, this node is also not connected to a more abstract risk node and will therefore have a low PageRank score. Risks that are connected to patterns that have a high authority (i.e., pointing to by many other links) are highly ranked within PageRank (Figure 3). The risk *black Swan*, for example, has only one pattern it occurs in, but this pattern can be filled by many other risks (e.g., *fire*, *regulations*). Hence, the PageRank score of the black swan is high similar to well known risks, such as *fraud*.

4.3 Risk alerting method

We compile the risk taxonomy into a trie automaton, and create a second trie for company names from the meta-data of our corpus. The Risk Monitor reads the two tries and uses the first to detect mentions of risks in the earning reports and the second one to tag company names, both using case-insensitive matching for better recall. Optionally, we can use Porter stemming during trie construction and matching to trade precision for even higher recall, but in the experiments reported here this is not used. Once a signal term or phrase matches, we look up its risk type in a hash table, take a note of the company that the current earnings report is about, and increase the frequency

liquidity IS-A financial risks
 credit IS-A financial risks
 direct risks IS-A financial risks
 fraud IS-A financial risks
 irregular activity IS-A operational risks
 process failure IS-A operational risks
 human error IS-A operational risks
 labor strikes IS-A operational risks
 customer acceptance IS-A IT market risks
 interest rate changes IS-A capital market risks
 uncertainty IS-A market risks
 volatility IS-A mean reverting market risks
 copyright infringement IS-A legal risks
 negligence IS-A other legal risks
 an unfair dismissal IS-A the legal risks
 Sarbanes IS-A legal risks
 government changes IS-A global political risks
 crime IS-A Social and political risks
 state intervention IS-A political risks
 terrorist acts IS-A geopolitical risks
 earthquakes IS-A natural disaster risks
 floods IS-A natural disaster risks
 global climate change IS-A environmental risks
 severe and extreme weather IS-A environmental risks
 internal cracking IS-A any technological risks
 GM technologies IS-A tech risks
 scalability issues IS-A technology risks
 viruses IS-A the technical risks

Figure 4: Selected financial risk tuples after Web validation.

count for this \langle company; risk type \rangle tuple, which we use for graphic rendering purposes.

4.4 Risk mapping method

To demonstrate the method presented here, we created a visualization that displays a *risk map*, which is a two dimensional table showing companies and the types of risk they are facing, together with bubble sizes proportional to the number of alerts that the Risk Monitor could discover in the corpus. The second option also permits the user to explore the detected risk mentions per company and by risk type.

5 Results

From the Web mining process, we obtain a set of pairs (Figure 4), from which the taxonomy is constructed. In one run with only 12 seeds (just the risk type names with variants), we obtained a taxonomy with 280 validated leave nodes that are connected transitively to the risks root node.

Our resulting system produces visualizations we call “risk maps”, because they graphically present the extracted risk types in aggregated form. A set of risk types can be selected for presentation as well as a set of companies of interest. A risk map display is then generated using either R (Figure 5) or an interactive Web page, depending on the user’s preference.

Qualitative error analysis. We inspected the output of the risk miner and observed the follow-



Figure 5: An Example Risk Map.

ing classes of issues: (a) chunker errors: if phrasal boundaries are placed at the wrong position, the taxonomy will include wrong relations. For example, deictic determiners such as “this” were a problem (e.g. that IS-A indirect risks) before we introduced a stop word filter that discards candidate tuples that contain no content words. Another prominent example is “short term” instead of the correct “short term risk”; (b) *semantic drift*³: due to polysemy, words and phrases can denote risk and non-risk meanings, depending on context. Talking about risks even a specific pattern such as “such as” [sic] is used by authors to induce a variety of perspectives on the topic of risk, and after several iterations negative effects of type (a) errors compound; (c) off-topic relations: the seeds are designed to induce a taxonomy specific to risk types. As a side effect, many (correct or incorrect) irrelevant relations are learned, e.g. credit and debit cards is-a money transfer. We currently discard these by virtue of ignoring all relations not transitively connected with the root node risks, so no formalized domain knowledge is required; (d) overlap: the concept space is divided up differently by different writers, both on the Web and in the risk management literature, and this is reflected by multiple category membership of many risks (e.g. is cash flow primarily an operational risk or a financial risk?). Currently, we do not deal with this phenomenon automatically; (e) redundant relations: at the time of writing, we do not cache all already extracted and validated risks/non-risks. This means there is room for improvement w.r.t. runtime, because we make more Web queries than strictly necessary. While we have not evaluated this system yet, we found by in-

³to use a term coined by Andy Lauriston

specting the output that our method is particularly effective for learning natural disasters and medical conditions, probably because they are well-covered by news sites and biomedical abstracts on the Web. We also found that some classes contain more noise than others, for example operational risk was less precise than financial risk, probably due to the lesser specificity of the former risk type.

6 Summary, Conclusions & Future Work

Summary of Contributions.

In this paper, we introduced the task of risk mining, which produces patterns that are useful in another task, risk alerting. Both tasks provide computational assistance to risk-related decision making in the financial sector. We described a special-purpose algorithm for inducing a risk taxonomy offline, which can then be used online to analyze earning reports in order to signal risks. In doing so, we have addressed two research questions of general relevance, namely how to extract *rare patterns*, for which frequency-based methods fail, and how to use the Web to bridge the *vocabulary gap*, i.e. how to match up terms and phrases in financial news prose with the more abstract language typically used in talking about risk in general.

We have described an implemented demonstrator system comprising an offline risk taxonomy miner, an online risk alerter and a visualization component that creates visual risk maps by company and risk type, which we have applied to a corpus of earnings call transcripts.

Future Work. Extracted negative and also positive risks can be used in many applications, ranging from *e-mail alerts* to determining *credit ratings*. Our preliminary work on risk maps can be put on a more theoretical footing (Hunter, 2000). After studying further how output of risk alerting correlates⁴ with non-textual signals like share price, risk detection signals could inform human or trading decisions.

Acknowledgments. We are grateful to Khalid Al-Kofahi, Peter Jackson and James Powell for supporting this work. Thanks to George Bonne, Ryan Roser, and Craig D'Alessio at Starmine, a Thomson Reuters company, for sharing the StreetEvents dataset with us, and to David Rosenblatt for discussions and to Jack Conrad for feedback on this paper.

⁴Our hypothesis is that risk patterns can outperform bag of words (Kogan et al., 2009).

References

- Marco Costantino. 1992. Financial information extraction using pre-defined and user-definable templates in the LOLITA system. *Proceedings of the Fifteenth International Conference on Computational Linguistics (COLING 1992)*, vol. 4, pages 241–255.
- Stijn De Saeger, Kentaro Torisawa, and Jun'ichi Kazama. 2008. Looking for trouble. In *Proceedings of the 22nd International Conference on Computational Linguistics (COLING 2008)*, pages 185–192, Morristown, NJ, USA. Association for Computational Linguistics.
- Oren Etzioni, Michael J. Cafarella, Doug Downey, Stanley Kok, Ana-Maria Popescu, Tal Shaked, Stephen Soderland, Daniel S. Weld, and Alexander Yates. 2004. Web-scale information extraction in KnowItAll: preliminary results. In Stuart I. Feldman, Mike Uretsky, Marc Najork, and Craig E. Wills, editors, *Proceedings of the 13th international conference on World Wide Web (WWW 2004)*, New York, NY, USA, May 17-20, 2004, pages 100–110. ACM.
- Andrew B. Goldberg, Nathanael Fillmore, David Andrzejewski, Zhiting Xu, Bryan Gibson, and Xiaojin Zhu. 2009. May all your wishes come true: A study of wishes and how to recognize them. In *Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 263–271, Boulder, Colorado, June. Association for Computational Linguistics.
- Marti Hearst. 1992. Automatic acquisition of hyponyms from large text corpora. In *Proceedings of the Fourteenth International Conference on Computational Linguistics (COLING 1992)*.
- Anthony Hunter. 2000. Ramification analysis using causal mapping. *Data and Knowledge Engineering*, 32:200–227.
- Shimon Kogan, Dmitry Levin, Bryan R. Routledge, Jacob S. Sagi, and Noah A. Smith. 2009. Predicting risk from financial reports with regression. In *Proceedings of the Joint International Conference on Human Language Technology and the Annual Meeting of the North American Chapter of the Association for Computational Linguistics (HLT-NAACL)*.
- Zornitsa Kozareva, Ellen Riloff, and Eduard Hovy. 2008. Semantic class learning from the web with hyponym pattern linkage graphs. In *Proceedings of ACL-HLT*, pages 1048–1056, Columbus, OH, USA. Association for Computational Linguistics.
- Marc Light, Xin Ying Qiu, and Padmini Srinivasan. 2004. The language of bioscience: Facts, speculations, and statements in between. In *BioLINK 2004: Linking Biological Literature, Ontologies and Databases*, pages 17–24. ACL.
- Nassim Nicholas Taleb. 2007. *The Black Swan: The Impact of the Highly Improbable*. Random House.

Speech-driven Access to the Deep Web on Mobile Devices

Taniya Mishra and Srinivas Bangalore

AT&T Labs - Research

180 Park Avenue

Florham Park, NJ 07932 USA.

{taniya,srini}@research.att.com.

Abstract

The Deep Web is the collection of information repositories that are not indexed by search engines. These repositories are typically accessible through web forms and contain dynamically changing information. In this paper, we present a system that allows users to access such rich repositories of information on mobile devices using spoken language.

1 Introduction

The World Wide Web (WWW) is the largest repository of information known to mankind. It is generally agreed that the WWW continues to significantly enrich and transform our lives in unprecedented ways. Be that as it may, the WWW that we encounter is limited by the information that is accessible through search engines. Search engines, however, do not index a large portion of WWW that is variously termed as the *Deep Web*, *Hidden Web*, or *Invisible Web*.

Deep Web is the information that is in proprietary databases. Information in such databases is usually more structured and changes at higher frequency than textual web pages. It is conjectured that the Deep Web is 500 times the size of the surface web. Search engines are unable to index this information and hence, unable to retrieve it for the user who may be searching for such information. So, the only way for users to access this information is to find the appropriate web-form, fill in the necessary search parameters, and use it to query the database that contains the information that is being searched for. Examples of such web forms include, movie, train and bus times, and airline/hotel/restaurant reservations.

Contemporaneously, the devices to access information have moved out of the office and home environment into the open world. The ubiquity of mobile devices has made information access an any time, any place activity. However, informa-

tion access using text input on mobile devices is tedious and unnatural because of the limited screen space and the small (or soft) keyboards. In addition, by the *mobile* nature of these devices, users often like to use them in hands-busy environments, ruling out the possibility of typing text. Filling web-forms using the small screens and tiny keyboards of mobile devices is neither easy nor quick.

In this paper, we present a system, *Qme!*, designed towards providing a spoken language interface to the Deep Web. In its current form, *Qme!* provides a unified interface onn iPhone (shown in Figure 1) that can be used by users to search for *static* and *dynamic* questions. *Static* questions are questions whose answers to these questions remain the same irrespective of when and where the questions are asked. Examples of such questions are *What is the speed of light?*, *When is George Washington's birthday?*. For *static* questions, the system retrieves the answers from an archive of human generated answers to questions. This ensures higher accuracy for the answers retrieved (if found in the archive) and also allows us to retrieve related questions on the user's topic of interest.

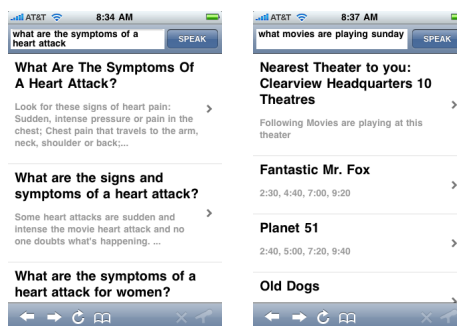


Figure 1: Retrieval results for static and dynamic questions using *Qme!*

Dynamic questions are questions whose answers depend on when and where they are asked. Examples of such questions are *What is the stock price of General Motors?*, *Who won the game last night?*, *What is playing at the theaters near me?*.

The answers to dynamic questions are often part of the Deep Web. Our system retrieves the answers to such dynamic questions by parsing the questions to retrieve pertinent search keywords, which are in turn used to query information databases accessible over the Internet using web forms. However, the internal distinction between dynamic and static questions, and the subsequent differential treatment within the system is seamless to the user. The user simply uses a single unified interface to ask a question and receive a collection of *answers* that potentially address her question directly.

The layout of the paper is as follows. In Section 2, we present the system architecture. In Section 3, we present bootstrap techniques to distinguish dynamic questions from static questions, and evaluate the efficacy of these techniques on a test corpus. In Section 4, we show how our system retrieves answers to dynamic questions. In Section 5, we show how our system retrieves answers to static questions. We conclude in Section 6.

2 Speech-driven Question Answer System

Speech-driven access to information has been a popular application deployed by many companies on a variety of information resources (Microsoft, 2009; Google, 2009; YellowPages, 2009; vlingo.com, 2009). In this prototype demonstration, we describe a speech-driven question-answer application. The system architecture is shown in Figure 2.

The user of this application provides a spoken language query to a mobile device intending to find an answer to the question. The speech recognition module of the system recognizes the spoken query. The result from the speech recognizer can be either a single-best string or a weighted word lattice.¹ This textual output of recognition is then used to classify the user query either as a dynamic query or a static query. If the user query is static, the result of the speech recognizer is used to search a large corpus of question-answer pairs to retrieve the relevant answers. The retrieved results are ranked using *tf.idf* based metric discussed in Section 5. If the user query is dynamic, the answers are retrieved by querying a web form from the appropriate web site (e.g. *www.fandango.com* for movie information). In Figure 1, we illustrate the answers that Qme! returns for static and dy-

¹For this paper, the ASR used to recognize these utterances incorporates an acoustic model adapted to speech collected from mobile devices and a four-gram language model that is built from the corpus of questions.

amic questions.

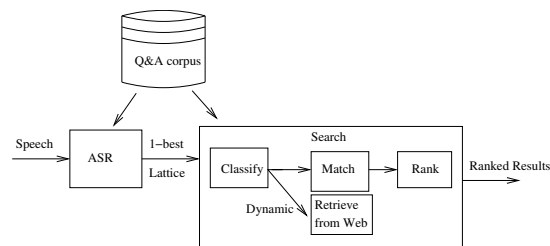


Figure 2: The architecture of the speech-driven question-answering system

2.1 Demonstration

In the demonstration, we plan to show the users static and dynamic query handling on an iPhone using spoken language queries. Users can use the iPhone and speak their queries using an interface provided by Qme!. A Wi-Fi access spot will make this demonstration more compelling.

3 Dynamic and Static Questions

As mentioned in the introduction, dynamic questions require accessing the hidden web through a web form with the appropriate parameters. Answers to dynamic questions cannot be preindexed as can be done for static questions. They depend on the time and geographical location of the question. In dynamic questions, there may be no explicit reference to time, unlike the questions in the TERQAS corpus (Radev and Sundheim., 2002) which explicitly refer to the temporal properties of the entities being questioned or the relative ordering of past and future events.

The time-dependency of a dynamic question lies in the temporal nature of its answer. For example, consider the question, *What is the address of the theater White Christmas is playing at in New York?* White Christmas is a seasonal play that plays in New York every year for a few weeks in December and January, but not necessarily at the same theater every year. So, depending when this question is asked, the answer will be different. If the question is asked in the summer, the answer will be “This play is not currently playing anywhere in NYC.” If the question is asked during December, 2009, the answer might be different than the answer given in December 2010, because the theater at which White Christmas is playing differs from 2009 to 2010.

There has been a growing interest in temporal analysis for question-answering since the late 1990’s. Early work on temporal expressions iden-

tification using a tagger culminated in the development of TimeML (Pustejovsky et al., 2001), a markup language for annotating temporal expressions and events in text. Other examples include, QA-by-Dossier with Constraints (Prager et al., 2004), a method of improving QA accuracy by asking auxiliary questions related to the original question in order to temporally verify and restrict the original answer. (Moldovan et al., 2005) detect and represent temporally related events in natural language using logical form representation. (Sagete et al., 2009) use the temporal relations in a question to decompose it into simpler questions, the answers of which are recomposed to produce the answers to the original question.

3.1 Question Classification: Dynamic and Static Questions

We automatically classify questions as dynamic and static questions. The answers to static questions can be retrieved from the QA archive. To answer dynamic questions, we query the database(s) associated with the topic of the question through web forms on the Internet. We first use a topic classifier to detect the topic of a question followed by a dynamic/static classifier trained on questions related to a topic, as shown in Figure 3. For the question *what movies are playing around me?*, we detect it is a movie related dynamic question and query a movie information web site (e.g. www.fandango.com) to retrieve the results based on the user’s GPS information.

Dynamic questions often contain temporal indexicals, i.e., expressions of the form *today*, *now*, *this week*, *two summers ago*, *currently*, *recently*, etc. Our initial approach was to use such signal words and phrases to automatically identify dynamic questions. The chosen signals were based on annotations in TimeML. We also included spatial indexicals, such as *here* and other clauses that were observed to be contained in dynamic questions such as *cost of*, and *how much is* in the list of signal phrases. These signals words and phrases were encoded into a regular-expression-based recognizer.

This regular-expression based recognizer identified 3.5% of our dataset – which consisted of several million questions – as dynamic. The type of questions identified were *What is playing in the movie theaters tonight?*, *What is tomorrow’s weather forecast for LA?*, *Where can I go to get Thai food near here?* However, random samplings of the same dataset, annotated by four independent human labelers, indicated that on average 13.5%

of the dataset is considered dynamic. This shows that the temporal and spatial indexicals encoded as a regular-expression based recognizer is unable to identify a large percentage of the dynamic questions.

This approach leaves out dynamic questions that do not contain temporal or spatial indexicals. For example, *What is playing at AMC Loew’s?*, or *What is the score of the Chargers and Dolphines game?* For such examples, considering the tense of the verb in question may help. The last two examples are both in the present continuous tense. But verb tense does not help for a question such as *Who got voted off Survivor?*. This question is certainly dynamic. The information that is most likely being sought by this question is what is the name of the person who got voted off the TV show Survivor *most recently*, and not what is the name of the person (or persons) who have gotten voted off the Survivor at some point in the past.

Knowing the broad topic (such as movies, current affairs, and music) of the question may be very useful. It is likely that there may be many dynamic questions about movies, sports, and finance, while history and geography may have few or none. This idea is bolstered by the following analysis. The questions in our dataset are annotated with a broad topic tag. Binning the 3.5% of our dataset identified as dynamic questions by their broad topic produced a long-tailed distribution. Of the 104 broad topics, the top-5 topics contained over 50% of the dynamic questions. These top five topics were sports, TV and radio, events, movies, and finance.

Considering the issues laid out in the previous section, our classification approach is to chain two machine-learning-based classifiers: a topic classifier chained to a dynamic/static classifier, as shown in Figure 3. In this architecture, we build one topic classifier, but several dynamic/static classifiers, each trained on data pertaining to one broad topic.

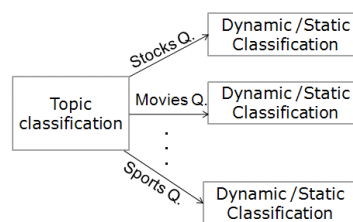


Figure 3: Chaining two classifiers

We used supervised learning to train the topic

classifier, since our entire dataset is annotated by human experts with topic labels. In contrast, to train a dynamic/static classifier, we experimented with the following three different techniques.

Baseline: We treat questions as dynamic if they contain temporal indexicals, e.g. *today, now, this week, two summers ago, currently, recently*, which were based on the TimeML corpus. We also included spatial indexicals such as *here*, and other substrings such as *cost of* and *how much is*. A question is considered static if it does not contain any such words/phrases.

Self-training with bagging: The general self-training with bagging algorithm (Banko and Brill, 2001). The benefit of self-training is that we can build a better classifier than that built from the small seed corpus by simply adding in the large unlabeled corpus without requiring hand-labeling.

Active-learning: This is another popular method for training classifiers when not much annotated data is available. The key idea in active learning is to annotate only those instances of the dataset that are most difficult for the classifier to learn to classify. It is expected that training classifiers using this method shows better performance than if samples were chosen randomly for the same human annotation effort.

We used the maximum entropy classifier in LLAMA (Haffner, 2006) for all of the above classification tasks. We have chosen the active learning classifier due to its superior performance and integrated it into the Qme! system. We provide further details about the learning methods in (Mishra and Bangalore, 2010).

3.2 Experiments and Results

3.2.1 Topic Classification

The topic classifier was trained using a training set consisting of over one million questions downloaded from the web which were manually labeled by human experts as part of answering the questions. The test set consisted of 15,000 randomly selected questions. Word trigrams of the question are used as features for a MaxEnt classifier which outputs a score distribution on all of the 104 possible topic labels. The error rate results for models selecting the top topic and the top two topics according to the score distribution are shown in Table 1. As can be seen these error rates are far lower than the baseline model of selecting the most frequent topic.

Model	Error Rate
Baseline	98.79%
Top topic	23.9%
Top-two topics	12.23%

Table 1: Results of topic classification

3.2.2 Dynamic/static Classification

As mentioned before, we experimented with three different approaches to bootstrapping a dynamic/static question classifier. We evaluated these methods on a 250 question test set drawn from the broad topic of *Movies*. The error rates are summarized in Table 2. We provide further details of this experiment in (Mishra and Bangalore, 2010).

Training approach	Lowest Error rate
Baseline	27.70%
“Supervised” learning	22.09%
Self-training	8.84%
Active-learning	4.02%

Table 2: Best Results of dynamic/static classification

4 Retrieving answers to dynamic questions

Following the classification step outlined in Section 3.1, we know whether a user query is static or dynamic, and the broad category of the question. If the question is dynamic, then our system performs a vertical search based on the broad topic of the question. In our system, so far, we have incorporated vertical searches on three broad topics: *Movies*, *Mass Transit*, and *Yellow Pages*.

For each broad topic, we have identified a few trusted content aggregator websites. For example, for dynamic questions related to *Movies*-related dynamic user queries, www.fandango.com is a trusted content aggregator website. Other such trusted content aggregator websites have been identified for *Mass Transit* related and for *Yellowpages* related dynamic user queries. We have also identified the web-forms that can be used to search these aggregator sites and the search parameters that these web-forms need for searching. So, given a user query, whose broad category has been determined and which has been classified as a dynamic query by the system, the next step is to parse the query to obtain pertinent search parameters.

The search parameters are dependent on the broad category of the question, the trusted content aggregator website(s), the web-forms associated with this category, and of course, the content

of the user query. From the search parameters, a search query to the associated web-form is issued to search the related aggregator site. For example, for a movie-related query, *What time is Twilight playing in Madison, New Jersey?*, the pertinent search parameters that are parsed out are *movie-name: Twilight, city: Madison, and state: New Jersey*, which are used to build a search string that Fandango’s web-form can use to search the Fandango site. For a yellow-pages type of query, *Where is the Saigon Kitchen in Austin, Texas?*, the pertinent search parameters that are parsed out are *business-name: Saigon Kitchen, city: Austin, and state: Texas*, which are used to construct a search string to search the Yellowpages website. These are just two examples of the kinds of dynamic user queries that we encounter. Within each broad category, there is a wide variety of the sub-types of user queries, and for each sub-type, we have to parse out different search parameters and use different web-forms. Details of this extraction are presented in (Feng and Bangalore, 2009).

It is quite likely that many of the dynamic queries may not have all the pertinent search parameters explicitly outlined. For example, a mass transit query may be *When is the next train to Princeton?*. The bare minimum search parameters needed to answer this query are a *from-location*, and a *to-location*. However, the *from-location* is not explicitly present in this query. In this case, the *from-location* is inferred using the GPS sensor present on the iPhone (on which our system is built to run). Depending on the web-form that we are querying, it is possible that we may be able to simply use the latitude-longitude obtained from the GPS sensor as the value for the *from-location* parameter. At other times, we may have to perform an intermediate latitude-longitude to city/state (or zip-code) conversion in order to obtain the appropriate search parameter value.

Other examples of dynamic queries in which search parameters are not explicit in the query, and hence, have to be deduced by the system, include queries such as *Where is XMen playing?* and *How long is Ace Hardware open?*. In each of these examples, the user has not specified a location. Based on our understanding of natural language, in such a scenario, our system is built to assume that the user wants to find a movie theatre (or, is referring to a hardware store) *near* where he is currently located. So, the system obtains the user’s location from the GPS sensor and uses it to search for a theatre (or locate the hardware store) within a five-mile radius of her location.

In the last few paragraphs, we have discussed how we search for answers to dynamic user queries from the hidden web by using web-forms. However, the search results returned by these web-forms usually cannot be displayed as is in our Qme! interface. The reason is that the results are often HTML pages that are designed to be displayed on a desktop or a laptop screen, not a small mobile phone screen. Displaying the results as they are returned from search would make readability difficult. So, we parse the HTML-encoded result pages to get just the answers to the user query and reformat it, to fit the Qme! interface, which is designed to be easily readable on the iPhone (as seen in Figure 1).²

5 Retrieving answers to static questions

Answers to static user queries – questions whose answers do not change over time – are retrieved in a different way than answers to dynamic questions. A description of how our system retrieves the answers to static questions is presented in this section.

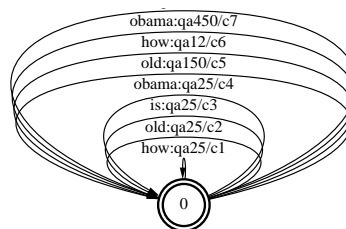


Figure 4: An example of an FST representing the search index.

5.1 Representing Search Index as an FST

To obtain results for static user queries, we have implemented our own search engine using finite-state transducers (FST), in contrast to using Lucene (Hatcher and Gospodnetic., 2004) as it is a more efficient representation of the search index that allows us to consider word lattices output by ASR as input queries.

The FST search index is built as follows. We index each question-answer (QA) pair from our repository $((q_i, a_i), qa_i$ for short) using the words (w_{q_i}) in question q_i . This index is represented as a weighted finite-state transducer (*SearchFST*) as shown in Figure 4. Here a word w_{q_i} (e.g. *old*) is the input symbol for a set of arcs whose output symbol is the index of the QA pairs where *old* appears

²We are aware that we could use SOAP (Simple Object Access Protocol) encoding to do the search, however not all aggregator sites use SOAP yet.

in the question. The weight of the arc $c_{(w_{q_i}, q_i)}$ is one of the similarity based weights discussed in Section 4.1. As can be seen from Figure 4, the words *how*, *old*, *is* and *obama* contribute a score to the question-answer pair *qa25*; while other pairs, *qa150*, *qa12*, *qa450* are scored by only one of these words.

5.2 Search Process using FSTs

A user's speech query, after speech recognition, is represented as a finite state automaton (FSA, either 1-best or WCN), *QueryFSA*. The *QueryFSA* is then transformed into another FSA (*NgramFSA*) that represents the set of n -grams of the *QueryFSA*. In contrast to most text search engines, where stop words are removed from the query, we weight the query terms with their *idf* values which results in a weighted *NgramFSA*. The *NgramFSA* is composed with the *SearchFST* and we obtain all the arcs $(w_q, qa_{w_q}, c_{(w_q, qa_{w_q})})$ where w_q is a query term, qa_{w_q} is a QA index with the query term and, $c_{(w_q, qa_{w_q})}$ is the weight associated with that pair. Using this information, we aggregate the weight for a QA pair (qa_q) across all query words and rank the retrieved QAs in the descending order of this aggregated weight. We select the top N QA pairs from this ranked list. The query composition, QA weight aggregation and selection of top N QA pairs are computed with finite-state transducer operations as shown in Equations 1 and 2³. An evaluation of this search methodology on word lattices is presented in (Mishra and Bangalore, 2010).

$$D = \pi_2(NgramFSA \circ SearchFST) \quad (1)$$

$$TopN = fsmbestpath(fsmdeinitialize(D), N) \quad (2)$$

6 Summary

In this demonstration paper, we have presented Qme!, a speech-driven question answering system for use on mobile devices. The novelty of this system is that it provides users with a single unified interface for searching both the visible and the hidden web using the most natural input modality for use on mobile phones – spoken language.

7 Acknowledgments

We would like to thank Junlan Feng, Michael Johnston and Mazin Gilbert for the help we received in putting this system together. We would

³We have dropped the need to convert the weights into the *real* semiring for aggregation, to simplify the discussion.

also like to thank ChaCha for providing us the data included in this system.

References

- M. Banko and E. Brill. 2001. Scaling to very very large corpora for natural language disambiguation. In *Proceedings of the 39th annual meeting of the association for computational linguistics: ACL 2001*, pages 26–33.
- J. Feng and S. Bangalore. 2009. Effects of word confusion networks on voice search. In *Proceedings of EACL-2009*, Athens, Greece.
- Google, 2009. <http://www.google.com/mobile>.
- P. Haffner. 2006. Scaling large margin classifiers for spoken language understanding. *Speech Communication*, 48(iv):239–261.
- E. Hatcher and O. Gospodnetic. 2004. *Lucene in Action (In Action series)*. Manning Publications Co., Greenwich, CT, USA.
- Microsoft, 2009. <http://www.live.com>.
- T. Mishra and S. Bangalore. 2010. Qme!: A speech-based question-answering system on mobile devices. In *Proceedings of NAACL-HLT*.
- D. Moldovan, C. Clark, and S. Harabagiu. 2005. Temporal context representation and reasoning. In *Proceedings of the 19th International Joint Conference on Artificial Intelligence*, pages 1009–1104.
- J. Prager, J. Chu-Carroll, and K. Czuba. 2004. Question answering using constraint satisfaction: Qa-by-dossier-with-constraints. In *Proceedings of the 42nd annual meeting of the association for computational linguistics: ACL 2004*, pages 574–581.
- J. Pustejovsky, R. Ingria, R. Saurí, J. Casta no, J. Littman, and R. Gaizauskas., 2001. *The language of time: A reader*, chapter The specification language – TimeML. Oxford University Press.
- D. Radev and B. Sundheim. 2002. Using timeml in question answering. Technical report, Brandies University.
- E. Saquete, J. L. Vicedo, P. Martínez-Barco, R. Muñoz, and H. Llorens. 2009. Enhancing qa systems with complex temporal question processing capabilities. *Journal of Artificial Intelligence Research*, 35:775–811.
- vlingo.com, 2009. <http://www.vlingomobile.com/downloads.html>.
- YellowPages, 2009. <http://www.speak4it.com>.

Tools for Multilingual Grammar-Based Translation on the Web

Aarne Ranta and Krasimir Angelov and Thomas Hallgren

Department of Computer Science and Engineering

Chalmers University of Technology and University of Gothenburg

aarne@chalmers.se, krasimir@chalmers.se, hallgren@chalmers.se

Abstract

This is a system demo for a set of tools for translating texts between multiple languages in real time with high quality. The translation works on restricted languages, and is based on semantic interlinguas. The underlying model is GF (Grammatical Framework), which is an open-source toolkit for multilingual grammar implementations. The demo will cover up to 20 parallel languages.

Two related sets of tools are presented: grammarian's tools helping to build translators for new domains and languages, and translator's tools helping to translate documents. The grammarian's tools are designed to make it easy to port the technique to new applications. The translator's tools are essential in the restricted language context, enabling the author to remain in the fragments recognized by the system.

The tools that are demonstrated will be applied and developed further in the European project MOLTO (Multilingual On-Line Translation) which has started in March 2010 and runs for three years.

1 Translation Needs for the Web

The best-known translation tools on the web are Google translate¹ and Systran². They are targeted to **consumers** of web documents: users who want to find out what a given document is about. For this purpose, **browsing quality** is sufficient, since the user has intelligence and good will, and understands that she uses the translation at her own risk.

Since Google and Systran translations can be grammatically and semantically flawed, they don't reach **publication quality**, and cannot hence be used by the **producers** of web documents. For instance, the provider of an e-commerce site cannot take the risk that the product descriptions or selling conditions have errors that change the original intentions.

There are very few automatic translation systems actually in use for producers of information. As already

noted by Bar-Hillel (1964), machine translation is one of those AI-complete tasks that involves a trade-off between coverage and precision, and the current mainstream systems opt for coverage. This is also what web users expect: they want to be able to throw just anything at the translation system and get something useful back. Precision-oriented approaches, the prime example of which is METEO (Chandioux 1977), have not been popular in recent years.

However, from the producer's point of view, large coverage is not essential: unlike the consumer's tools, their input is predictable, and can be restricted to very specific domains, and to content that the producers themselves are creating in the first place. But even in such tasks, two severe problems remain:

- The **development cost problem**: a large amount of work is needed for building translators for new domains and new languages.
- The **authoring problem**: since the method does not work for all input, the author of the source text of translation may need special training to write in a way that can be translated at all.

These two problems have probably been the main obstacles to making high-quality restricted language translation more wide-spread in tasks where it would otherwise be applicable. We address these problems by providing tools that help developers of translation systems on the one hand, and authors and translators—i.e. the users of the systems—on the other.

In the MOLTO project (Multilingual On-Line Translation)³, we have the goal to improve both the development and use of restricted language translation by an order of magnitude, as compared with the state of the art. As for development costs, this means that a system for many languages and with adequate quality can be built in a matter of days rather than months. As for authoring, this means that content production does not require the use of manuals or involve trial and error, both of which can easily make the work ten times slower than normal writing.

In the proposed system demo, we will show how some of the building blocks for MOLTO can already now be used in web-based translators, although on a

¹www.google.com/translate

²www.systransoft.com

³www.molto-project.eu

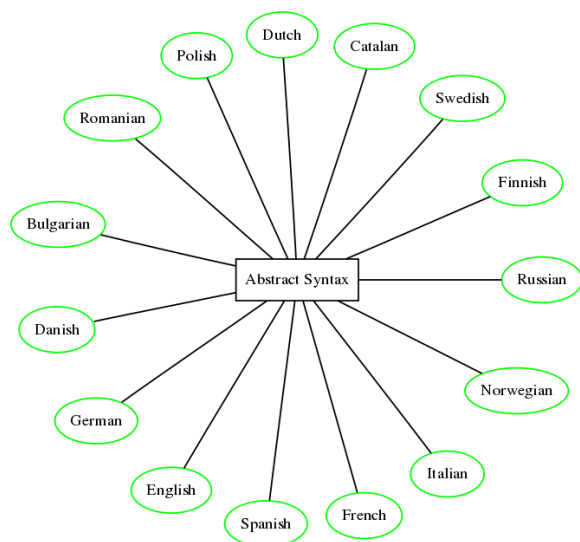


Figure 1: A multilingual GF grammar with reversible mappings from a common abstract syntax to the 15 languages currently available in the GF Resource Grammar Library.

smaller scale as regards languages and application domains. A running demo system is available at <http://grammaticalframework.org:41296>.

2 Multilingual Grammars

The translation tools are based on **GF, Grammatical Framework**⁴ (Ranta 2004). GF is a **grammar formalism**—that is, a mathematical model of natural language, equipped with a formal notation for writing grammars and a computer program implementing parsing and generation which are declaratively defined by grammars. Thus GF is comparable with formalism such as HPSG (Pollard and Sag 1994), LFG (Bresnan 1982) or TAG (Joshi 1985). The novel feature of GF is the notion of **multilingual grammars**, which describe several languages simultaneously by using a common representation called **abstract syntax**; see Figure 1.

In a multilingual GF grammar, meaning-preserving translation is provided as a composition of parsing and generation via the abstract syntax, which works as an **interlingua**. This model of translation is different from approaches based on other comparable grammar formalisms, such as synchronous TAGs (Shieber and Schabes 1990), Pargram (Butt & al. 2002, based on LFG), LINGO Matrix (Bender and Flickinger 2005, based on HPSG), and CLE (Core Language Engine, Alshawi 1992). These approaches use **transfer rules** between individual languages, separate for each pair of languages.

Being interlingua-based, GF translation scales up linearly to new languages without the quadratic blow-up of transfer-based systems. In transfer-based sys-

tems, as many as $n(n - 1)$ components (transfer functions) are needed to cover all language pairs in both directions. In an interlingua-based system, $2n + 1$ components are enough: the interlingua itself, plus translations in both directions between each language and the interlingua. However, in GF, $n + 1$ components are sufficient, because the mappings from the abstract syntax to each language (the **concrete syntaxes**) are **reversible**, i.e. usable for both generation and parsing.

Multilingual GF grammars can be seen as an implementation of Curry’s distinction between **tectogrammatical** and **phenogrammatical** structure (Curry 1961). In GF, the tectogrammatical structure is called abstract syntax, following standard computer science terminology. It is defined by using a **logical framework** (Harper & al. 1993), whose mathematical basis is in the **type theory** of Martin-Löf (1984). Two things can be noted about this architecture, both showing improvements over state-of-the-art grammar-based translation methods.

First, the translation interlingua (the abstract syntax) is a powerful logical formalism, able to express semantical structures such as context-dependencies and anaphora (Ranta 1994). In particular, dependent types make it more expressive than the type theory used in Montague grammar (Montague 1974) and employed in the Rosetta translation project (Rosetta 1998).

Second, GF uses a **framework for interlinguas**, rather than one universal interlingua. This makes the interlingual approach more light-weight and feasible than in systems assuming one universal interlingua, such as Rosetta and UNL, Universal Networking Language⁵. It also gives more precision to special-purpose translation: the interlingua of a GF translation system (i.e. the abstract syntax of a multilingual grammar) can encode precisely those structures and distinctions that are relevant for the task at hand. Thus an interlingua for mathematical proofs (Hallgren and Ranta 2000) is different from one for commands for operating an MP3 player (Perera and Ranta 2007). The expressive power of the logical framework is sufficient for both kinds of tasks.

One important source of inspiration for GF was the WYSIWYM system (Power and Scott 1998), which used domain-specific interlinguas and produced excellent quality in multilingual generation. But the generation components were hard-coded in the program, instead of being defined declaratively as in GF, and they were not usable in the direction of parsing.

3 Grammars and Ontologies

Parallel to the first development efforts of GF in the late 1990’s, another framework idea was emerging in web technology: XML, Extensible Mark-up Language, which unlike HTML is not a single mark-up language but a framework for creating custom mark-up lan-

⁴www.grammaticalframework.org

⁵www.undl.org

guages. The analogy between GF and XML was seen from the beginning, and GF was designed as a formalism for multilingual rendering of semantic content (Dymetman and al. 2000). XML originated as a format for structuring documents and structured data serialization, but a couple of its descendants, RDF(S) and OWL, developed its potential to formally express the semantics of data and content, serving as the fundamentals of the emerging Semantic Web.

Practically any meaning representation format can be converted into GF’s abstract syntax, which can then be mapped to different target languages. In particular the OWL language can be seen as a syntactic sugar for a subset of Martin-Löf’s type theory so it is trivial to embed it in GF’s abstract syntax.

The translation problem defined in terms of an ontology is radically different from the problem of translating plain text from one language to another. Many of the projects in which GF has been used involve precisely this: a meaning representation formalized as GF abstract syntax. Some projects build on previously existing meaning representation and address mathematical proofs (Hallgren and Ranta 2000), software specifications (Beckert & al. 2007), and mathematical exercises (the European project WebALT⁶). Other projects start with semantic modelling work to build meaning representations from scratch, most notably ones for dialogue systems (Perera and Ranta 2007) in the European project TALK⁷. Yet another project, and one closest to web translation, is the multilingual Wiki system presented in (Meza Moreno and Bringert 2008). In this system, users can add and modify reviews of restaurants in three languages (English, Spanish, and Swedish). Any change made in any of the languages gets automatically translated to the other languages.

To take an example, the OWL-to-GF mapping translates OWL’s classes to GF’s categories and OWL’s properties to GF’s functions that return propositions. As a running example in this and the next section, we will use the class of integers and the two-place property of being divisible (“ x is divisible by y ”). The correspondences are as follows:

```
Class(pp:integer ...)
  ⇕
cat integer
ObjectProperty(pp:div
  domain(pp:integer)
  range(pp:integer))
  ⇕
fun div :
  integer -> integer -> prop
```

4 Grammar Engineer’s Tools

In the GF setting, building a multilingual translation system is equivalent to building a multilingual GF

grammar, which in turn consists of two kinds of components:

- a language-independent abstract syntax, giving the semantic model via which translation is performed;
- for each language, a concrete syntax mapping abstract syntax trees to strings in that language.

While abstract syntax construction is an extra task compared to many other kinds of translation methods, it is technically relatively simple, and its cost is moreover amortized as the system is extended to new languages. Concrete syntax construction can be much more demanding in terms of programming skills and linguistic knowledge, due to the complexity of natural languages. This task is where GF claims perhaps the highest advantage over other approaches to special-purpose grammars. The two main assets are:

- **Programming language support:** GF is a modern functional programming language, with a powerful type system and module system supporting modular and collaborative programming and reuse of code.
- **RGL, the GF Resource Grammar Library,** implementing the basic linguistic details of languages: **inflectional morphology** and **syntactic combination functions**.

The RGL covers fifteen languages at the moment, shown in Figure 1; see also Khagai 2006, El Dada and Ranta 2007, Angelov 2008, Ranta 2009a,b, and Enache *et al.* 2010. To give an example of what the library provides, let us first consider the inflectional morphology. It is presented as a set of lexicon-building functions such as, in English,

```
mkV : Str -> V
```

i.e. function `mkV`, which takes a string (`Str`) as its argument and returns a verb (`V`) as its value. The verb is, internally, an inflection table containing all forms of a verb. The function `mkV` derives all these forms from its argument string, which is the infinitive form. It predicts all regular variations: (`mkV "walk"`) yields the purely agglutinative forms *walk-walks-walked-walked-walking* whereas (`mkV "cry"`) gives *cries-cried-cried-crying*, and so on. For irregular English verbs, RGL gives a three-argument function taking forms such as *sing,sang,sung*, but it also has a fairly complete lexicon of irregular verbs, so that the normal application programmer who builds a lexicon only needs the regular `mkV` function.

Extending a lexicon with domain-specific vocabulary is typically the main part of the work of a concrete syntax author. Considerable work has been put into RGL’s inflection functions to make them as “intelligent” as possible and thereby ease the work of the

⁶EDC-22253, webalt.math.helsinki.fi

⁷IST-507802, 2004–2006, www.talk-project.org

users of the library, who don't know the linguistic details of morphology. For instance, even Finnish, whose verbs have hundreds of forms and are conjugated in accordance with around 50 conjugations, has a one-argument function `mkV` that yields the correct inflection table for 90% of Finnish verbs.

As an example of a syntactic combination function of RGL, consider a function for predication with two-place adjectives. This function takes three arguments: a two-place adjective, a subject noun phrase, and a complement noun phrase. It returns a sentence as value:

```
pred : A2 -> NP -> NP -> S
```

This function is available in all languages of RGL, even though the details of sentence formation are vastly different in them. Thus, to give the concrete syntax of the abstract (semantical) predicate `div x y` ("x is divisible by y"), the English grammarian can write

```
div x y = pred
  (mkA2 "divisible" "by") x y
```

The German grammarian can write

```
div x y = pred
  (mkA2 "teilbar" durch_Prep) x y
```

which, even though superficially using different forms from English, generates a much more complex structure: the complement preposition `durch_Prep` takes care of rendering the argument `y` in the accusative case, and the sentence produced has three forms, as needed in grammatically different positions (*x ist teilbar durch y* in main clauses, *ist x teilbar durch y* after adverbs, and *x durch y teilbar ist* in subordinate clauses).

The syntactic combinations of the RGL have their own abstract syntax, but this abstract syntax is *not* the interlingua of translation: it is only used as a library for implementing the semantic interlingua, which is based on an ontology and abstracts away from syntactic structure. Thus the translation equivalents in a multilingual grammar need not use the same syntactic combinations in different languages. Assume, for the sake of argument, that *x is divisible by y* is expressed in Swedish by the transitive verb construction *y delar x* (literally, "y divides x"). This can be expressed easily by using the transitive verb predication function of the RGL and switching the subject and object,

```
div x y = pred (mkV2 "dela") y x
```

Thus, even though GF translation is interlingua-based, there is a component of transfer between English and Swedish. But this transfer is performed at compile time. In general, the use of the large-coverage RGL as a library for restricted grammars is called **grammar specialization**. The way GF performs grammar specialization is based on techniques for optimizing functional programming languages, in particular **partial evaluation** (Ranta 2004, 2007). GF also gives a possibility to run-time transfer via semantic actions on abstract syntax trees, but this option has rarely been needed in previous applications, which helps to keep translation systems simple and efficient.

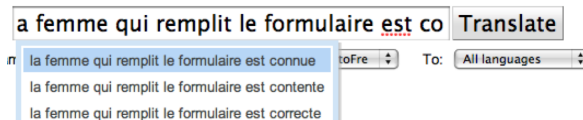


Figure 2: French word prediction in GF parser, suggesting feminine adjectives that agree with the subject *la femme*.

As shown in Figure 1, the RGL is currently available for 15 languages, of which 12 are official languages of the European Union. A similar number of new languages are under construction in this collaborative open-source project. Implementing a new language is an effort of 3–6 person months.

5 Translator's Tools

For the translator's tools, there are three different use cases:

- restricted source
 - production of source in the first place
 - modifying source produced earlier
- unrestricted source

Working with restricted source language recognizable by a GF grammar is straightforward for the translating tool to cope with, except when there is ambiguity in the text. The real challenge is to help the author to keep inside the restricted language. This help is provided by **predictive parsing**, a technique recently developed for GF (Angelov 2009)⁸. Incremental parsing yields **word predictions**, which guide the author in a way similar to the T9 method⁹ in mobile phones. The difference from T9 is that GF's word prediction is sensitive to the grammatical context. Thus it does not suggest all existing words, but only those words that are grammatically correct in the context. Figure 2 shows an example of the parser at work. The author has started a sentence as *la femme qui remplit le formulaire est co* ("the woman who fills the form is *co*"), and a menu shows a list of words beginning with *co* that are given in the French grammar and possible in the context at hand; all these words are adjectives in the feminine form. Notice that the very example shown in Figure 2 is one that is difficult for n-gram-based statistical translators: the adjective is so far from the subject with which it agrees that it cannot easily be related to it.

Predictive parsing is a good way to help users produce translatable content in the first place. When modifying the content later, e.g. in a wiki, it may not be optimal, in particular if the text is long. The text can

⁸ Parsing in GF is polynomial with an arbitrary exponent in the worst case, but, as shown in Angelov 2009, linear in practice with realistic grammars.

⁹ www.t9.com

```

Pred known_A (Rel woman_N (Compl
fill_V2 form_N))
the woman who fills the form is known
la femme qui remplit le formulaire est connue
—→
Pred known_A (Rel man_N (Compl fill_V2
form_N))
the man who fills the form is known
l' homme qui remplit le formulaire est connu

```

Figure 3: Change in one word (boldface) propagated to other words depending on it (italics).

contain parts that depend on each other but are located far apart. For instance, if the word *femme* (“woman”) in the previous example is changed to *homme*, the preceding article *la* has to be changed to *l’*, and the adjective has to be changed to the masculine form: thus *connue* (“known”) would become *connu*, and so on. Such changes are notoriously difficult even for human authors and translators, and can easily leave a document in an inconsistent state. This is where another utility of the abstract syntax comes in: in the abstract syntax tree, all that is changed is the noun, and the regenerated concrete syntax string automatically obeys all the agreement rules. The process is shown in Figure 3. The one-word change generating the new set of documents can be performed by editing any of the three representations: the tree, the English version, or the French version. This functionality is implemented in the GF **syntax editor** (Khegai & al. 2003).

Restricted languages in the sense of GF are close to **controlled languages**, such as Attempto (Fuchs & al. 2008); the examples shown in this section are actually taken from a GF implementation that generalizes Attempto Controlled English to five languages (Angelov and Ranta 2009). However, unlike typical controlled languages, GF does not require the absence of ambiguity. In fact, when a controlled language is generalized to new languages, lexical ambiguities in particular are hard to avoid.

The predictive parser of GF does not try to resolve ambiguities, but simply returns all alternatives in the parse chart. If the target language has exactly the same ambiguity, it remains hidden in the translation. But if the ambiguity does make a difference in translation, it has to be resolved, and the system has to provide a possibility of manual disambiguation by the user to guarantee high quality.

The translation tool snapshot in Figure 2 is from an actual web-based prototype. It shows a slot in an HTML page, built by using JavaScript via the Google Web Toolkit (Bringert & al. 2009). The translation is performed using GF in a server, which is called via HTTP. Also client-side translators, with similar user interfaces, can be built by converting the whole GF grammar to JavaScript (Meza Moreno and Bringert 2008).

6 The Demo

In the demo, we will show

- how a simple translation system is built and compiled by using the GF grammar compiler and the resource grammar library
- how the translator is integrated in a web page
- how the translator is used in a web browser by means of an integrated incremental parser

A preliminary demo can be seen in <http://grammaticalframework.org:41296>. All the demonstrated tools are available as open-source software from <http://grammaticalframework.org>.

The work reported here is supported by MOLTO (Multilingual On-Line Translation. FP7-ICT-247914).

References

- Alshawi, H. (1992). *The Core Language Engine*. Cambridge, Ma: MIT Press.
- Angelov, K. (2008). Type-Theoretical Bulgarian Grammar. In B. Nordström and A. Ranta (Eds.), *Advances in Natural Language Processing (GOTAL 2008)*, Volume 5221 of *LNCS/LNAI*, pp. 52–64. URL <http://www.springerlink.com/content/978-3-540-85286-5/>.
- Angelov, K. (2009). Incremental Parsing with Parallel Multiple Context-Free Grammars. In *Proceedings of EACL’09, Athens*.
- Angelov, K. and A. Ranta (2009). Implementing Controlled Languages in GF. In *Proceedings of CNL-2009, Marettimo*, LNCS. to appear.
- Bar-Hillel, Y. (1964). *Language and Information*. Reading, MA: Addison-Wesley.
- Beckert, B., R. Hähnle, and P. H. Schmitt (Eds.) (2007). *Verification of Object-Oriented Software: The KeY Approach*. LNCS 4334. Springer-Verlag.
- Bender, E. M. and D. Flickinger (2005). Rapid prototyping of scalable grammars: Towards modularity in extensions to a language-independent core. In *Proceedings of the 2nd International Joint Conference on Natural Language Processing IJCNLP-05 (Posters/Demos)*, Jeju Island, Korea. URL <http://faculty.washington.edu/ebender/papers/modules05.pdf>.
- Bresnan, J. (Ed.) (1982). *The Mental Representation of Grammatical Relations*. MIT Press.
- Bringert, B., K. Angelov, and A. Ranta (2009). Grammatical Framework Web Service. In *Proceedings of EACL’09, Athens*.

- Butt, M., H. Dyvik, T. H. King, H. Masuichi, and C. Rohrer (2002). The Parallel Grammar Project. In *COLING 2002, Workshop on Grammar Engineering and Evaluation*, pp. 1–7. URL <http://www2.parc.com/isl/groups/nltp/pargram/buttetal-coling02.pdf>.
- Chandioux, J. (1976). MÉTÉO: un système opérationnel pour la traduction automatique des bulletins météorologiques destinés au grand public. *META* 21, 127–133.
- Curry, H. B. (1961). Some logical aspects of grammatical structure. In R. Jakobson (Ed.), *Structure of Language and its Mathematical Aspects: Proceedings of the Twelfth Symposium in Applied Mathematics*, pp. 56–68. American Mathematical Society.
- Dada, A. E. and A. Ranta (2007). Implementing an Open Source Arabic Resource Grammar in GF. In M. Mughazy (Ed.), *Perspectives on Arabic Linguistics XX*, pp. 209–232. John Benjamin's.
- Dean, M. and G. Schreiber (2004). OWL Web Ontology Language Reference. URL <http://www.w3.org/TR/owl-ref/>.
- Dymetman, M., V. Lux, and A. Ranta (2000). XML and multilingual document authoring: Convergent trends. In *COLING, Saarbrücken, Germany*, pp. 243–249. URL <http://www.cs.chalmers.se/~aarne/articles/coling2000.ps.gz>.
- Enache, R., A. Ranta, and K. Angelov (2010). An Open-Source Computational Grammar for Romanian. In A. Gelbukh (Ed.), *Intelligent Text Processing and Computational Linguistics (CICLing-2010), Iasi, Romania, March 2010, LNCS*, to appear.
- Fuchs, N. E., K. Kaljurand, and T. Kuhn (2008). Attempto Controlled English for Knowledge Representation. In C. Baroglio, P. A. Bonatti, J. Małuszyński, M. Marchiori, A. Polleres, and S. Schaffert (Eds.), *Reasoning Web, Fourth International Summer School 2008*, Number 5224 in Lecture Notes in Computer Science, pp. 104–124. Springer.
- Hallgren, T. and A. Ranta (2000). An extensible proof text editor. In M. Parigot and A. Voronkov (Eds.), *LPAR-2000*, Volume 1955 of *LNCS/LNAI*, pp. 70–84. Springer. URL <http://www.cs.chalmers.se/~aarne/articles/lpar2000.ps.gz>.
- Harper, R., F. Honsell, and G. Plotkin (1993). A Framework for Defining Logics. *JACM* 40(1), 143–184.
- Joshi, A. (1985). Tree-adjointing grammars: How much context-sensitivity is required to provide reasonable structural descriptions. In D. Dowty, L. Karttunen, and A. Zwicky (Eds.), *Natural Language Parsing*, pp. 206–250. Cambridge University Press.
- Khegai, J. (2006). GF Parallel Resource Grammars and Russian. In *Coling/ACL 2006*, pp. 475–482.
- Khegai, J., B. Nordström, and A. Ranta (2003). Multilingual Syntax Editing in GF. In A. Gelbukh (Ed.), *Intelligent Text Processing and Computational Linguistics (CICLing-2003), Mexico City, February 2003*, Volume 2588 of *LNCS*, pp. 453–464. Springer-Verlag. URL <http://www.cs.chalmers.se/~aarne/articles/mexico.ps.gz>.
- Martin-Löf, P. (1984). *Intuitionistic Type Theory*. Napoli: Bibliopolis.
- Meza Moreno, M. S. and B. Bringert (2008). Interactive Multilingual Web Applications with Grammatical Framework. In B. Nordström and A. Ranta (Eds.), *Advances in Natural Language Processing (GoTAL 2008)*, Volume 5221 of *LNCS/LNAI*, pp. 336–347. URL <http://www.springerlink.com/content/978-3-540-85286-5/>.
- Montague, R. (1974). *Formal Philosophy*. New Haven: Yale University Press. Collected papers edited by Richmond Thomason.
- Perera, N. and A. Ranta (2007). Dialogue System Localization with the GF Resource Grammar Library. In *SPEECHGRAM 2007: ACL Workshop on Grammar-Based Approaches to Spoken Language Processing, June 29, 2007, Prague*. URL <http://www.cs.chalmers.se/~aarne/articles/perera-ranta.pdf>.
- Pollard, C. and I. Sag (1994). *Head-Driven Phrase Structure Grammar*. University of Chicago Press.
- Power, R. and D. Scott (1998). Multilingual authoring using feedback texts. In *COLING-ACL*.
- Ranta, A. (1994). *Type Theoretical Grammar*. Oxford University Press.
- Ranta, A. (2004). Grammatical Framework: A Type-Theoretical Grammar Formalism. *The Journal of Functional Programming* 14(2), 145–189. URL <http://www.cs.chalmers.se/~aarne/articles/gf-jfp.ps.gz>.
- Ranta, A. (2009a). Grammars as Software Libraries. In Y. Bertot, G. Huet, J.-J. Lévy, and G. Plotkin (Eds.), *From Semantics to Computer Science*. Cambridge University Press. URL <http://www.cs.chalmers.se/~aarne/articles/libraries-kahn.pdf>.
- Ranta, A. (2009b). The GF Resource Grammar Library. In *Linguistic Issues in Language Technology*, Vol. 2. URL <http://elanguage.net/journals/index.php/lilt/article/viewFile/214/158>.
- Rosetta, M. T. (1994). *Compositional Translation*. Dordrecht: Kluwer.
- Shieber, S. M. and Y. Schabes (1990). Synchronous tree-adjointing grammars. In *COLING*, pp. 253–258.

Demonstration of a prototype for a Conversational Companion for reminiscing about images

Yorick Wilks

IHMC, Florida
ywilks@ihmc.us

Alexiei Dingli

University of Malta, Malta
alexiei.dingli@um.edu.mt

Roberta Catizone

University of Sheffield, UK
r.catizone@dcs.shef.ac.uk

Weiwei Cheng

University of Sheffield, UK
w.cheng@dcs.shef.ac.uk

Abstract

This paper describes an initial prototype demonstrator of a Companion, designed as a platform for novel approaches to the following: 1) The use of Information Extraction (IE) techniques to extract the content of incoming dialogue utterances after an Automatic Speech Recognition (ASR) phase, 2) The conversion of the input to Resource Descriptor Format (RDF) to allow the generation of new facts from existing ones, under the control of a Dialogue Manger (DM), that also has access to stored knowledge and to open knowledge accessed in real time from the web, all in RDF form, 3) A DM implemented as a stack and network virtual machine that models mixed initiative in dialogue control, and 4) A tuned dialogue act detector based on corpus evidence. The prototype platform was evaluated, and we describe this briefly; it is also designed to support more extensive forms of emotion detection carried by both speech and lexical content, as well as extended forms of machine learning.

1. Introduction

This demonstrator Senior Companion (SC) was built during the initial phase of the Companions project and aims to change the way we think about the relationships of people to computers and the internet by developing a virtual conversational 'Companion that will be an agent or 'presence' that stays with the user for long periods of time, developing a relationship and 'knowing its owners' preferences and wishes. The Companion communicates with the user primarily through speech, but also using other technologies such as touch screens and sensors.

This paper describes the functionality and system modules of the Senior Companion, one of two initial prototypes built in the first two years of the project. The SC provides a multimodal interface for eliciting, retrieving and inferring personal information from elderly users by means of conversation about their photographs. The Companion, through conversation, elicits life memo-

ries and reminiscences, often prompted by discussion of their photographs; the aim is that the Companion should come to know a great deal about its user, their tastes, likes, dislikes, emotional reactions etc, through long periods of conversation. It is assumed that most life information will soon be stored on the internet (as in the Memories for Life project: <http://www.memoriesforlife.org/>) and we have linked the SC directly to photo inventories in Facebook (see below). The overall aim of the SC project (not yet achieved) is to produce a coherent life narrative for its user from conversations about personal photos, although its short-term goals, reported here, are to assist, amuse and entertain the user.

The technical content of the project is to use a number of types of machine learning (ML) to achieve these ends in original ways, initially using a methodology developed in earlier research: first, by means of an Information Extraction (IE) approach to deriving content from user input utterances; secondly, using a training method for attaching Dialogue Acts to these utterance and, lastly, using a specific type of dialogue manager (DM) that uses Dialogue Action Forms (DAF) to determine the context of any utterance. A stack of these DAFs is the virtual machine that models the ongoing dialogue by means of shared user and Companion initiative and generates appropriate responses. In this description of the demo, we shall:

- describe the current SC prototype's functionality;
- set out its architecture and modules, focusing on the Natural Language Understanding module and the Dialogue Manager.

A mini-version of the demo running in real time can be seen at: [URL](#)

<http://www.youtube.com/watch?v=-Xx5hgjD-Mw>

2. The Senior Companion System

The Senior Companion prototype (Wilks, 2010) was designed to make a rapid advance in the first two years of the project so as to be basis for a second round of prototypes embodying more advanced ML. This strategy was deliberately chosen to avoid a well-known problem with experimental AI systems: that a whole project is spent in design so that a prototype never emerges until the very end, but which is then never fully evaluated and, most importantly, nothing is ever built on the experience obtained in its construction. The central function of the SC is engaging the user in discussion about their photographs: where and when they were taken, details about the people in them and their relationship to the user and each other.

Once a photo is loaded, it is processed with face recognition software to identify any faces in it. The recognition software, OpenCV, provides positional information by identifying the face coordinates and this information is exploited in the Dialogue Manager by making explicit reference to the position of people in the photograph (the person on the left, right, center, etc.) as well as recognizing when there are groups of people. The system discusses properties of the photo as well as properties and relationships of the people in the photos.

The news reading feature adds an interesting accompaniment to the photo domain and demonstrates the ability of the system to handle more than one kind of application at a time, and news has, of course, an unconstrained vocabulary. The following is a fairly typical example of its current capacity, depending of course on the images loaded, and comes from the middle part of a sample dialogue generated by the system when discussing a group photo from a vacation in Zanzibar, and after the system had already stored, from an earlier conversation, that the user has a daughter Zoe (using researchers' children to avoid any problems):

SC: What is your relationship to the doctor?
User: There is no relationship.
SC: Please tell me how you know the doctor
U: Uh, he's a friend.
SC: How old was the doctor when this photo was taken?
U: I'm not sure but I think about 30.

SC: How do you know Octavia?
U: She's my daughter.
SC: Aha! I see. Octavia's sister is Zoe.
SC: How old was your daughter Octavia when this photo was taken?
U: Octavia was 19 years old.
SC: What was the occasion of this photo?
U: This was another family holiday.
SC: Where was this photo taken?
U: This photo was taken in Tanzania.
SC: have you ever been to Zanzibar?
U: Yes, I have.
SC: Great place, I hope you saw the beaches. When was this photo taken?
U: last year.

It is a key feature of the SC that it can make simple inferences from family relationships it already knows (e.g. that daughters of the same parent are siblings) and that it can access real-time information about places to show that it has some knowledge of what is being talked about, in this case the beaches on Zanzibar, and how this is done is discussed below. This real time access to unconstrained place information on the internet is an attempt to break out of classic AI systems that only know the budget of facts they have been primed with.

This basic system provides the components for future development of the SC, as well as its main use as a device to generate more conversation data for machine learning research in the future. Key features of the SC are listed below followed by a description of the system architecture and modules. The SC:

- Contains a visually appealing multi-modal interface with a character avatar to mediate the system's functionality to the user.
- Interacts with the user using multiple modalities – speech and touch.
- Includes face detection software for identifying the position of faces in the photos.
- Accepts pre-annotated (XML) photo inventories as a means for creating richer dialogues more quickly.
- Engages in conversation with the user about topics within the photo domain: when and where the photo was taken, discussion of the people in the photo including their relationships to the user.
- Reads news from three categories: politics, business and sports.

- Tells jokes taken from an internet-based joke website.
- Retains all user input for reference in repeat user sessions, in addition to the knowledge base that has been updated by the Dialogue Manager on the basis of what was said.
- Contains a fully integrated Knowledge Base for maintaining user information including:
 - Ontological information which is exploited by the Dialogue Manager and provides domain-specific relations between fundamental concepts.
 - A mechanism for storing information in a triple store (Subject-Predicate-Object) - the RDF Semantic Web format - for handling unexpected user input that falls outside of the photo domain, e.g. arbitrary locations in which photos might have been taken.
 - A reasoning module for reasoning over the Knowledge Base and world knowledge obtained in RDF format from the internet; the SC is thus a primitive Semantic Web device (see reference8, 2008)
- Contains basic photo management capability allowing the user, in conversation, to select photos as well as display a set of photos with a particular feature.

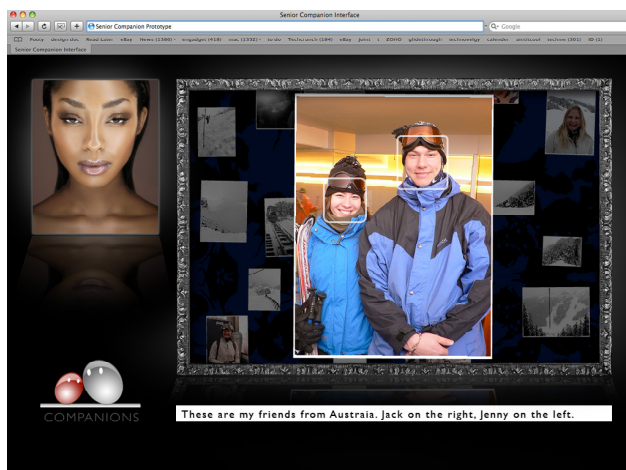


Figure 1: The Senior Companion Interface

3. System Architecture

In this section we will review the components of the SC architecture. As can be seen from Figure 2, the architecture contains three abstract level components – Connectors, Input Handlers and Application Services –together with the Dialogue Manager and the Natural Language Understander (NLU).

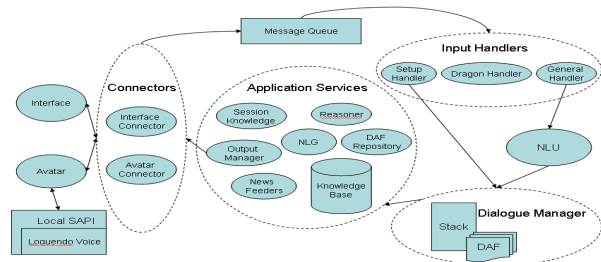


Figure 2: Senior Companion system architecture

Connectors form a communication bridge between the core system and external applications. The external application refers to any modules or systems which provide a specific set of functionalities that might be changed in the future. There is one connector for each external application. It hides the underlying complex communication protocol details and provides a general interface for the main system to use. This abstraction decouples the connection of external and internal modules and makes changing and adding new external modules easier. At this moment, there are two connectors in the system – Napier Interface Connector and CrazyTalk Avatar Connector. Both of them are using network sockets to send/receive messages.

Input Handlers are a set of modules for processing messages according to message types. Each handler deals with a category of messages where categories are coarse-grained and could include one or more message types. The handlers separate the code handling inputs into different places and make the code easier to locate and change. Three handlers have been implemented in the Senior Companion system – Setup Handler, Dragon Events Handler and General Handler. The Setup Handler is responsible for loading the photo annotations if any, performing face detection if no annotation file is associated with the photo and checking the Knowledge Base in case

the photo being processed has been discussed in earlier sessions. Dragon Event Handler deals with dragon speech recognition commands sent from the interface while the General Handler processes user utterances and photo change events of the interface.

Application Services are a group of internal modules which provide interfaces for the Dialogue Action Forms (DAF) to use. It has an easy-to-use high-level interface for general DAF designers to code associated tests and actions as well as a low level interface for advanced DAFs. It also provides the communication link between DAFs and the internal system and enables DAFs to access system functionalities. Following is a brief summary of modules grouped into Application Services.

News Feeders are a set of RSS Feeders for fetching news from the internet. Three different news feeders have been implemented for fetching news from BBC website Sports, Politics and Business channels. There is also a Jokes Feeder to fetch Jokes from internet in a similar way. During the conversation, the user can request news about particular topics and the SC simply reads the news downloaded through the feeds.

The DAF Repository is a list of DAFs loaded from files generated by the DAF Editor.

The Natural Language Generation (NLG) module is responsible for randomly selecting a system utterance from a template. An optional variable can be passed when calling methods on this module. The variable will be used to replace special symbols in the text template if applicable.

Session Knowledge is the place where global information for a particular running session is stored. For example, the name of the user who is running the session, the list of photos being discussed in this session and the list of user utterances etc.

The Knowledge Base is the data store of persistent knowledge. It is implemented as an RDF triplestore using a Jena implementation. The triplestore API is a layer built upon a traditional relational database. The application can save/retrieve information as RDF triples rather than table records. The structure of knowledge represented in RDF triples is discussed later.

The Reasoner is used to perform inference on existing knowledge in the Knowledge Base (see example in next section).

The Output Manager deals with sending messages to external applications. It has been implemented in a publisher/subscriber fashion. There are three different channels in the system: the text channel, the interface command channel and the avatar command channel. Those channels could be subscribed to by any connectors and handled respectively.

4. Dialogue understanding and inference

Every utterance is passed through the Natural Language Understanding (NLU) module for processing. This module uses a set of well-established natural language processing tools such as those found in the GATE (Cunningham, et al., 1997) system. The basic processes carried out by GATE are: tokenizing, sentence splitting, POS tagging, parsing and Named Entity Recognition. These components have been further enhanced for the SC system by adding 1) new and improved gazetteers including family relations and 2) accompanying extraction rules. The Named Entity (NE) recognizer is a key part of the NLU module and recognizes the significant entities required to process dialogue in the photo domain: PERSON NAMES, LOCATION NAMES, FAMILY RELATIONS and DATES. Although GATE recognizes basic entities, more complex entities are not handled. Apart from the gazetteers mentioned earlier and the hundreds of extraction rules already present in GATE, about 20 new extraction rules using the JAPE rule language were also developed for the SC module. These included rules which identify complex dates, family relationships, negations and other information related to the SC domain. The following is an example of a simple rule used to identify relationship in utterances such as “Mary is my sister”:

```
Macro: RELATIONSHIP_IDENTIFIER
(
  ({{Token.category=="PRP$"}}{{Token.category=="PRP"}}{{Lookup.majorType=="person_first"}}):person2
  ({{Token.string=="is"}})
  ({{Token.string=="my"}}):person1
  ({{Lookup.minorType=="Relationship"}}):relationship)

```


Using this rule with the example mentioned earlier, the rule interprets person1 as referring to the speaker so, if the name of the user speaking is John (which was known from previous conversations), it is utilized. Person 2 is then the name of the person mentioned, i.e. Mary. This name is recognised by using the gazetteers we have in the system (which contain about 40,000 first names). The relationship is once again identified using the almost 800 unique relationships added to the gazetteer. With this information, the NLU module identifies Information Extraction patterns in the dialogue that represent significant content with respect to a user's life and photos.

The information obtained (such as Mary=sister-of John) is passed to the Dialogue Manager (DM) and then stored in the knowledge base (KB). The DM filters what to include and exclude from the KB. Given, in the example above, that Mary is the sister of John, the NLU knows that sister is a relationship between two people and is a key relationship. However, the NLU also discovers syntactical information such as the fact the both Mary and John are nouns. Even though this information is important, it is too low level to be of any use by the SC with respect to the user, i.e. the user is not interested in the parts-of-speech of a word. Thus, this information is discarded by the DM and not stored in the KB. The NLU module also identifies a Dialogue Act Tag for each user utterance based on the DAMSL set of DA tags and prior work done jointly with the University of Albany (Webb et al., 2008).

The KB is a long-term store of information which makes it possible for the SC to retrieve information stored between different sessions. The information can be accessed anytime it is needed by simply invoking the relevant calls. The structure of the data in the database is an RDF triple, and the KB is more commonly referred to as a triple store. In mathematical terms, a triple store is nothing more than a large database of interconnected graphs. Each triple is made up of a subject, a predicate and an object. So, if we took the previous example, Mary sister-of John; Mary would be the subject, sister-of would be the predicate and John would be the object. The inference engine is an important part of the system because it allows us to discover new facts beyond what is elicited from the conversation with the user.

Uncle Inference Rule:
 (?a sisterOf ?b),
 (?x sonOf ?a),
 (?b gender male) -> (?b uncleOf ?x)

Triples:
 (Mary sisterOf John)
 (Tom sonOf Mary)

Triples produced automatically by ANNIE (the semantic tagger):
 (John gender male)

Inference:
 (Mary sisterOf John)
 (Tom sonOf Mary)
 (John gender male)
 ->
 (John uncleOf Tom)

This kind of inference is already used by the SC and we have about 50 inference rules aimed at producing new data on the relationships domain. This combination of triple store, inference engine and inference rules makes a system which is weak but powerful enough to mimic human reasoning in this domain and thus simulate basic intelligence in the SC. For our prototype, we are using the JENA Semantic Web Framework for the inference engine together with a MySQL database as the knowledge base. However, this system of family relationships is not enough to cover all the possible topics which can crop up during a conversation and, in such circumstances, the DM switches to an open-world model and instructs the NLU to seek further information online.

5. The Hybrid-world approach

When the DM requests further information on a particular topic, the NLU first checks with the KB whether the topic is about something known. At this stage, we have to keep in mind that any topic requested by the DM should be already in the KB since it was preprocessed by the NLU when it was mentioned in the utterance. So, if the user informs the system that the photograph was taken in Paris, (in response to a system question asking where the photo was taken), the utterance is first processed by the NLU which discovers that "Paris" is a location using its semantic tagger ANNIE (A Nearly New Information Extraction engine). The semantic tagger makes use of gazetteers and IE rules in order to accomplish

this task. It also goes through the KB and retrieves any triples related to “Paris”. Inference is then performed on this data and the new information generated by this process is stored back in the KB.

Once the type of information is identified, the NLU can use various predefined strategies: In the case of LOCATIONS, one of the strategies used is to seek for information in Wiki-Travel or Virtual Tourists. The system already knows how to query these sites and interpret their output by using predefined wrappers. This is then used to extract relevant information from the mentioned sites webpages by sending an online query to these sites and storing the information retrieved in the triple-store. This information is then used by the DM to generate a reply. In the previous example, the system manages to extract the best sightseeing spots in Paris. The NLU would then store in the KB triples such as [Paris, sightseeing, Eiffel Tower] and the DM with the help of the NLG would ask the user “I’ve heard that the X is a very famous spot. Have you seen it while you were there?” Obviously in this case, X would be replaced by the “Eiffel Tower”.

On the other hand, if the topic requested by the DM is unknown, or the semantic tagger is not capable of understanding the semantic category, the system uses a normal search engine (and this is what we call “hybrid-world”: the move outside the world the system already knows). A query containing the unknown term in context is sent to standard engines and the top pages are retrieved. These pages are then processed using ANNIE and their tagged attributes are analyzed. The standard attributes returned by ANNIE include information about Dialogue Acts, Polarity (i.e. whether a sentence has positive, negative or neutral connotations), Named Entities, Semantic Categories (such as dates and currency), etc. The system then filters the information collected by using more generic patterns and generates a reply from the resultant information. ANNIE’s polarity methods have been shown to be an adequate implementation of the general word-based polarity methods pioneered by Wiebe and her colleagues (see e.g. Akkaya et al., 2009).

6. Evaluation

The notion of companionship is not yet one with any agreed evaluation strategy or metric, though developing one is part of the main project itself.

Again, there are established measures for the assessment of dialogue programs but they have all been developed for standard task-based dialogues and the SC is not of that type: there is no specific task either in reminiscing conversations, nor in the elicitation of the content of photos, that can be assessed in standard ways, since there is no clear point at which an informal dialogue need stop, having been completed. Conventional dialogue evaluations often use measures like “stickiness” to determine how much a user will stay with or stick with a dialogue system and not leave it, presumably because they are disappointed or find it lacking in some feature. But it is hard to separate that feature out from a task rapidly and effectively completed, where stickiness would be low not high. Traum (Traum et al., 2004) has developed a methodology for dialogue evaluation based on “appropriateness” of responses and the Companions project has developed a model of evaluation for the SC based on that (Benyon et al., 2008).

Acknowledgement

This work was funded by the Companions project (2006-2009) sponsored by the European Commission as part of the Information Society Technologies (IST) programme under EC grant number IST-FP6-034434.

References

- David Benyon, Prem Hansen and Nick Webb, 2008. Evaluating Human-Computer Conversation in Companions. In: *Proc. 4th International Workshop on Human-Computer Conversation*, Bellagio, Italy.
- Cem Akkaya, Jan Wiebe, and Rada Mihalcea, 2009. Subjectivity Word Sense Disambiguation, In: *EMNLP 2009*.
- Hamish Cunningham, Kevin Humphreys, Robert Gaizauskas, and Yorick Wilks, 1997. GATE -- a TIPSTER based General Architecture for Text Engineering. In: *Proceedings of the TIPSTER Text Program (Phase III) 6 Month Workshop*. Morgan Kaufmann, CA.
- David Traum, Susan Robinson, and Jens Stephan. 2004. Evaluation of multi-party virtual reality dialogue interaction, In: *Proceedings of Fourth International Conference on Language Resources and Evaluation (LREC 2004)*, pp.1699-1702
- Yorick Wilks (ed.) 2010. *Artificial Companions in Society: scientific, economic, psychological and philosophical perspectives*. John Benjamins: Amsterdam.

It Makes Sense: A Wide-Coverage Word Sense Disambiguation System for Free Text

Zhi Zhong and Hwee Tou Ng
Department of Computer Science
National University of Singapore

13 Computing Drive
Singapore 117417

{zhongzhi, nght}@comp.nus.edu.sg

Abstract

Word sense disambiguation (WSD) systems based on supervised learning achieved the best performance in SenseEval and SemEval workshops. However, there are few publicly available open source WSD systems. This limits the use of WSD in other applications, especially for researchers whose research interests are not in WSD.

In this paper, we present IMS, a supervised English all-words WSD system. The flexible framework of IMS allows users to integrate different preprocessing tools, additional features, and different classifiers. By default, we use linear support vector machines as the classifier with multiple knowledge-based features. In our implementation, IMS achieves state-of-the-art results on several SenseEval and SemEval tasks.

1 Introduction

Word sense disambiguation (WSD) refers to the task of identifying the correct sense of an ambiguous word in a given context. As a fundamental task in natural language processing (NLP), WSD can benefit applications such as machine translation (Chan et al., 2007a; Carpuat and Wu, 2007) and information retrieval (Stokoe et al., 2003).

In previous SenseEval workshops, the supervised learning approach has proven to be the most successful WSD approach (Palmer et al., 2001; Snyder and Palmer, 2004; Pradhan et al., 2007). In the most recent SemEval-2007 English all-words tasks, most of the top systems were based on supervised learning methods. These systems used a set of knowledge sources drawn from sense-annotated data, and achieved significant improvements over the baselines.

However, developing such a system requires much effort. As a result, very few open source WSD systems are publicly available – the only other publicly available WSD system that we are aware of is SenseLearner (Mihalcea and Csomai, 2005). Therefore, for applications which employ WSD as a component, researchers can only make use of some baselines or unsupervised methods. An open source supervised WSD system will promote the use of WSD in other applications.

In this paper, we present an English all-words WSD system, IMS (It Makes Sense), built using a supervised learning approach. IMS is a Java implementation, which provides an extensible and flexible platform for researchers interested in using a WSD component. Users can choose different tools to perform preprocessing, such as trying out various features in the feature extraction step, and applying different machine learning methods or toolkits in the classification step. Following Lee and Ng (2002), we adopt support vector machines (SVM) as the classifier and integrate multiple knowledge sources including parts-of-speech (POS), surrounding words, and local collocations as features. We also provide classification models trained with examples collected from parallel texts, SEMCOR (Miller et al., 1994), and the DSO corpus (Ng and Lee, 1996).

A previous implementation of the IMS system, NUS-PT (Chan et al., 2007b), participated in SemEval-2007 English all-words tasks and ranked first and second in the coarse-grained and fine-grained task, respectively. Our current IMS implementation achieves competitive accuracies on several SenseEval/SemEval English lexical-sample and all-words tasks.

The remainder of this paper is organized as follows. Section 2 gives the system description, which introduces the system framework and the details of the implementation. In Section 3, we present the evaluation results of IMS on SenseE-

val/SemEval English tasks. Finally, we conclude in Section 4.

2 System Description

In this section, we first outline the IMS system, and introduce the default preprocessing tools, the feature types, and the machine learning method used in our implementation. Then we briefly explain the collection of training data for content words.

2.1 System Architecture

Figure 1 shows the system architecture of IMS. The system accepts any input text. For each content word w (noun, verb, adjective, or adverb) in the input text, IMS disambiguates the sense of w and outputs a list of the senses of w , where each sense s_i is assigned a probability according to the likelihood of s_i appearing in that context. The sense inventory used is based on WordNet (Miller, 1990) version 1.7.1.

IMS consists of three independent modules: preprocessing, feature and instance extraction, and classification. Knowledge sources are generated from input texts in the preprocessing step. With these knowledge sources, instances together with their features are extracted in the instance and feature extraction step. Then we train one classification model for each word type. The model will be used to classify test instances of the corresponding word type.

2.1.1 Preprocessing

Preprocessing is the step to convert input texts into formatted information. Users can integrate different tools in this step. These tools are applied on the input texts to extract knowledge sources such as sentence boundaries, part-of-speech tags, etc. The extracted knowledge sources are stored for use in the later steps.

In IMS, preprocessing is carried out in four steps:

- Detect the sentence boundaries in a raw input text with a sentence splitter.
- Tokenize the split sentences with a tokenizer.
- Assign POS tags to all tokens with a POS tagger.
- Find the lemma form of each token with a lemmatizer.

By default, the sentence splitter and POS tagger in the OpenNLP toolkit¹ are used for sentence splitting and POS tagging. A Java version of Penn TreeBank tokenizer² is applied in tokenization. JWNL³, a Java API for accessing the WordNet (Miller, 1990) thesaurus, is used to find the lemma form of each token.

2.1.2 Feature and Instance Extraction

After gathering the formatted information in the preprocessing step, we use an instance extractor together with a list of feature extractors to extract the instances and their associated features.

Previous research has found that combining multiple knowledge sources achieves high WSD accuracy (Ng and Lee, 1996; Lee and Ng, 2002; Decadt et al., 2004). In IMS, we follow Lee and Ng (2002) and combine three knowledge sources for all content word types⁴:

- *POS Tags of Surrounding Words* We use the POS tags of three words to the left and three words to the right of the target ambiguous word, and the target word itself. The POS tag feature cannot cross sentence boundary, which means all the associated surrounding words should be in the same sentence as the target word. If a word crosses sentence boundary, the corresponding POS tag value will be assigned as *null*.

For example, suppose we want to disambiguate the word *interest* in a POS-tagged sentence “My/PRP\$ brother/NN has/VBZ always/RB taken/VBN a/DT keen/JJ interest/NN in/IN my/PRP\$ work/NN ./.”. The 7 POS tag features for this instance are $\langle VBN, DT, JJ, NN, IN, PRP$, NN \rangle$.

- *Surrounding Words* Surrounding words features include all the individual words in the surrounding context of an ambiguous word w . The surrounding words can be in the current sentence or immediately adjacent sentences.

However, we remove the words that are in a list of stop words. Words that contain no alphabetic characters, such as punctuation

¹<http://opennlp.sourceforge.net/>

²<http://www.cis.upenn.edu/~treebank/tokenizer.sed>

³<http://jwordnet.sourceforge.net/>

⁴Syntactic relations are omitted for efficiency reason.

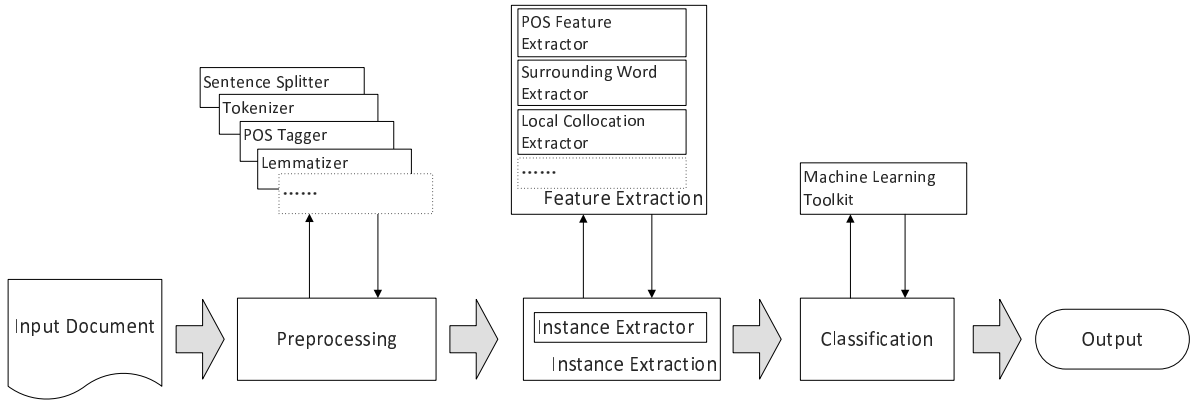


Figure 1: IMS system architecture

symbols and numbers, are also discarded. The remaining words are converted to their lemma forms in lower case. Each lemma is considered as one feature. The feature value is set to be 1 if the corresponding lemma occurs in the surrounding context of w , 0 otherwise.

For example, suppose there is a set of surrounding words features $\{account, economy, rate, take\}$ in the training data set of the word *interest*. For a test instance of *interest* in the sentence “My brother has always taken a keen interest in my work .”, the surrounding word feature vector will be $\langle 0, 0, 0, 1 \rangle$.

- *Local Collocations* We use 11 local collocations features including: $C_{-2,-2}$, $C_{-1,-1}$, $C_{1,1}$, $C_{2,2}$, $C_{-2,-1}$, $C_{-1,1}$, $C_{1,2}$, $C_{-3,-1}$, $C_{-2,1}$, $C_{-1,2}$, and $C_{1,3}$, where $C_{i,j}$ refers to an ordered sequence of words in the same sentence of w . Offsets i and j denote the starting and ending positions of the sequence relative to w , where a negative (positive) offset refers to a word to the left (right) of w .

For example, suppose in the training data set, the word *interest* has a set of local collocations $\{“account .”, “of all”, “in my”, “to be”\}$ for $C_{1,2}$. For a test instance of *interest* in the sentence “My brother has always taken a keen interest in my work .”, the value of feature $C_{1,2}$ will be “in my”.

As shown in Figure 1, we implement one feature extractor for each feature type. The IMS software package is organized in such a way that users can easily specify their own feature set by im-

plementing more feature extractors to exploit new features.

2.1.3 Classification

In IMS, the classifier trains a model for each word type which has training data during the training process. The instances collected in the previous step are converted to the format expected by the machine learning toolkit in use. Thus, the classification step is separate from the feature extraction step. We use LIBLINEAR⁵ (Fan et al., 2008) as the default classifier of IMS, with a linear kernel and all the parameters set to their default values. Accordingly, we implement an interface to convert the instances into the LIBLINEAR feature vector format.

The utilization of other machine learning software can be achieved by implementing the corresponding module interfaces to them. For instance, IMS provides module interfaces to the WEKA machine learning toolkit (Witten and Frank, 2005), LIBSVM⁶, and MaxEnt⁷.

The trained classification models will be applied to the test instances of the corresponding word types in the testing process. If a test instance word type is not seen during training, we will output its predefined default sense, i.e., the WordNet first sense, as the answer. Furthermore, if a word type has neither training data nor predefined default sense, we will output “U”, which stands for the missing sense, as the answer.

⁵<http://www.bwaldvogel.de/liblinear-java/>

⁶<http://www.csie.ntu.edu.tw/~cjlin/libsvm/>

⁷<http://maxent.sourceforge.net/>

2.2 The Training Data Set for All-Words Tasks

Once we have a supervised WSD system, for the users who only need WSD as a component in their applications, it is also important to provide them the classification models. The performance of a supervised WSD system greatly depends on the size of the sense-annotated training data used. To overcome the lack of sense-annotated training examples, besides the training instances from the widely used sense-annotated corpus SEMCOR (Miller et al., 1994) and DSO corpus (Ng and Lee, 1996), we also follow the approach described in Chan and Ng (2005) to extract more training examples from parallel texts.

The process of extracting training examples from parallel texts is as follows:

- Collect a set of sentence-aligned parallel texts. In our case, we use six English-Chinese parallel corpora: Hong Kong Hansards, Hong Kong News, Hong Kong Laws, Sinorama, Xinhua News, and the English translation of Chinese Treebank. They are all available from the Linguistic Data Consortium (LDC).
- Perform tokenization on the English texts with the Penn TreeBank tokenizer.
- Perform Chinese word segmentation on the Chinese texts with the Chinese word segmentation method proposed by Low et al. (2005).
- Perform word alignment on the parallel texts using the GIZA++ software (Och and Ney, 2000).
- Assign Chinese translations to each sense of an English word w .
- Pick the occurrences of w which are aligned to its chosen Chinese translations in the word alignment output of GIZA++.
- Identify the senses of the selected occurrences of w by referring to their aligned Chinese translations.

Finally, the English side of these selected occurrences together with their assigned senses are used as training data.

We only extract training examples from parallel texts for the top 60% most frequently occurring polysemous content words in Brown Corpus

(BC), which includes 730 nouns, 190 verbs, and 326 adjectives. For each of the top 60% nouns and adjectives, we gather a maximum of 1,000 training examples from parallel texts. For each of the top 60% verbs, we extract not more than 500 examples from parallel texts, as well as up to 500 examples from the DSO corpus. We also make use of the sense-annotated examples from SEMCOR as part of our training data for all nouns, verbs, adjectives, and 28 most frequently occurring adverbs in BC.

POS	noun	verb	adj	adv
# of types	11,445	4,705	5,129	28

Table 1: Statistics of the word types which have training data for WordNet 1.7.1 sense inventory

The frequencies of word types which we have training instances for WordNet sense inventory version 1.7.1 are listed in Table 1. We generated classification models with the IMS system for over 21,000 word types which we have training data. On average, each word type has 38 training instances. The total size of the models is about 200 megabytes.

3 Evaluation

In our experiments, we evaluate our IMS system on SensEval and SemEval tasks, the benchmark data sets for WSD. The evaluation on both lexical-sample and all-words tasks measures the accuracy of our IMS system as well as the quality of the training data we have collected.

3.1 English Lexical-Sample Tasks

	SensEval-2	SensEval-3
IMS	65.3%	72.6%
Rank 1 System	64.2%	72.9%
Rank 2 System	63.8%	72.6%
MFS	47.6%	55.2%

Table 2: WSD accuracies on SensEval lexical-sample tasks

In SensEval English lexical-sample tasks, both the training and test data sets are provided. A common baseline for lexical-sample task is to select the most frequent sense (MFS) in the training data as the answer.

We evaluate IMS on the SensEval-2 and SensEval-3 English lexical-sample tasks. Table 2 compares the performance of our system to the top

two systems that participated in the above tasks (Yarowsky et al., 2001; Mihalcea and Moldovan, 2001; Mihalcea et al., 2004). Evaluation results show that IMS achieves significantly better accuracies than the MFS baseline. Comparing to the top participating systems, IMS achieves comparable results.

3.2 English All-Words Tasks

In SensEval and SemEval English all-words tasks, no training data are provided. Therefore, the MFS baseline is no longer suitable for all-words tasks. Because the order of senses in WordNet is based on the frequency of senses in SEMCOR, the WordNet first sense (WNs1) baseline always assigns the first sense in WordNet as the answer. We will use it as the baseline in all-words tasks.

Using the training data collected with the method described in Section 2.2, we apply our system on the SensEval-2, SensEval-3, and SemEval-2007 English all-words tasks. Similarly, we also compare the performance of our system to the top two systems that participated in the above tasks (Palmer et al., 2001; Snyder and Palmer, 2004; Pradhan et al., 2007). The evaluation results are shown in Table 3. IMS easily beats the WNs1 baseline. It ranks first in SensEval-3 English fine-grained all-words task and SemEval-2007 English coarse-grained all-words task, and is also competitive in the remaining tasks. It is worth noting that because of the small test data set in SemEval-2007 English fine-grained all-words task, the differences between IMS and the best participating systems are not statistically significant.

Overall, IMS achieves good WSD accuracies on both all-words and lexical-sample tasks. The performance of IMS shows that it is a state-of-the-art WSD system.

4 Conclusion

This paper presents IMS, an English all-words WSD system. The goal of IMS is to provide a flexible platform for supervised WSD, as well as an all-words WSD component with good performance for other applications.

The framework of IMS allows us to integrate different preprocessing tools to generate knowledge sources. Users can implement various feature types and different machine learning methods or toolkits according to their requirements. By default, the IMS system implements three kinds

of feature types and uses a linear kernel SVM as the classifier. Our evaluation on English lexical-sample tasks proves the strength of our system. With this system, we also provide a large number of classification models trained with the sense-annotated training examples from SEMCOR, DSO corpus, and 6 parallel corpora, for all content words. Evaluation on English all-words tasks shows that IMS with these models achieves state-of-the-art WSD accuracies compared to the top participating systems.

As a Java-based system, IMS is platform independent. The source code of IMS and the classification models can be found on the homepage: <http://nlp.comp.nus.edu.sg/software> and are available for research, non-commercial use.

Acknowledgments

This research is done for CSIDM Project No. CSIDM-200804 partially funded by a grant from the National Research Foundation (NRF) administered by the Media Development Authority (MDA) of Singapore.

References

- Marine Carpuat and Dekai Wu. 2007. Improving statistical machine translation using word sense disambiguation. In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 61–72, Prague, Czech Republic.
- Yee Seng Chan and Hwee Tou Ng. 2005. Scaling up word sense disambiguation via parallel texts. In *Proceedings of the 20th National Conference on Artificial Intelligence (AAAI)*, pages 1037–1042, Pittsburgh, Pennsylvania, USA.
- Yee Seng Chan, Hwee Tou Ng, and David Chiang. 2007a. Word sense disambiguation improves statistical machine translation. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 33–40, Prague, Czech Republic.
- Yee Seng Chan, Hwee Tou Ng, and Zhi Zhong. 2007b. NUS-PT: Exploiting parallel texts for word sense disambiguation in the English all-words tasks. In *Proceedings of the Fourth International Workshop on Semantic Evaluations (SemEval-2007)*, pages 253–256, Prague, Czech Republic.
- Bart Decadt, Veronique Hoste, and Walter Daelemans. 2004. GAMBL, genetic algorithm optimization of memory-based WSD. In *Proceedings of the Third*

	SensEval-2	SensEval-3	SemEval-2007	
	Fine-grained	Fine-grained	Fine-grained	Coarse-grained
IMS	68.2%	67.6%	58.3%	82.6%
Rank 1 System	69.0%	65.2%	59.1%	82.5%
Rank 2 System	63.6%	64.6%	58.7%	81.6%
WNs1	61.9%	62.4%	51.4%	78.9%

Table 3: WSD accuracies on SensEval/SemEval all-words tasks

- International Workshop on Evaluating Word Sense Disambiguation Systems (SensEval-3)*, pages 108–112, Barcelona, Spain.
- Rong-En Fan, Kai-Wei Chang, Cho-Jui Hsieh, Xiang-Rui Wang, and Chih-Jen Lin. 2008. LIBLINEAR: A library for large linear classification. *Journal of Machine Learning Research*, 9:1871–1874.
- Yoong Keok Lee and Hwee Tou Ng. 2002. An empirical evaluation of knowledge sources and learning algorithms for word sense disambiguation. In *Proceedings of the 2002 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 41–48, Philadelphia, Pennsylvania, USA.
- Jin Kiat Low, Hwee Tou Ng, and Wenyuan Guo. 2005. A maximum entropy approach to Chinese word segmentation. In *Proceedings of the Fourth SIGHAN Workshop on Chinese Language Processing*, pages 161–164, Jeju Island, Korea.
- Rada Mihalcea and Andras Csomai. 2005. SenseLearner: Word sense disambiguation for all words in unrestricted text. In *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL) Interactive Poster and Demonstration Sessions*, pages 53–56, Ann Arbor, Michigan, USA.
- Rada Mihalcea and Dan Moldovan. 2001. Pattern learning and active feature selection for word sense disambiguation. In *Proceedings of the Second International Workshop on Evaluating Word Sense Disambiguation Systems (SensEval-2)*, pages 127–130, Toulouse, France.
- Rada Mihalcea, Timothy Chklovski, and Adam Kilgarriff. 2004. The SensEval-3 English lexical sample task. In *Proceedings of the Third International Workshop on Evaluating Word Sense Disambiguation Systems (SensEval-3)*, pages 25–28, Barcelona, Spain.
- George Miller, Martin Chodorow, Shari Landes, Claudia Leacock, and Robert Thomas. 1994. Using a semantic concordance for sense identification. In *Proceedings of ARPA Human Language Technology Workshop*, pages 240–243, Morristown, New Jersey, USA.
- George Miller. 1990. Wordnet: An on-line lexical database. *International Journal of Lexicography*, 3(4):235–312.
- Hwee Tou Ng and Hian Beng Lee. 1996. Integrating multiple knowledge sources to disambiguate word sense: An exemplar-based approach. In *Proceedings of the 34th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 40–47, Santa Cruz, California, USA.
- Franz Josef Och and Hermann Ney. 2000. Improved statistical alignment models. In *Proceedings of the 38th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 440–447, Hong Kong.
- Martha Palmer, Christiane Fellbaum, Scott Cotton, Lauren Delfs, and Hoa Trang Dang. 2001. English tasks: All-words and verb lexical sample. In *Proceedings of the Second International Workshop on Evaluating Word Sense Disambiguation Systems (SensEval-2)*, pages 21–24, Toulouse, France.
- Sameer Pradhan, Edward Loper, Dmitriy Dligach, and Martha Palmer. 2007. SemEval-2007 task-17: English lexical sample, SRL and all words. In *Proceedings of the Fourth International Workshop on Semantic Evaluations (SemEval-2007)*, pages 87–92, Prague, Czech Republic.
- Benjamin Snyder and Martha Palmer. 2004. The English all-words task. In *Proceedings of the Third International Workshop on Evaluating Word Sense Disambiguation Systems (SensEval-3)*, pages 41–43, Barcelona, Spain.
- Christopher Stokoe, Michael P. Oakes, and John Tait. 2003. Word sense disambiguation in information retrieval revisited. In *Proceedings of the Twenty-Sixth Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, pages 159–166, Toronto, Canada.
- Ian H. Witten and Eibe Frank. 2005. *Data Mining: Practical Machine Learning Tools and Techniques*. Morgan Kaufmann, San Francisco, 2nd edition.
- David Yarowsky, Radu Florian, Siviú Cucerzan, and Charles Schafer. 2001. The Johns Hopkins SensEval-2 system description. In *Proceedings of the Second International Workshop on Evaluating Word Sense Disambiguation Systems (SensEval-2)*, pages 163–166, Toulouse, France.

Author Index

- Adolphs, Peter, 36
Angelov, Krasimir, 66
- Bangalore, Srinivas, 60
Bender, Emily M., 1
Blunsom, Phil, 7
Byrne, William, 48
- Callaway, Charles B., 13
Campbell, Gwendolyn, 13
Catizone, Roberta, 72
Cheng, Weiwei, 72
Cheng, Xiwen, 36
- Dines, John, 48
Dingli, Alexiei, 72
Drellishak, Scott, 1
Dyer, Chris, 7
Dzikovska, Myroslava O., 13
- Eidelman, Vladimir, 7
- Farrow, Elaine, 13
Fokkens, Antske, 1
- Ganitkevitch, Juri, 7
Garner, Philip N., 48
Gibson, Matthew, 48
Goodman, Michael Wayne, 1
Guan, Yong, 48
- Hallgren, Thomas, 66
Henrich, Verena, 19
Hinrichs, Erhard, 19, 25
Hinrichs, Marie, 25
Hirsimäki, Teemu, 48
- Jurgens, David, 30
- Karhila, Reima, 48
King, Simon, 48
Klüwer, Tina, 36
Kouylekov, Milen, 42
Kurimo, Mikko, 48
- Leidner, Jochen, 54
- Liang, Hui, 48
Lopez, Adam, 7
- Mills, Daniel P., 1
Mishra, Taniya, 60
Moore, Johanna D., 13
- Negri, Matteo, 42
Ng, Hwee Tou, 78
- Oura, Keiichiro, 48
- Poulson, Laurie, 1
- Ranta, Aarne, 66
Resnik, Philip, 7
- Saheer, Lakshmi, 48
Saleem, Safiyyah, 1
Schilder, Frank, 54
Setiawan, Hendra, 7
Shannon, Matt, 48
Shiota, Sayaki, 48
Steinhauser, Natalie, 13
Stevens, Keith, 30
- Tian, Jilei, 48
Ture, Ferhan, 7
- Uszkoreit, Hans, 36
- Weese, Jonathan, 7
Wilks, Yorick, 72
- Xu, Feiyu, 36
- Zastrow, Thomas, 25
Zhong, Zhi, 78

ACL 2010

**The 48th Annual Meeting of the
Association for Computational Linguistics**

Tutorial Abstracts

11 July 2010
Uppsala University
Uppsala, Sweden

Production and Manufacturing by
Taberg Media Group AB
Box 94, 562 02 Taberg
Sweden

©2010 The Association for Computational Linguistics

Order copies of this and other ACL proceedings from:

Association for Computational Linguistics (ACL)
209 N. Eighth Street
Stroudsburg, PA 18360
USA
Tel: +1-570-476-8006
Fax: +1-570-476-0860
acl@aclweb.org

Introduction

A total of 20 tutorial proposals were submitted to the ACL 2010 Tutorials track, from which 6 were finally accepted. We are very grateful to the ACL-NLP community for the large amount of exciting, diverse, and high-quality proposals we received. This guaranteed a strong final program, but at the same time made the selection process very difficult. It was really hard to reject some of the proposals. All 20 proposals were reviewed by both co-chairs, and the final selection was approved by the conference General Chair. We also sought expertise from external reviewers where necessary.

Based on the following criteria we selected the maximum number of proposals allowed: 1) Quality: the content and scope of the proposal, and the competence and experience of the presenters; 2) Diversity: We sought a range of different topics and approaches; 3) Appeal: Whether the tutorial topic would be likely to attract a reasonable number of participants; and 4) Novelty: Tutorial topics featured at very recent ACL events were dispreferred (unless the content was clearly novel and different).

The final tutorial programme covers a wide range of NLP topics, including Annotation, Grammars, Discourse Structure, Structured Prediction, Semantic Parsing and Machine Translation.

Sincere thanks to all authors for preparing and sending all the information and materials timely. We know that this is not always easy because of the tight schedule and other competing commitments.

We are very indebted to previous tutorial chairs, especially to Diana McCarthy and Chengqing Zong for sharing with us their experience and already developed materials. They provided also valuable advice all throughout the process. We are also equally grateful to the ACL 2010 General Chair, and the Local/Publicity/Publications Chairs for their help and advice in the organization of the Tutorials program and materials, and for sending always the appropriate reminders!

Finally, we only hope that the participation will be as successful as all the previous steps.

To all the attendees, enjoy the ACL 2010 tutorials!

Lluís Màrquez (Technical University of Catalonia, Spain)
Haifeng Wang (Baidu, Inc., China)
ACL 2010 Tutorial Co-Chairs

Tutorial Co-Chairs:

Lluís Màrquez, Technical University of Catalonia, Spain
Haifeng Wang, Baidu, Inc., China

Table of Contents

<i>Wide-Coverage NLP with Linguistically Expressive Grammars</i>	
Julia Hockenmaier, Yusuke Miyao and Josef van Genabith	1
<i>Tree-Based and Forest-Based Translation</i>	
Yang Liu and Liang Huang	2
<i>Discourse Structure: Theory, Practice and Use</i>	
Bonnie Webber, Markus Egg and Valia Kordoni	3
<i>Annotation</i>	
Eduard Hovy	4
<i>From Structured Prediction to Inverse Reinforcement Learning</i>	
Hal Daumé III	5
<i>Semantic Parsing: The Task, the State of the Art and the Future</i>	
Rohit J. Kate and Yuk Wah Wong	6

Tutorial Program

Sunday July 11, 2010

Morning

9:00-12:30

Wide-Coverage NLP with Linguistically Expressive Grammars

Julia Hockenmaier, Yusuke Miyao and Josef van Genabith

Tree-Based and Forest-Based Translation

Yang Liu and Liang Huang

Discourse Structure: Theory, Practice and Use

Bonnie Webber, Markus Egg and Valia Kordoni

Afternoon

2:00-5:30

Annotation

Eduard Hovy

From Structured Prediction to Inverse Reinforcement Learning

Hal Daumé III

Semantic Parsing: The Task, the State of the Art and the Future

Rohit J. Kate and Yuk Wah Wong

Wide-coverage NLP with Linguistically Expressive Grammars

Julia Hockenmaier

Department of Computer Science,
University of Illinois
juliahmr@illinois.edu

Yusuke Miyao

National Institute of Informatics
yusuke@nii.ac.jp

Josef van Genabith

Centre for Next Generation Localisation,
School of Computing,
Dublin City University
josef@computing.dcu.ie

1 Introduction

In recent years, there has been a lot of research on wide-coverage statistical natural language processing with linguistically expressive grammars such as Combinatory Categorical Grammars (CCG), Head-driven Phrase-Structure Grammars (HPSG), Lexical-Functional Grammars (LFG) and Tree-Adjoining Grammars (TAG). But although many young researchers in natural language processing are very well trained in machine learning and statistical methods, they often lack the necessary background to understand the linguistic motivation behind these formalisms. Furthermore, in many linguistics departments, syntax is still taught from a purely Chomskian perspective. Additionally, research on these formalisms often takes place within tightly-knit, formalism-specific subcommunities. It is therefore often difficult for outsiders as well as experts to grasp the commonalities of and differences between these formalisms.

2 Content Overview

This tutorial overviews basic ideas of TAG/CCG/LFG/HPSG, and provides attendees with a comparison of these formalisms from a linguistic and computational point of view. We start from stating the motivation behind using these expressive grammar formalisms for NLP, contrasting them with shallow formalisms like context-free grammars. We introduce a common set of examples illustrating various linguistic constructions that elude context-free grammars, and reuse them when introducing each formalism: bounded and unbounded non-local dependencies that arise through extraction and coordination, scrambling, mappings to meaning representations, etc. In the

second half of the tutorial, we explain two key technologies for wide-coverage NLP with these grammar formalisms: grammar acquisition and parsing models. Finally, we show NLP applications where these expressive grammar formalisms provide additional benefits.

3 Tutorial Outline

1. Introduction: Why expressive grammars
2. Introduction to TAG
3. Introduction to CCG
4. Introduction to LFG
5. Introduction to HPSG
6. Inducing expressive grammars from corpora
7. Wide-coverage parsing with expressive grammars
8. Applications
9. Summary

References

- Aoife Cahill, Michael Burke, Ruth O'Donovan, Stefan Riezler, Josef van Genabith and Andy Way. 2008. Wide-Coverage Deep Statistical Parsing using Automatic Dependency Structure Annotation. *Computational Linguistics*, 34(1). pp.81-124, MIT Press.
- Yusuke Miyao and Jun'ichi Tsujii. 2008. Feature Forest Models for Probabilistic HPSG Parsing. *Computational Linguistics*, 34(1). pp.35-80, MIT Press.
- Julia Hockenmaier and Mark Steedman. 2007. CCG-bank: A Corpus of CCG Derivations and Dependency Structures Extracted from the Penn Treebank. *Computational Linguistics*, 33(3). pp.355-396, MIT Press.

Tree-based and Forest-based Translation

Yang Liu

Institute of Computing Technology
Chinese Academy of Sciences
yliu@ict.ac.cn

Liang Huang

Information Sciences Institute
University of Southern California
lhuang@isi.edu

1 Introduction

The past several years have witnessed rapid advances in syntax-based machine translation, which exploits natural language syntax to guide translation. Depending on the type of input, most of these efforts can be divided into two broad categories: (a) **string-based systems** whose input is a string, which is simultaneously parsed and translated by a synchronous grammar (Wu, 1997; Chiang, 2005; Galley et al., 2006), and (b) **tree-based systems** whose input is already a parse tree to be directly converted into a target tree or string (Lin, 2004; Ding and Palmer, 2005; Quirk et al., 2005; Liu et al., 2006; Huang et al., 2006).

Compared with their string-based counterparts, tree-based systems offer many attractive features: they are much faster in decoding (linear time vs. cubic time), do not require sophisticated binarization (Zhang et al., 2006), and can use separate grammars for parsing and translation (e.g. a context-free grammar for the former and a tree substitution grammar for the latter).

However, despite these advantages, most tree-based systems suffer from a major drawback: they only use 1-best parse trees to direct translation, which potentially introduces translation mistakes due to parsing errors (Quirk and Corston-Oliver, 2006). This situation becomes worse for resource-poor source languages without enough Treebank data to train a high-accuracy parser.

This problem can be alleviated elegantly by using packed forests (Huang, 2008), which encodes exponentially many parse trees in a polynomial space. Forest-based systems (Mi et al., 2008; Mi and Huang, 2008) thus take a packed forest instead of a parse tree as an input. In addition, packed forests could also be used for translation rule extraction, which helps alleviate the propagation of parsing errors into rule set. Forest-based translation can be regarded as a compromise between the string-based and tree-based methods, while com-

binning the advantages of both: decoding is still fast, yet does not commit to a single parse. Surprisingly, translating a forest of millions of trees is even faster than translating 30 individual trees, and offers significantly better translation quality. This approach has since become a popular topic.

2 Content Overview

This tutorial surveys tree-based and forest-based translation methods. For each approach, we will discuss the two fundamental tasks: *decoding*, which performs the actual translation, and *rule extraction*, which learns translation rules from real-world data automatically. Finally, we will introduce some more recent developments to tree-based and forest-based translation, such as tree sequence based models, tree-to-tree models, joint parsing and translation, and faster decoding algorithms. We will conclude our talk by pointing out some directions for future work.

3 Tutorial Overview

1. Tree-based Translation

- Motivations and Overview
- Tree-to-String Model and Decoding
- Tree-to-String Rule Extraction
- Language Model-Integrated Decoding: Cube Pruning

2. Forest-based Translation

- Packed Forest
- Forest-based Decoding
- Forest-based Rule Extraction

3. Extensions

- Tree-Sequence-to-String Models
- Tree-to-Tree Models
- Joint Parsing and Translation
- Faster Decoding Methods

4. Conclusion and Open Problems

Discourse Structure: Theory, Practice and Use

Bonnie Webber,[♡] Markus Egg,[◇] Valia Kordoni[♠]

♡ University of Edinburgh ◇ Humboldt University ♠ Saarland University
bonnie@inf.ed.ac.uk markus.egg@anglistik.hu-berlin.de kordoni@dfki.de

1 Introduction

This tutorial aims to provide attendees with a clear notion of how discourse structure is relevant for language technology (LT), what is needed for exploiting discourse structure, what methods and resources are available to support its use, and what more could be done in the future.

2 Content Overview

This tutorial consists of four parts. Part I starts with a brief introduction to different bases for discourse structuring, properties of discourse structure that are relevant to LT, and accessible evidence for discourse structure.

For discourse structure to be useful for language technologies, one must be able to automatically recognize or generate with it. Hence, Part II surveys computational approaches to recognizing and generating discourse structure, both manually-authored approaches and ones developed through Machine Learning.

Part III of the tutorial describes applications of discourse structure recognition and generation in LT, as well as discourse-related resources being made available in English, German, Turkish, Hindi, Czech, Arabic and Chinese. Part IV concludes with a list of future possibilities.

3 Tutorial Outline

1. PART I – General Overview
 - (a) Bases for structure in monologic, dialogic and multiparty discourse
 - (b) Aspects of discourse structure relevant to Language Technology
 - (c) Evidence for discourse structure
2. PART II – Computational Recognition and Generation of discourse structure

- (a) Discourse chunking and parsing
- (b) Recognizing arguments and sense of discourse connectives
- (c) Recognizing and generating entity-based discourse structure
- (d) Dialogue parsing

3. PART III – Applications and Resources

- (a) Applications to Language Technology
- (b) Discourse structure resources (monolingual and multilingual)

4. PART IV – Future Developments

4 References

- Regina Barzilay and Lillian Lee (2004). Catching the Drift: Probabilistic Content Models, with Applications to Generation and Summarization. *Proc. 2nd Human Language Technology Conference and Annual Meeting of the North American Chapter, Association for Computational Linguistics*, pp. 113-120.
- Regina Barzilay and Mirella Lapata (2008). Modeling Local Coherence: An Entity-based Approach. *Computational Linguistics* 34(1), pp. 1-34.
- Daniel Marcu (2000). *The theory and practice of discourse parsing and summarization*. Cambridge: MIT Press.
- Umangi Oza, Rashmi Prasad, Sudheer Kolachina, Dipti Misra Sharma and Aravind Joshi (2009). The Hindi Discourse Relation Bank. *Proc. Third Linguistic Annotation Workshop (LAW III)*. Singapore.
- Rashmi Prasad, Nikhil Dinesh, Alan Lee, Eleni Miltsakaki et al. (2008). The Penn Discourse TreeBank 2.0. *Proc. 6th Int'l Conference on Language Resources and Evaluation*.
- Manfred Stede (2008). RST revisited: Disentangling nuclearity. In Cathrine Fabricius-Hansen and Wiebke Ramm (eds.), *Subordination versus Coordination in Sentence and Text*. Amsterdam: John Benjamins.
- Ben Wellner (2008). *Sequence Models and Ranking Methods for Discourse Parsing*. Brandeis University.
- Deniz Zeyrek, Ümit Deniz Turan, Cem Bozsahin, Ruket Çakici et al. (2009). Annotating Subordinators in the Turkish Discourse Bank. *Proc. Third Linguistic Annotation Workshop (LAW III)*. Singapore.

Annotation

Eduard Hovy

Information Sciences Institute
University of Southern California
email: hovy@isi.edu

1. Introduction

As researchers seek to apply their machine learning algorithms to new problems, corpus annotation is increasingly gaining importance in the NLP community. But since the community currently has no general paradigm, no textbook that covers all the issues (though Wilcock's book published in Dec 2009 covers some basic ones very well), and no accepted standards, setting up and performing small-, medium-, and large-scale annotation projects remains something of an art.

To attend, no special expertise in computation or linguistics is required.

2. Content Overview

This tutorial is intended to provide the attendee with an in-depth look at the procedures, issues, and problems in corpus annotation, and highlights the pitfalls that the annotation manager should avoid. The tutorial first discusses why annotation is becoming increasingly relevant for NLP and how it fits into the generic NLP methodology of train-evaluate-apply. It then reviews currently available resources, services, and frameworks that support someone wishing to start an annotation project easily. This includes the QDAP annotation center, Amazon's Mechanical Turk, annotation facilities in GATE, and other resources such as UIMA. It then discusses the seven major open issues at the heart of annotation for which there are as yet no standard and fully satisfactory answers or methods. Each issue is described in detail and current practice is shown. The seven issues are: 1. How does one decide what specific phenomena to annotate? How does one adequately capture the theory behind the phenomenon/a and express it in simple annotation instructions? 2. How does one obtain a balanced corpus to annotate, and when is a corpus balanced (and representative)? 3. When hiring annotators, what characteristics are important? How does one ensure that they are adequately (but not over- or under-) trained? 4. How does one

establish a simple, fast, and trustworthy annotation procedure? How and when does one apply measures to ensure that the procedure remains on track? How and where can active learning help? 5. What interface(s) are best for each type of problem, and what should one know to avoid? How can one ensure that the interfaces do not influence the annotation results? 6. How does one evaluate the results? What are the appropriate agreement measures? At which cutoff points should one redesign or re-do the annotations? 7. How should one formulate and store the results? When, and to whom, should one release the corpus? How should one report the annotation effort and results for best impact?

The notes include several pages of references and suggested readings.

3. Tutorial Overview

1. Toward a Science of Annotation
 - a. What is Annotation, and Why do We Need It?
2. Setting up an Annotation Project
 - a. The Basic Steps
 - b. Useful Resources and Services
3. Examples of Annotation Projects
4. The Seven Questions of Annotation
 - a. Instantiating the Theory
 - b. Selecting the Corpus
 - c. Designing the Annotation Interface
 - d. Selecting and Training Annotators
 - e. Specifying the Annotation Procedure
 - f. Evaluation and Validation
 - g. Distribution and Maintenance
5. Closing: The Future of Annotation in NLP

Discourse Structure: Theory, Practice and Use

Bonnie Webber,[♡] Markus Egg,[◇] Valia Kordoni[♠]

♡ University of Edinburgh ◇ Humboldt University ♠ Saarland University
bonnie@inf.ed.ac.uk markus.egg@anglistik.hu-berlin.de kordoni@dfki.de

1 Introduction

This tutorial aims to provide attendees with a clear notion of how discourse structure is relevant for language technology (LT), what is needed for exploiting discourse structure, what methods and resources are available to support its use, and what more could be done in the future.

2 Content Overview

This tutorial consists of four parts. Part I starts with a brief introduction to different bases for discourse structuring, properties of discourse structure that are relevant to LT, and accessible evidence for discourse structure.

For discourse structure to be useful for language technologies, one must be able to automatically recognize or generate with it. Hence, Part II surveys computational approaches to recognizing and generating discourse structure, both manually-authored approaches and ones developed through Machine Learning.

Part III of the tutorial describes applications of discourse structure recognition and generation in LT, as well as discourse-related resources being made available in English, German, Turkish, Hindi, Czech, Arabic and Chinese. Part IV concludes with a list of future possibilities.

3 Tutorial Outline

1. PART I – General Overview
 - (a) Bases for structure in monologic, dialogic and multiparty discourse
 - (b) Aspects of discourse structure relevant to Language Technology
 - (c) Evidence for discourse structure
2. PART II – Computational Recognition and Generation of discourse structure

- (a) Discourse chunking and parsing
- (b) Recognizing arguments and sense of discourse connectives
- (c) Recognizing and generating entity-based discourse structure
- (d) Dialogue parsing

3. PART III – Applications and Resources

- (a) Applications to Language Technology
- (b) Discourse structure resources (monolingual and multilingual)

4. PART IV – Future Developments

4 References

- Regina Barzilay and Lillian Lee (2004). Catching the Drift: Probabilistic Content Models, with Applications to Generation and Summarization. *Proc. 2nd Human Language Technology Conference and Annual Meeting of the North American Chapter, Association for Computational Linguistics*, pp. 113-120.
- Regina Barzilay and Mirella Lapata (2008). Modeling Local Coherence: An Entity-based Approach. *Computational Linguistics* 34(1), pp. 1-34.
- Daniel Marcu (2000). *The theory and practice of discourse parsing and summarization*. Cambridge: MIT Press.
- Umangi Oza, Rashmi Prasad, Sudheer Kolachina, Dipti Misra Sharma and Aravind Joshi (2009). The Hindi Discourse Relation Bank. *Proc. Third Linguistic Annotation Workshop (LAW III)*. Singapore.
- Rashmi Prasad, Nikhil Dinesh, Alan Lee, Eleni Miltsakaki et al. (2008). The Penn Discourse TreeBank 2.0. *Proc. 6th Int'l Conference on Language Resources and Evaluation*.
- Manfred Stede (2008). RST revisited: Disentangling nuclearity. In Cathrine Fabricius-Hansen and Wiebke Ramm (eds.), *Subordination versus Coordination in Sentence and Text*. Amsterdam: John Benjamins.
- Ben Wellner (2008). *Sequence Models and Ranking Methods for Discourse Parsing*. Brandeis University.
- Deniz Zeyrek, Ümit Deniz Turan, Cem Bozsahin, Ruket Çakici et al. (2009). Annotating Subordinators in the Turkish Discourse Bank. *Proc. Third Linguistic Annotation Workshop (LAW III)*. Singapore.

Semantic Parsing: The Task, the State-of-the-Art and the Future

Rohit J. Kate

Department of Computer Science
The University of Texas at Austin
Austin, TX 78712, USA
rjkate@cs.utexas.edu

Yuk Wah Wong

Google Inc.
Pittsburgh, PA 15213, USA
ywwong@google.com

1 Introduction

Semantic parsing is the task of mapping natural language sentences into complete formal meaning representations which a computer can execute for some domain-specific application. This is a challenging task and is critical for developing computing systems that can understand and process natural language input, for example, a computing system that answers natural language queries about a database, or a robot that takes commands in natural language. While the importance of semantic parsing was realized a long time ago, it is only in the past few years that the state-of-the-art in semantic parsing has been significantly advanced with more accurate and robust semantic parser learners that use a variety of statistical learning methods. Semantic parsers have also been extended to work beyond a single sentence, for example, to use discourse contexts and to learn domain-specific language from perceptual contexts. Some of the future research directions of semantic parsing with potentially large impacts include mapping entire natural language documents into machine processable form to enable automated reasoning about them and to convert natural language web pages into machine processable representations for the Semantic Web to support automated high-end web applications.

This tutorial will introduce the semantic parsing task and will bring the audience up-to-date with the current research and state-of-the-art in semantic parsing. It will also provide insights about semantic parsing and how it relates to and differs from other natural language processing tasks. It will point out research challenges and some promising future directions for semantic parsing.

2 Content Overview

The proposed tutorial on semantic parsing will start with an introduction to the task, giving ex-

amples of some application domains and meaning representation languages. It will also point out its distinctions from and relations to other NLP tasks. Next, it will talk in depth about various semantic parsers that have been built, starting with earlier hand-built systems to the current state-of-the-art statistical semantic parser learners. It will point out the underlying commonalities and differences between the learners. The next section of the tutorial will talk about the recent advances in extending semantic parsing to work beyond parsing a single sentence. Finally, the tutorial will point out the current research challenges and some promising future directions for semantic parsing.

3 Outline

1. Introduction to the task of semantic parsing
 - (a) Definition of the task
 - (b) Examples of application domains and meaning representation languages
 - (c) Distinctions from and relations to other NLP tasks
2. Semantic parsers
 - (a) Earlier hand-built systems
 - (b) Learning for semantic parsing
 - i. Semantic parsing learning task
 - ii. Non-statistical semantic parser learners
 - iii. Statistical semantic parser learners
 - iv. Exploiting syntax for semantic parsing
 - v. Various forms of supervision: semi-supervision, ambiguous supervision
 - (c) Underlying commonality and differences between different semantic parser learners
3. Semantic parsing beyond a sentence
 - (a) Using discourse contexts for semantic parsing
 - (b) Learning language from perceptual contexts
4. Research challenges and future directions
 - (a) Machine reading of documents: Connecting with knowledge representation
 - (b) Applying semantic parsing techniques to the Semantic Web
 - (c) Future research directions
5. Conclusions

Author Index

Daumé III, Hal, 5

Egg, Markus, 3

Genabith, Josef van, 1

Hockenmaier, Julia, 1

Hovy, Eduard, 4

Huang, Liang, 2

Kate, Rohit J., 6

Kordoni, Valia, 3

Liu, Yang, 2

Miyao, Yusuke, 1

Webber, Bonnie, 3

Wong, Yuk Wah, 6