

[Day-2] Unsupervised Learning for Natural Language Processing

(Part IV: Potential Traps & Source of Problems)

Keh-Yih Su and Jing-Shin Chang
kysu@bdc.com.tw, shin@nlp.csie.ncnu.edu.tw

(2002/8/18)

Behavior Design Corporation (BDC)
No. 5, 2F, Industrial East Road IV, Science-Based Industrial Park
Hsinchu, Taiwan 30077, R.O.C.

Department of Computer Science and Information Engineering
National Chi-Nan University
Puli, Nantou, Taiwan 545, R.O.C.

Day-2: Unsupervised Learning for Natural Language Processing

- Part I: Introduction
 - ◆ What and When for Unsupervised Learning, Why it is getting popular
- Part II: Basic Concepts and Background (using EM as an example)
 - ◆ Incomplete Data Space
 - ◆ Learnability
- Part III: Typical Unsupervised Learning Algorithms: Viterbi & EM
 - ◆ Procedures, Characteristics
- **Part IV: Potential Traps & Source of Problems**
 - ◆ **Various Mismatches, Model Deficiencies, Local Maximum, and Over-fitting**
- Part V: Suggested Strategies for Better Performance
 - ◆ Lessons Learned from Past Experience
 - ◆ Recommended Procedures for Unsupervised Learning
- Part VI: Advanced Topic: Co-Training
 - ◆ Basic Principles
 - ◆ Example: Chinese New Word Extraction
- References

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

2

Part IV: Advanced Topics: Potential Traps, Sources of Problems, and Why

- Criteria Mismatch: Human Preference in Testing Set vs. Model Fitting in Training Set
 - ◆ Mismatch of Measuring Functions
 - ◆ Mismatch of Measuring Sources
 - ◆ Implied Assumptions During Problem Solving
- Sources Causing Mismatch
 - ◆ Model Deficiency
 - ◆ Local Traps
 - ◆ Insufficient Training Data
 - ◆ Statistical Characteristics Variation
- Methods to Reduce Mismatch Effect
 - ◆ Reduce Measuring Function Mismatch Effect
 - ◆ Reduce Measuring Source Mismatch Effect

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

3

Fundamental Problem with Unsupervised Learning -- Criteria Mismatch (1)

- Criteria Mismatch: Human Preference in Testing Set vs. Model Fitting in Training Set
 - ◆ System Performance: Error rate in the testing set
 - ◆ Error rate measures the fitting for human preference
 - ◆ Unsupervised Learning Convergence Direction: Maximum of Likelihood Values in the training set
 - ◆ Likelihood value measures the fitting for the adopted model
 - ◆ Two measures are not necessarily to be highly correlated, if not under proper setting
 - ◆ Sources Resulting Mismatch
 - ◆ Adopting different measuring functions
 - ◆ Sampling from different sources

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

4

Criteria Mismatch (2):

- Unsupervised Learning Wish:
 - ◆ System Performance is getting improved iteration by iteration
 - ◆ The iteration process will finally converge to the point in the parameter space which possesses the minimum error rate performance (measured in the testing set)
- Implied Assumption for the success of unsupervised learning:
Increasing Training Set Likelihood Values \Rightarrow Decreasing Testing Set Error Rate
 - ◆ Monotonically increasing of likelihood value in the training set (no problem, it is guaranteed)
 - ◆ Likelihood Value Increases \Rightarrow Error Rate Decreases (in both training set and testing set)
 - ◆ Maximizing Likelihood Value \Rightarrow Minimizing Error Rate (in both training set and testing set)

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

5

Criteria Mismatch (3):

- Implied Assumption for the success of unsupervised learning (cont.):
 - ◆ Increasing Training Set Likelihood Value \Rightarrow Increasing Testing Set Likelihood Value
 - ◆ Maximizing Training Set Likelihood Values \Rightarrow Maximizing Testing Set Likelihood Values
- Implied Conclusion:
 - ◆ Increasing Likelihood Value in Training Set \Rightarrow Decreasing the Error Rate in the Testing Set
 - ◆ Maximizing Training Set Likelihood Values \Rightarrow Minimizing Testing Set Error Rate

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

6

Mismatch between Measuring Functions (1):

- *Model Fitting* (Maximizing Likelihood Value) versus *Preference Finding* (Minimizing Error Rate)
- Many learning methods (designed for the recognition task) pursue “minimizing error rate” indirectly via training the model with other “optimizing criteria”
 - ◆ Possible criteria
 - ◆ Minimal Sum of Square Error (e.g., Clustering, VQ, etc.)
 - ◆ Minimal Inter-Cluster Distance (e.g., Clustering, VQ, etc.)
 - ◆ Maximal Likelihood Value (e.g., EM, Viterbi)
 - ◆ Maximum Entropy (e.g., IBM Maximum Entropy approach), etc.
 - ◆ Implicit Assumption: the model that can optimize the chosen Criterion can also achieve the minimum error rate performance

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

7

Mismatch between Measuring Functions (2):

- Mismatch Result:
 - ◆ The parameter set that maximizes the likelihood in the training set is not the one which can really minimize the error rate in the training set.
 - ◆ Indirectly adjusting parameters is relatively ineffective (and sometimes awkward)
- However, it could still be used as a good starting point in supervised learning, and no other better way in unsupervised learning
 - ◆ There is still no good statistical model that can directly pursue the correct ranking order so far.
 - ◆ Bayesian framework is sound and relatively good (comparing to other approaches)
 - ◆ It just needs a little twist for fine tuning in supervised learning

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

8

Mismatch between Measuring Functions (3):

- For remedying the drawbacks above mentioned, Discriminative Training was proposed to directly pursue “Minimizing error rate” in supervised learning
 - ◆ Approximate each error by an analytical Loss Function (e.g., arctan or sigmoid)
 - ◆ Searching the parameter space for minimizing the corresponding Risk Function
 - ◆ Under discriminative training, minimizing risk function in the training set does imply minimizing error rate in the training set
 - ◆ Result: Better performance, more effective in adjusting parameters

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

9

Mismatch between Measuring Functions (4):

- In unsupervised learning, however, human preference is not known; therefore, discriminative training cannot be applied in many situations
 - ◆ Errors can no longer be perceived in the training set
 - ◆ Therefore, the error rate cannot be used as the searching criterion
- Mismatch between measuring functions is thus unavoidable
 - ◆ Result: optimizing the chosen criterion in the training set does not imply we can also minimize the error rate in the training set

$$\hat{\Lambda}_{MLE}(TR) \neq \hat{\Lambda}_{Err}(TR); \quad \hat{\Lambda}_{MLE}(TS) \neq \hat{\Lambda}_{Err}(TS)$$

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

10

Mismatch between Measuring Functions (5):

- Example: using a bi-gram model for part of speech tagging; assume each word in the corpus has exactly two tags (e.g., noun and verb)
 - ◆ Switching noun and verb of the best (human preferred) tag sequence results in the same likelihood value: just an exchange of the labels; however, it would result 100% and 0% accuracy rates, respectively
- To make unsupervised learning work, a high correlation between those two measures (likelihood in the training set & error rate in the training set) must be inherited (or implied) from the model
 - ◆ The higher the degree of correlation, the better the chance for obtaining good performance
 - ◆ However, these two measures will not automatically closely correlate with each other if not under proper setting

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

11

Mismatch between Sampling Sources:

- Mismatch of Sampling Sources: *Training Set* vs. *Testing Set*
 - ◆ Statistical learning methods implicitly assume that the parameters obtained from the training set are also applicable to the testing set
 - ◆ Implicit Assumption:
 - ◆ Both the training set and the testing set have identical statistical characteristics
 - ◆ The parameter estimation error is negligible (i.e., the training set have almost infinitive sampling size)
 - ◆ Implied Conclusion:
 - ◆ Maximizing Training Set Likelihood Value => Maximizing Testing Set Likelihood Value; however, it is not guaranteed
 - ◆ Minimizing Training Set Error Rate => Minimizing Testing Set Error Rate; again, it is not guaranteed

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

12

Mismatch between Sampling Sources (cont.):

■ Factors that cause mismatch:

- ◆ Statistical characteristics variation (possible sources: sampling from different domains) between training set and testing set
- ◆ Finite sampling size: causing estimation error (Note: the estimation error cannot be perceived in the training set)

■ Result:

- ◆ The parameter set that can maximize the *likelihood value* in the *training* set might not be the one that can also do the same in the *testing* set
- ◆ The parameter set that can minimize the *error rate* in the *training* set might not be the one that can also minimize the error rate in the *testing* set

$$\hat{\Lambda}_{MLE}(TR) \neq \hat{\Lambda}_{MLE}(TS); \quad \hat{\Lambda}_{Err}(TR) \neq \hat{\Lambda}_{Err}(TS)$$

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

13

Sources for Causing Mismatch:

■ Model Deficiency

- ◆ Inappropriate Feature Set
- ◆ Inappropriate Feature Dependency Relationship
- ◆ Which would cause the mismatch between two measuring functions

■ Local Traps

- ◆ Multiple local optimum points inherent in the parameter space
- ◆ Which would cause the mismatch between two measuring functions

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

14

Sources for Causing Mismatch (Cont.):

■ Insufficient Training Data

- ◆ Large estimation error perceived in the testing set
- ◆ Which would cause the mismatch between two sampling sources

■ Statistical Characteristic Variation

- ◆ Different statistical characteristics between training set and testing set
- ◆ Which would cause the mismatch between two sampling sources

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

15

Inappropriate Feature Set (1)

■ The selected Feature Space decides Performance Upper Bound

- ◆ Once the feature space is specified, the best reachable performance is also determined for the given task. The system designers can only try to find a good discriminator to approach the upper bound.
- ◆ Feature selection is probably the most important step
 - ◆ Problem Analysis is usually required (versus black-box approach)

■ Feature Set Mismatch:

- ◆ Using naive raw features instead of preference-based features:
 - ◆ Surface-level features (e.g., words) are used, instead of deeper level features, in the adopted stochastic language model
 - ◆ Unable to catch underlying linguistic units based on which human really uses to make preference
- ◆ Causing the mismatch between two measuring functions

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

16

Inappropriate Feature Set (2)

- Naive stochastic language model usually fails to catch Long-distance Dependency (frequently adopted by the human preference model) implied in the deep structure
 - ◆ Surface word N-gram was usually adopted
 - ◆ With heuristically determined window size (to avoid exponential explosion of the number of parameters)
 - ◆ Believe “Data is the King”: just get more data (better to be annotated)
 - ◆ Usually require a huge number of parameters (as no classes are adopted)
 - ◆ At the cost of lower performance by ignoring the features that the long distance dependency requires

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

17

Inappropriate Feature Set (3)

- Mismatch of Feature Set Examples:
 - ◆ Semantic tags assignment: Semantic Markov Chain (which adopts Semantic Tag N-gram) versus Head-Features
 - ◆ with heuristically determined window size
 - ◆ Aligning bi-lingual sentences: using length-based feature instead of transfer dictionary
 - ◆ IBM Machine Translation Model (I): Free-order word-string versus BDC BehaviorTran linguistic structure
 - ◆ OCR/OLCR: crossing counts versus Chinese strokes

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

18

Inappropriate Feature Dependencies (1)

- Dependency Relationship Mismatch will make the measuring functions to be unmatched
- Inappropriate Markov Assumption is widely assumed
 - ◆ Most Markov models only keep a few nearest adjacent neighbors, and drop those constituents that are relatively farther (i.e., only handle local dependency)
 - ◆ May not reflect real dependencies among constituents (i.e., the human preference network in which long distance dependency is usually implied)
 - ◆ Example: use bi-gram model to predict the next word when the next word really depends on a head word that is ten words away.
 - ✦ The prediction power, implied by the dependency, provided by the head word will attenuate to almost nothing after 10-step state transitions

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

19

Inappropriate Feature Dependencies (2)

- Conditional Independence is inappropriately assumed
 - ◆ Assuming features are conditional independent (which is frequently used to drop terms) while they are actually highly correlated
 - ✦ Example: $P(f_1, f_2, f_3 | c_i) \cong P(f_1 | c_i) \times P(f_2 | c_i) \times P(f_3 | c_i)$
 - ◆ Some features in the adopted feature set are highly correlated, the strong dependency should be utilized in the model
 - ✦ For example, it might be better reduced to:

$$P(f_1, f_2, f_3 | c_i) \cong P(f_1 | f_2, c_i) \times P(f_3 | f_2, c_i) \times P(f_2 | c_i)$$

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

20

Inappropriate Feature Dependencies (3)

- Example: Parse tree selection with Stochastic Context Free Grammar versus Context-Sensitive Layered-Scoring Function [Su 88, Chiang 96]
 - ◆ A language can be represented by a context-free grammar does not imply that its constituents can be mixed in a context-free manner (most constituents have selection restriction on its context)
 - ◆ Normalization Issue: parse trees with less nodes get a higher score (introducing errors un-related to the linguistics characteristics)

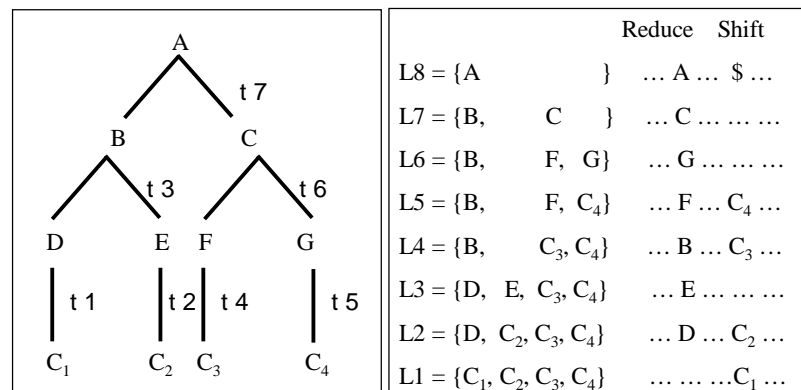
2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

21

Inappropriate Feature Dependencies (4)

- Syntactic Score [Su 88]
- $S_{\text{syn}} \equiv P(\text{Syn}, \text{Lex} \mid \text{Wrd})$
- Decomposition of Syntax Tree (into phrase levels for score computation in bottom-up GLR parser)



2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

22

Inappropriate Feature Dependencies (5)

■ Basic Context-Sensitive Formulation

$$P(L_j | L_1^{j-1}, c_1^n, w_1^n) \approx P(L_j | L_{j-1}) \\ \equiv P(\{l_A, A, r_A\} | \{l_A, X_1, X_2, \dots, X_m, r_A\})$$

- ◆ L_j : the j-th phrase level

$$L_j : \{l_A, A, r_A\} \leftarrow L_{j-1} : \{l_A, X_1, X_2, \dots, X_m, r_A\}$$

- ◆ Encode context-sensitivity within context-free framework
- ◆ Evaluated after each **reduce** action

■ Example:

$$\begin{aligned} \text{SCORE}_{\text{syn}}(\text{Syn}_A) \\ &= P(L_8, L_7, \dots, L_2 | L_1) \\ &= P(L_8 | L_7, \dots, L_1) \times P(L_7 | L_6, \dots, L_1) \times \dots \times P(L_2 | L_1) \\ &\approx P(L_8 | L_7) \times P(L_7 | L_6) \times \dots \times P(L_2 | L_1) \\ &\approx P(\{A\} | \{l_7, B, C, r_7\}) \times P(\{C\} | \{l_6, F, G, r_6\}) \times \dots \times P(\{D\} | \{l_1, c_1, r_1\}) \end{aligned}$$

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

23

Inappropriate Feature Dependencies (6)

■ Run-Time Formulation (Normalization Form)

- ◆ Compact multiple highly correlated phrase levels when evaluating score (Head-Lexicon in L_1 can be kept in the following example)
- ◆ Evaluated after each **shift** action
- ◆ Avoid "normalization problem" (same input with different number of transition probabilities)

■ Example:

$$\begin{aligned} \text{SCORE}_{\text{syn}}(\text{Syn}_A) \\ &= P(L_8, L_7, \dots, L_2 | L_1) \\ &= P(L_8, L_7, L_6 | L_1^5) \times P(L_5 | L_1^4) \times P(L_4, L_3 | L_1^2) \times P(L_2 | L_1) \\ &\approx P(L_8, L_7, L_6 | L_5) \times P(L_5 | L_4) \times P(L_4, L_3 | L_2) \times P(L_2 | L_1) \\ &\approx P(L_8 | L_5) \times P(L_5 | L_4) \times P(L_4 | L_2) \times P(L_2 | L_1) \quad [\text{transition between shifts}] \end{aligned}$$

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

24

Local Maximum Trap

- Multiple local maximums or non-unique global maximum points in the parameter space trap the searching process frequently
- Poor initial guess might cause the searching process converges to an undesired local maximum not preferred by the human
 - ◆ Causing the mismatch between two measuring functions
 - ◆ Seed corpus can be used to provide a better starting point
- Example: using a bi-gram model for part of speech tagging; however, each word in the corpus has exactly two tags (e.g., noun and verb)
 - ◆ Switching noun and verb generates the same likelihood value
 - ◆ May be trapped to the completely reversed (& the worst) candidate if not guided by human preference

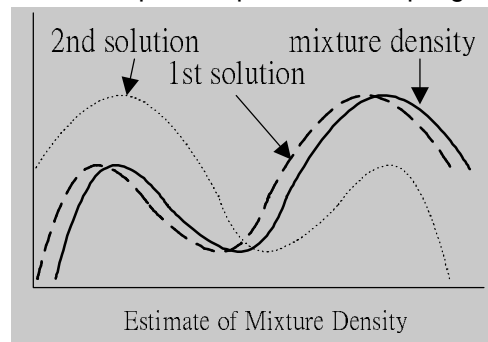
2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

25

Local Maximum Trap (Cont.)

- Complicated tasks usually have many local maximum points
 - ◆ Task complexity can be measured by the perplexity factor
 - ◆ Less chance for the unsupervised learning process converging to the desired local maximum point in complicated tasks
 - ◆ Need implicit or explicit hints if unsupervised learning is adopted
- Local Maximum Trap: Example of non-unique global maximum



2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

26

Insufficient Training Data (1)

- The size of the training samples might not be large enough to support the complexity of the adopted model
- Causing the mismatch between two sampling sources (resulted from the problem of Over Fitting)
 - ◆ Cross Entropy always increases through non-trivially refining the features (i.e., increase the dimensionality of the feature vector; thus, it also increases the number of parameters to be estimated) in the training set;
 - ◆ On the other hand, perplexity always decreases with the same procedure (as prediction capability enhanced)
 - ◆ Decreasing Modeling Error in the training set might Increase the Estimation Error in the testing set, as the size of the available training data is fixed.
 - ◆ The extra errors induced (by the increase of the estimation error) in the testing set might out run those errors to be wiped out (by the decrease of the Modeling Error)

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

27

Insufficient Training Data (2)

- Example I: increase “N” in an N-gram model
 - ◆ Increasing “N” also increases the maximum likelihood value that we can obtain in the training set
 - ◆ It also decreases the error rate in the training set under the supervised mode, as the modeling error will be reduced too (through covering wider context)
 - ◆ However, the error rate in the testing set will go up eventually if you keep increasing the “N”.

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

28

Insufficient Training Data (3)

■ Example II: Line Fitting

- ◆ Assume data are really generated from a linear model with noise independently added
- ◆ A high order polynomial function ($y = a x^{99} + b x^{98} + \dots + c x + d$) is adopted as the model
- ◆ Now trained with 3 data points:
 - ◆ Modeling error would be observed in the training set for the linear model
 - ◆ Obtain zero modeling error in the training set for any high order model (perfectly fitted by the quadratic curve of the form $y = a' x^2 + b' x + c'$)
 - ◆ BUT, the linear model enjoys smaller error in the the testing set

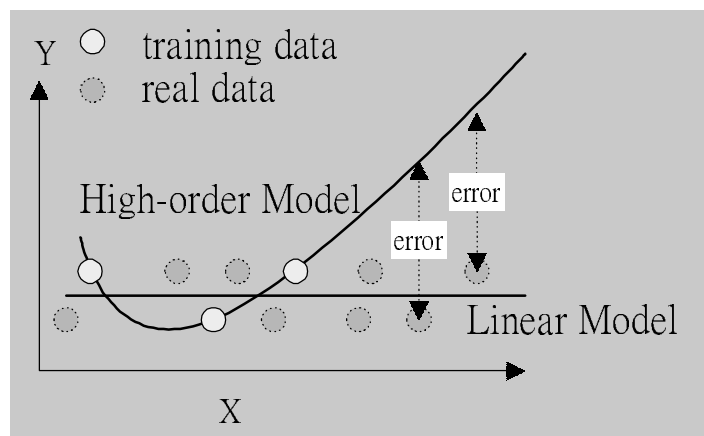
2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

29

Over Fitting: (Example - Line Fitting)

■ Training Set and Testing Set Errors in Fitting Lines



2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

30

Insufficient Training Data (4)

■ Over-Tuning Effect

- ◆ Usually, in the first few iterations, the performance in the testing set goes up; however, as the iteration keeps going, the performance in the testing set starts to deteriorate with iteration (although the associated training set likelihood value still keeps increasing)
- ◆ The main reason is that if the model tuned to fit the training set data too much, the data sampling variation between the training set and the testing set will be unveiled
- ◆ This effect is very similar to the over-tuning effect in the adaptive learning of supervised learning
- ◆ A cross-validation set can be adopted to help decide when to stop

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

31

Insufficient Training Data (5)

■ Model Resolution versus Coverage Rate in the Feature Space

- ◆ Increasing the model resolution (by increasing the model complexity, or by reducing the model covering scope) usually decreases the coverage rate in the testing set
 - ◆ Increasing the model resolution increases the discrimination power in the training set
 - ◆ However, if the local description function gets sharper, the scope that it can cover gets smaller
 - ◆ No information would be available on those uncovered regions. Thus, it would induce low coverage rate on the real data (testing set)
 - ◆ Example: Regard each word as a class (IBM first statistical MT) !

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

32

Insufficient Training Data (6)

■ Model Resolution versus Coverage Rate (Cont.)

◆ Example: Histogram and Kernel functions (data-driven approaches)

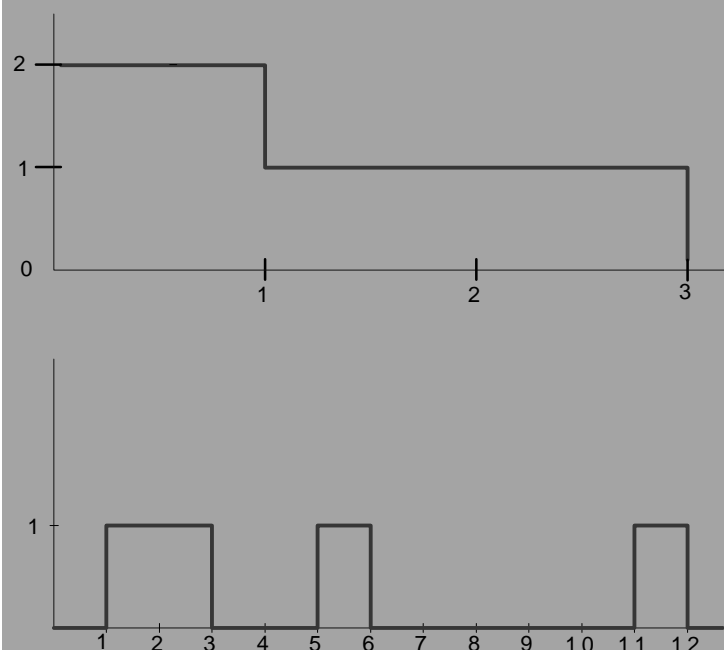
- ◆ If you divide a histogram into too many divisions, many cells will be empty (and they tell us almost nothing about the real distribution)
- ◆ Please see following figures

2002/08/18

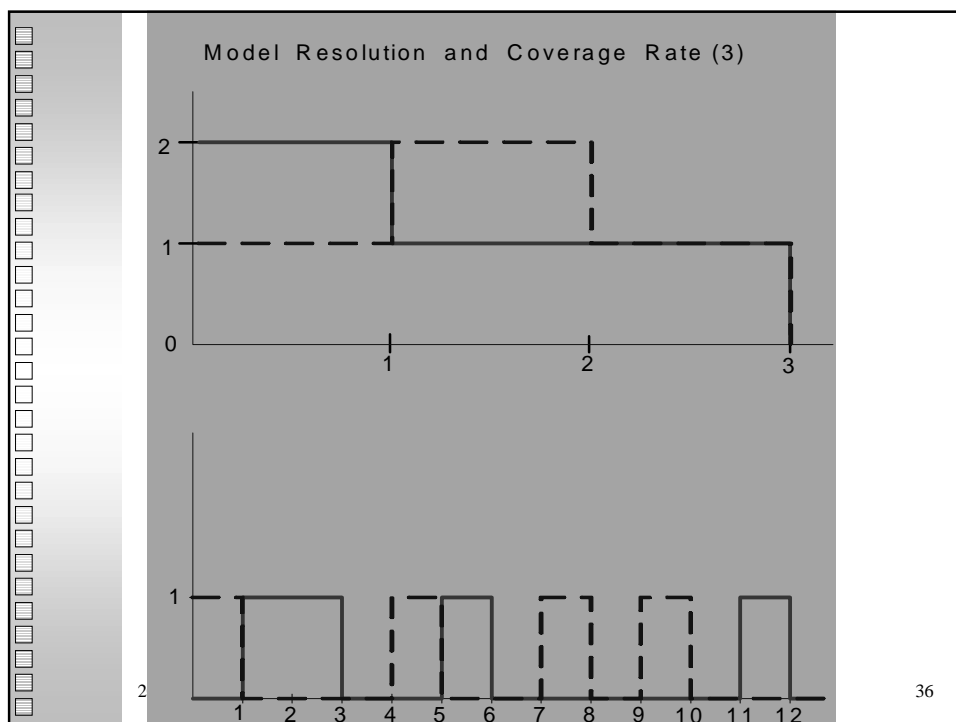
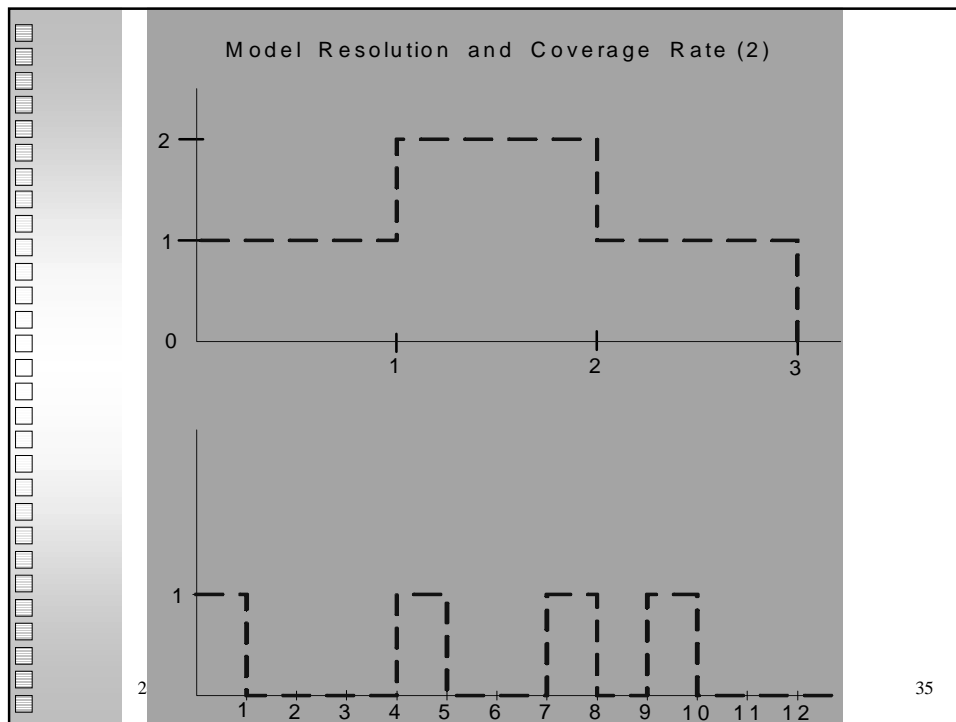
Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

33

Model Resolution and Coverage Rate (1)



34



Insufficient Training Data (7)

■ How much is enough?

- ◆ Usually 5 to 10 times is considered to be enough (i.e., similar performance will be observed in the testing set) for most applications
- ◆ However, the cases that use much less data (typically less than 1 times) to train their NLP models are not rare
- ◆ The suitable size actually depends on the problems and the models adopted
- ◆ Class-based approach and back-off smoothing can greatly relieve the adverse data sparseness phenomenon

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

37

General Trends

- Increasing Model Complexity (in the same family) always increase the likelihood in the training set
- First rising then falling of the performance curve (in the testing set) are frequently observed, if we keep increasing the model complexity
- Coverage Rate decreases while Model Complexity increases
- Coverage Rate decreases while the Corpus-size of the training set decreases

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

38

Statistic Characteristics Mismatch

- Caused by adopting the testing set with different domains or styles (via sampling from different sources/ locations, at different time, etc.)
 - ◆ Language usage is usually very dynamic in the real world (very difficult to precisely predicate every possible situation that will occur in the real applications)
 - ◆ Pre-assumed conditions rarely can hold long
- Generating the mismatch between two sampling sources
 - ◆ Mismatch between Lexicon usage statistics (mainly in domain mismatch)
 - ◆ Mismatch between other syntactic (or semantic) patterns statistics (e.g. Style)

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

39

Statistic Characteristics Mismatch (Cont.)

- Model Sensitivity (versus Characteristics Variation) is low
 - ◆ If the adopted features are invulnerable (e.g., having large inter-class distance, and small intra-class variance)
 - ◆ If the adopted estimation method is robust (e.g., adopting smoothing techniques, discarding outliers, etc.)
- Model Sensitivity usually goes up when the model complexity goes up
 - ◆ Simple is beautiful (if it can provide the similar training set performance, then it will usually deliver better testing set performance) !
 - ◆ Less parameters is better (if both give the similar training set performance)

2002/08/18

Keh-Yih Su / Jing-Shin Chang Statistical NLP D2-Part-IV

40

Methods to Reduce Mismatch Effect

■ Reduce Measuring Functions Mismatch Effect

- ◆ Adopting good language model that is closely related to the human preference model
- ◆ Adopting heuristic initial guess (or adopting seed corpus) to avoid local trap

■ Reduce Measuring Sources Mismatch Effect

- ◆ Adopting the language models and the estimation methods that are robust (i.e., insensitive) to the statistical characteristics variation (also the sampling variation) between the training set and the testing set.
 - ◆ The features that peoples really use for understanding utterance are robust (e.g., four-tones in Chinese); otherwise, he cannot be understood
- ◆ Adopting class-based approaches, if necessary, and smoothing techniques to lessen the effect caused by the finite sampling size (remember, the estimation error cannot be perceived in the training set)