

## [Day-2] Unsupervised Learning for Natural Language Processing

### (Part I: Introduction)

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-I

2

## Part I: Introduction

- What is Unsupervised Learning
  - ◆ Characteristics & differences with supervised learning
- Clustering and Hidden Markov Model
  - ◆ Introduction and examples
- Cross-Entropy for Feature Selection
  - ◆ Definition and E-Set example
- When Should Unsupervised Learning be Used
  - ◆ Problem characteristics and suitable situations for unsupervised learning
- Why Unsupervised Learning is Becoming Popular
  - ◆ Environmental factors & paradigm shift

2002/08/18

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

3

## Modes of Learning

- Supervised Learning
  - ◆ Learning from Annotated Examples
  - ◆ Advantage: capable to achieve better performance (as more information is carried by the annotation) given the same amount of training data
  - ◆ Disadvantage: human annotation is usually time-consuming and expensive
- Unsupervised Learning: Clustering, Viterbi, EM
  - ◆ Learning with Un-annotated Examples
  - ◆ Advantage: human annotation is not required
  - ◆ Disadvantage: performance achieved usually is inferior to that of supervised learning
- Bootstrapping: Co-training
  - ◆ Learning with Un-annotated Training Data, however, start from an Annotated *Seed Corpus*
  - ◆ A compromise between the supervised learning and un-supervised learning
  - ◆ Provide most cost effective solution, if used appropriately

2002/08/18

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

4

## Clustering for forming Classes

### ■ Goal:

- ◆ To collect different tokens having similar behavior into various equivalent classes to improve the robustness of the adopted statistical language model

### ■ Applications:

- ◆ Word sense disambiguation, Machine Translation, etc.

### ■ Typical Procedures

- ◆ Define feature space (e.g., frequent words occurrence vector, etc.)
- ◆ Define similarity measure (e.g., distance, angle difference, etc.)
- ◆ Cluster data iteratively according to the given criterion (e.g., minimum mean square error, maximin-distance, maximum-likelihood-value, etc.)
- ◆ Stop when desired criteria are matched.

2002/08/18

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

5

## Clustering for forming Classes (Cont.)

### ■ Commonly used Clustering Algorithms

- ◆ Dynamic Clustering: find the best K clusters (for a given K)
- ◆ Hierarchical Clustering: #clusters decrement/increment by one per iteration
  - ◆ Agglomerative (Bottom-up/Clumping approach): n singletons => K clusters
  - ◆ Divisive (Top-Down/Splitting approach): single cluster => K clusters

2002/08/18

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

6

## Dynamic Clustering

### ■ Procedure:

- ◆ Employs an iterative algorithm to optimize a clustering criterion function.
- ◆ At each iteration, data points are assigned to clusters.
- ◆ Then, the cluster representatives are updated to reflect any change in the data point assignment.
- ◆ The new cluster models are used in the next iteration.
- ◆ Continue until a stable partition is obtained.
- The number of clusters is known beforehand.

### ■ Example: K-means clustering

1. Choose the number of classes,  $K$
2. Choose initial class means:  $\mu_1, \mu_2, \dots, \mu_K$ .
3. Classify each data  $x_i$  to one of the  $K$  classes.
4. Re-compute the estimates for  $\mu_i$  using the results of 3.
5. If the  $\mu_i$  are consistent then STOP; otherwise go to (3) and (4).

2002/08/18

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

7

## Hierarchical Clustering [Bottom-Up]

### ■ Advantage: Number of Clusters need not be pre-specified

### ■ Procedure:

- ◆ Initially, every point in the data set is considered as a separate cluster.
- ◆ At any stage of a hierarchical clustering algorithm, the two of the existing clusters which are most *similar* are merged to create a new cluster, thus reducing the number of potential clusters by one.
- ◆ Terminate when the desired criteria are matched.
- The number of clusters is unknown beforehand.

### ■ Examples [Brown 92]

- ◆ Friday Monday Thursday ...
- ◆ people guys folks fellows ...
- ◆ water gas coal liquid acid ...
- ◆ man woman boy girl ...
- ◆ head body hands eyes ...

2002/08/18

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

8

## Hidden Markov Model (HMM)

- In most NLP applications, HMM is a Markov Chain with its associated state-label in each stage unknown
  - ◆ The associated state-sequence is unknown, it can only be indirectly guessed through the given observation sequence
  - ◆ Likelihood of observation: ( $\Leftrightarrow$  words in tagging task)

$$P(o_1^n | \lambda) = \sum_{s_1^n} P(o_1^n, s_1^n | \lambda) = \sum_{s_1^n} P(o_1^n | s_1^n, \lambda) P(s_1^n | \lambda)$$

$$P(s_1^n | \lambda) = \prod_i P(s_i | s_{i-1}, \lambda)$$

- ◆ Likelihood of state sequence: ( $\Leftrightarrow$  tags in tagging task)

$$\arg \max_{s_1^n} P(s_1^n | o_1^n, \lambda) = \arg \max_{s_1^n} \prod_i P(o_i | s_i, \lambda) P(s_i | s_{i-1}, \lambda)$$

[Solution: Viterbi Decoding Algorithm]

2002/08/18

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

9

## Hidden Markov Model (HMM)

- Tagging example: The associated tag-sequence ( $C_1$  to  $C_n$ ) is unknown (and it is to be found out), only the word-sequence ( $w_1$  to  $w_n$ ) is given, and

$$P(C_1^n | w_1^n, \lambda) \approx \left[ \prod_i P(w_i | C_i, \lambda) P(C_i | C_{i-1}, \lambda) \right] / P(w_1^n)$$

- Three HMM problems [Rabiner 93]:
  - ◆ How to compute output probability of observed output symbols
  - ◆ How to find the best corresponding state-sequence  $\{s_1, \dots, s_n\}$
  - ◆ How to estimate the associated parameter set  $\Lambda = (\mathbf{A}, \mathbf{B}, \Pi)$  from the given observation sequence (Unsupervised)
  - ◆ (How to efficiently resolve the above problems)

2002/08/18

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

10

## Hidden Markov Model (1)

### ■ Parameter set $\Lambda = (\mathbf{A}, \mathbf{B}, \Pi)$

- ◆ State Transition Probability Matrix  $\mathbf{A}$ :

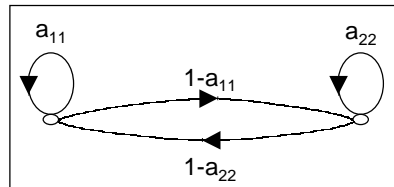
$$a_{ij} = P(q_t = j | q_{t-1} = i), \quad a_{ij} \geq 0, \sum_j a_{ij} = 1$$

- ◆ Observation Symbol Probability Matrix  $\mathbf{B}$ :

$$b_j(k) = P(o_t = k | q_t = j)$$

- ◆ Initial State Distribution Probability Vector  $\Pi$ :

$$\pi_i = P(q_1 = i)$$



$\pi_1=0.4$        $\pi_2=0.6$   
 $P(H) = P_1$      $P(H) = P_2$   
 $P(T) = 1-P_1$     $P(T) = 1-P_2$   
 $O = HHTTHTHTHTH \dots$   
 $S = 2 \ 1 \ 1 \ 2 \ 2 \ 2 \ 1 \ 2 \ 2 \ 1 \ 2 \dots$

2002/08/18

Keh-Yih Su / Jing-Shin Chang

Statistical NLP D2-Part-I

11

## Hidden Markov Model (2)

### ■ Tagging Example

$$P(C_1^n | w_1^n) = \left[ \prod_{i=1}^n P(C_i | C_{i-1}) P(w_i | C_i) \right] / P(w_1^n)$$

$$\pi_i = P(C_{1,i}), \quad a_{ij} \Leftrightarrow P(C_{t,j} | C_{t-1,i}), \quad b_j(k) \Leftrightarrow P(w_{t,j,k} = w_t | C_{t,j})$$

[subscript 't': time index]

### ■ Usually adopt Dynamic Programming and A\* Searching Algorithms [Winston 92]

- ◆ Dynamic Programming is first adopted to perform forward searching for finding the best candidate (state sequence)
- ◆ A\* algorithm is then applied backwards to get desired Top-N candidates

2002/08/18

Keh-Yih Su / Jing-Shin Chang

Statistical NLP D2-Part-I

12

## HMM Problem I (Decoding)

### - Compute Observation Probability

#### ■ Direct computation:

- ◆ Introducing hidden state variables

$$\begin{aligned} P(O | \lambda) &= \sum_Q P(O, Q | \lambda) \\ &= \sum_Q P(O | Q, \lambda) P(Q | \lambda) \end{aligned}$$

$$P(Q | \lambda) = \pi_{q_1} a_{q_1 q_2} a_{q_2 q_3} \dots a_{q_{T-1} q_T}$$

$$P(O | Q, \lambda) = b_{q_1}(o_1) b_{q_2}(o_2) \dots b_{q_T}(o_T)$$

$$P(O | \lambda) = \sum_{Q=q_1^T} \pi_{q_1} b_{q_1}(o_1) \prod_{t=1, T-1} a_{q_t q_{t+1}} b_{q_{t+1}}(o_{t+1})$$

2002/08/18

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

13

## HMM Problem II

### - Finding Best (Hidden) State Sequence

#### ■ Goal:

$$Q^* = \arg \max_Q P(Q | O, \lambda) = \arg \max_Q P(Q, O | \lambda)$$

- #### ■ Partial Solution:
- maximum probability of state sequences after seeing partial observation

$$\delta_t(i) = \max_{q_1 \dots q_{t-1}} P(q_1^{t-1}, o_1^{t-1}, q_t = i, o_t)$$

- #### ■ Recursive Relationship: between partial solutions

$$\begin{aligned} \delta_{t+1}(j) &= \max_i \delta_t(i) a_{ij} b_j(o_{t+1}) \\ \psi_{t+1}(j) &= \arg \max_i \delta_t(i) a_{ij} b_j(o_{t+1}) \end{aligned}$$

- #### ■ Termination and Tracing back best state sequence

$$q_T^* = \arg \max_i \delta_T(i) \quad q_t^* = \psi_{t+1}(q_{t+1}^*) \quad P^*(O, Q | \lambda) = \max_i \delta_T(i)$$

2002/08/18

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

14

## HMM Problem III

### - Parameter Estimation (Unsupervised, EM)

- Transition probability from state  $i$  to state  $j$  (at time  $t$ )

$$\xi_t(i, j) \equiv P(q_t = i, q_{t+1} = j | O, \lambda) = \frac{\alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}{P(O | \lambda) = \sum_s \alpha_t(s) \beta_t(s)}$$

$$\gamma_t(i) = \sum_j \xi_t(i, j), \text{ with } \alpha_t(i) \equiv P(o_1^t, q_t = i | \lambda) \text{ \& } \beta_t(i) \equiv P(o_{t+1}^T | q_t = i, \lambda)$$

- Estimating probability as expected number of transitions / probabilities: Iteratively starting from a guess on  $(\mathbf{A}, \mathbf{B}, \Pi)$

$$\bar{\pi}_i = \gamma_1(i) = \frac{\text{expected transition probability}}{\text{from state } i \text{ at time } (t = 1)}$$

$$\bar{a}_{ij} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)} = \frac{\text{expected \# transitions from state } i \text{ to } j}{\text{expected \# transitions from state } i}$$

$$\bar{b}_{ik} = \frac{\sum_{t: o_t = k} \gamma_t(i)}{\sum_{t=1}^T \gamma_t(i)} = \frac{\text{expected \# transition from state } i \text{ with output } k}{\text{expected \# transition from state } i}$$

2002

## More NLP Examples with Unsupervised Learning: Machine Translation System

$$P(T_i | S_i) = \sum_{I_i} P(T_i, I_i | S_i)$$

$$\equiv \sum_{I_i} \{ [P(T_i | PT_t(i)) \times P(PT_t(i) | NF1_t(i)) \times P(NF1_t(i) | NF2_t(i))] \cdots (1)$$

$$\times [P(NF2_t(i) | NF2_s(i))] \cdots (2)$$

$$\times [P(NF2_s(i) | NF1_s(i)) \times P(NF1_s(i) | PT_s(i)) \times P(PT_s(i) | S_i)] \cdots (3)$$

where

- ◆  $S_i$ : Source Sentence
- ◆  $T_i$ : Target Sentence
- ◆  $I_i$ : intermediate forms for the source-target pair
- ◆ PT: parse tree (s: source & t: target)
- ◆ NF1: level-1 normal form
- ◆ NF2: level-2 normal form
- ◆ (1) generation score (2) transfer score (3) analysis score

2002/08/18

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

16



## Cross-Entropy for Feature Selection

- In supervised learning, the error rate is usually the criterion for feature selection, which is, however, not available in unsupervised learning in many cases
- Other highly correlated measures (with the error rate) should be adopted
  - ◆ Characteristics of good features: possess large inter-cluster distance and with small intra-cluster variance
  - ◆ Or, in other words: with great probabilistic distribution mis-match (less overlapping portion in probabilistic distribution)
  - ◆ Cross-Entropy (also known as Divergence) is a probabilistic distribution mis-match measure [Tou&Gonzalez 74, and Blahut 87]

2002/08/18

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

17

## Cross-Entropy for Feature Selection (Cont.)

- Cross-Entropy Definition
  - ◆ Cross-Entropy (Discriminating Information, Kullback-Leibler Distance):
$$L(\mathbf{q}_0; \mathbf{q}_1) = \sum_{k=1}^K q_{0k} \log \frac{q_{0k}}{q_{1k}} = \sum_{k=1}^K q_{0k} l(k)$$
    - ◆ Expected value of log-likelihood ratio between two sets of probability vectors of observing features from two classes  $C_0$  and  $C_1$
    - ◆  $k$ : an observed feature,  $l(k)$ : log-likelihood ratio for feature  $k$
    - ◆  $q_{0k}$ : the probability that  $k$  is from  $C_0$ ,  $q_{1k}$ : the probability that  $k$  is from  $C_1$
  - ◆ Discrimination Information always increases *via* adding nontrivial features
    - ◆ The discrimination Information of each feature can be evaluated with the aid of a seed corpus (or a cross-validation set)
  - ◆ Symmetrical form is known as the Divergence:  $L(\mathbf{q}_0; \mathbf{q}_1) + L(\mathbf{q}_1; \mathbf{q}_0)$

2002/08/18

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

18

## Cross-Entropy Example: E-Set Recognition [Su&Lee 94]

- E-Set: {b, c, d, e, g, p, t, v, z}, total nine English letters with the same ending vowel "E" (a confusion set).
- Speech Recognition: conducted in a multi-speaker, isolated-word mode.
- Training Set: 900 tokens from 100 speakers
- Testing Set: another 900 tokens from the same 100 speakers
- Each letter is modeled by a 5-state left-to-right HMM
- HMM Baseline Recognition Rates: 61.7% for the testing set, and 80.2% for the training set.

2002/08/18

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

19

## Cross-Entropy Example (1)

- Input observation vector  $X$ :
  - ◆ where  $S_{ij}$  is either the averaged or the accumulated state log-likelihood (score) for state  $j$  of word  $i$ . There are 45 elements in this *score vector*.

$$X^T = [S_1^1, S_2^1, S_3^1, \dots, S_3^9, S_4^9, S_5^9]$$

- Subspace Projection:

- ◆ Maximin Algorithm is proposed for feature selection:

- ◆ Divergence  $D_{ij}(k)$  between the class  $i$  and the class  $j$ , for  $1 \leq i, j \leq 9; i < j$ ,  $1 \leq k \leq 45$ .

$$D_{ij}(k) = \int_{Y_k} (P_i(Y_k) - P_j(Y_k)) \cdot \ln \left( \frac{P_i(Y_k)}{P_j(Y_k)} \right) dY_k$$

- ◆ Find  $D_{\min}(k) = \min_{i,j} D_{ij}(k)$  for  $1 \leq i, j \leq 9; i < j$ ,  $1 \leq k \leq 45$ .
- ◆ Sort  $D_{\min}(k)$  in descending order, then the sequence of subspaces is obtained

2002/08/18

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

20

## Cross-Entropy Example (2)

**TABLE I**

THE 45 FEATURE INDICES LISTED IN THE ORDER OF DESCENDING DIVERGENCE VALUES (LEFT TO RIGHT FIRST, THEN TOP TO BOTTOM)

1 (B <sub>1</sub> )	6 (C <sub>1</sub> )	31 (T <sub>1</sub> )	41 (Z <sub>1</sub> )	26 (P <sub>1</sub> )
21 (G <sub>1</sub> )	17 (E <sub>2</sub> )	37 (V <sub>2</sub> )	36 (V <sub>1</sub> )	16 (E <sub>1</sub> )
22 (G <sub>2</sub> )	18 (E <sub>3</sub> )	11 (D <sub>1</sub> )	38 (V <sub>3</sub> )	2 (B <sub>2</sub> )
12 (D <sub>2</sub> )	10 (C <sub>5</sub> )	4 (B <sub>4</sub> )	28 (P <sub>3</sub> )	43 (Z <sub>3</sub> )
30 (P <sub>5</sub> )	20 (E <sub>5</sub> )	34 (T <sub>4</sub> )	27 (P <sub>2</sub> )	7 (C <sub>2</sub> )
32 (T <sub>2</sub> )	19 (E <sub>4</sub> )	14 (D <sub>4</sub> )	33 (T <sub>3</sub> )	45 (Z <sub>5</sub> )
13 (D <sub>3</sub> )	9 (C <sub>4</sub> )	40 (V <sub>5</sub> )	39 (V <sub>4</sub> )	42 (Z <sub>2</sub> )
25 (G <sub>5</sub> )	15 (D <sub>5</sub> )	24 (G <sub>4</sub> )	8 (C <sub>3</sub> )	29 (P <sub>4</sub> )
5 (B <sub>5</sub> )	3 (P <sub>3</sub> )	23 (G <sub>3</sub> )	35 (T <sub>5</sub> )	44 (Z <sub>4</sub> )

2002/08/18

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

21

## Cross-Entropy Example (3)

**TABLE II**

THE 45 CORRESPONDING DIVERGENCE VALUES LISTED IN DESCENDING ORDER

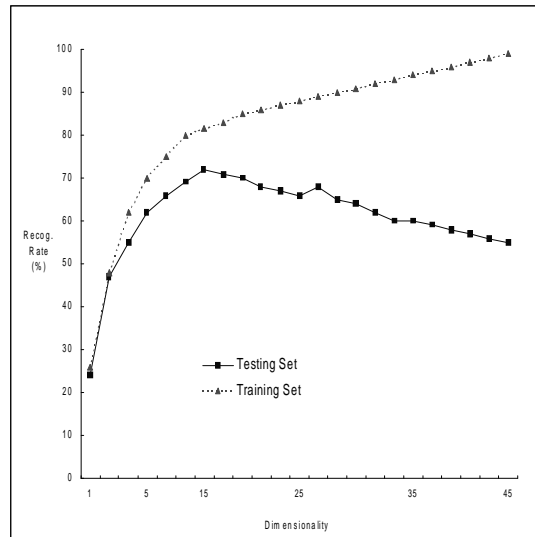
0.099410	0.066085	0.042937	0.038249	0.037496
0.028690	0.020272	0.019982	0.017487	0.011872
0.009741	0.005572	0.005305	0.004210	0.004037
0.004005	0.003803	0.003728	0.003656	0.003621
0.003214	0.003200	0.003085	0.002916	0.002690
0.002147	0.002064	0.001937	0.001900	0.001661
0.001651	0.001637	0.001496	0.001221	0.000969
0.000969	0.000596	0.000562	0.000384	0.000381
0.000229	0.000194	0.000187	0.000076	0.000042

2002/08/18

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

22

## Subspace-based Recognizer: Recognition Rate vs. Dimensionality



2002/08/18

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

23

## Characteristics of Unsupervised Learning

- Error rate of training set is not available
  - ◆ Other object measure is needed to guide both feature selection and parameter searching processes
- Human preference is not unveiled
  - ◆ Criterion mismatch is more serious
- The form of Language model is more critical
  - ◆ It would favor the case with less probability terms from iteration to iteration
- Many local maximums exist even for likelihood measure
  - ◆ The corresponding form is usually not in the regular exponential family
  - ◆ Local trap in searching space is more severe

2002/08/18

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

24

## Characteristics of Unsupervised Learning (Cont.)

- Any constraint would significantly help
  - ◆ Speech recognition: the syllables state segmentation sequence is unknown however, the corresponding text string is known
  - ◆ Sentence alignment: each alignment-passage is unknown, however, two documents are known to correspond to each other
  - ◆ Machine Translation: the detailed structure mapping is unknown, however, two sentences are known to be the corresponding translation.
- Learning efficiency is relatively lower comparing with supervised learning
  - ◆ A larger corpus is required to achieve the similar performance (if it is possible) of the supervised learning, as it is operating without the information implied by the annotation which is available in supervised learning

2002/08/18

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

25

## When to Use Unsupervised Learning (1)

- Problem Characteristics
  - ◆ The task at hand requires huge amount of fine grained knowledge to achieve acceptable performance
    - ◆ Usually imply that providing supervised learning examples might not be affordable
  - ◆ Inherent Constraints (or implied dependency) among the linguistic units are strong
    - ◆ Have a better chance to predict the best candidate through good language model
    - ◆ Bilingual corpus can help for imposing constraint
  - ◆ Training data have enough explicit Inherent anchor points
    - ◆ Help to impose contextual constrains on their neighbors (e.g., unambiguous words in tagging part-of-speech, paragraph markers in sentence alignment)
    - ◆ Make the task for resolving ambiguity easier

2002/08/18

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

26

## When to Use Unsupervised Learning (2)

### ■ Problem Characteristics (Cont.)

- ◆ Annotation Task is difficult to proceed (highly confusing in assigning labels)
  - ◆ For example, part of speech tagging versus word sense disambiguation (even human have difficulty to assign appropriate word senses according to WordNet definition)
  - ◆ Annotation cost would be high, also the consistency and quality will be hard to maintain
  - ◆ The entropy reduction (i.e., information gain) obtained from annotating the corpus would not be large, which would imply that supervised learning might not get significantly better result
  - ◆ Therefore, annotating corpus may not be worth

2002/08/18

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

27

## When to Use Unsupervised Learning (3)

### ■ Resource Scarceness: Mass amount of un-annotated data is available

- ◆ Information concentration (advantage own by the supervised-learning) can be compensated by the extra amount of training data
- ◆ With huge amount of data, constraints can be imposed to the corpus for helping to select more certain parts (e.g., selecting the boundary sentences of paragraphs for training the sentence segmentation model)
- ◆ Implied information might be enough to cover the unobservable knowledge which would be, otherwise, directly provided in the supervised-learning case (from less amount of training data)

2002/08/18

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

28

## When to Use Unsupervised Learning (4)

### ■ Cost for Preparing Learning Samples is high

- ◆ Cost for Collecting Training Samples:
  - ◆ Have people do the data entry work (e.g., LEXIS-NEXIS)
  - ◆ From public resources (e.g., LCD, ROCLING, etc.)
  - ◆ From the web/news/bbs with automatic tools
- ◆ Cost for Annotating Training Samples :
  - ◆ This is usually the bottleneck for supervised learning: requiring number of qualified persons for annotating the corpus and doing consistency check; besides, it also requires a long period of time for large scale projects
- ◆ Choose unsupervised learning if you cannot afford the cost required

2002/08/18

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

29

## When to Use Unsupervised Learning (5)

### ■ Features, domains/tasks updated frequently

- ◆ Re-annotation & re-training are frequently required, if supervised learning is adopted
  - ◆ Example 1: uncertain in classification hierarchy to be adopted
  - ◆ Example 2: uncertain in discriminative features to be adopted
- ◆ Will dimensionality, values, or labels of the feature vector change frequently?
- ◆ Do you expect the tasks or requirements (e.g., US/NIST MT Evaluation) to keep changing in each phase?
- ◆ Unsupervised-learning is prefer if the frequency is high

2002/08/18

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

30

## When to Use Unsupervised Learning (6)

- Size of available resource keep increasing with time:
  - ◆ System parameters could be updated frequently for incremental improvement, which would be difficult for supervised learning
- Good Language Model that can echo the human preference is available
  - ◆ Without the annotation to unveil the human preference (as the supervised learning possesses), the human preference must be implicitly implied in the given language model
  - ◆ With a good language model, it would have a better chance to learn language parameters that are capable to achieve satisfactory performance
- In summary: when the projected cost for annotating corpus to achieve the desired performance is high, and it is expected that unsupervised learning can achieve competitive performance by adopting better language model, choose unsupervised learning

2002/08/18

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

31

## Why unsupervised learning is getting popular in NLP (1)

- The fact that NLP requires huge and fine-grained knowledge and the infeasibility to annotate a big corpus is increasingly perceived
- In general, better performance requires deeper analyses; however, the annotating task gets more and more difficult when the analysis gets deeper
  - ◆ The increase of inherent non-determinism make the task of assigning tags more difficult.
- Many public corpora only provide minimum degree of annotation (partially anchored); the cost for further annotation is beyond the reach of most researchers
  - ◆ The environment and supports for supervised learning is limited

2002/08/18

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

32



## Why unsupervised learning is getting popular in NLP (2)

- On the other hand, the cost for possessing an un-annotated corpus diminishes to almost nothing
  - ◆ Almost free through resource sharing (e.g., LDC, ROCLING, etc.) or through acquiring from WWW; however, annotated corpora are still rare
- The size of on-line corpora increases rapidly in Internet age
  - ◆ The degree of knowledge concentration is no more essential.
  - ◆ Multi-lingual corpus is more available
    - ◆ Implicit constraints and implied annotations can reduce the degree of ambiguity
    - ◆ WWW provides an abundant resource

2002/08/18

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

33

## Why unsupervised learning is getting popular in NLP (3)

- Corpus-Based Statistic-Oriented (CBSO) approaches prevail in NLP community; however, it is a rapidly changing field
  - ◆ New model is tried in a fast pace (an exciting field)
  - ◆ New features, classes are tested and refined rapidly
  - ◆ Very difficult to keep the associated annotation updated accordingly, if the supervised-learning approach is adopted
- New applications emerge, and the capability to be customizable and self-learnable is increasingly emphasized
  - ◆ Annotating associated corpora for various domains is usually not affordable

2002/08/18

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

34