

Word Sense Disambiguation

Hsu Ting-Wei

Outline

- **Introduction**
 - 7.1 Methodological Preliminaries
 - 7.1.1 Supervised and Unsupervised learning
 - 7.1.2 Pseudowords
 - 7.1.3 Upper and lower bounds on performance
- **Methods for Disambiguation**
 - 7.2 Supervised Disambiguation
 - 7.2.1 Bayesian classification
 - 7.2.2 An information-theoretic approach
 - 7.3 Dictionary-based Disambiguation
 - 7.3.1 Disambiguation based on sense definitions
 - 7.3.2 Thesaurus-based disambiguation
 - 7.3.3 Disambiguation based on translations in a second-language corpus
 - 7.3.4 One sense per discourse, one sense per collocation
 - 7.4 Unsupervised Disambiguation

Introduction

- The task of disambiguation is to determine which of the **senses** of an ambiguous word is invoked in a particular use of the word.
- This is done by **looking at the context** of the word's use.
- Ex: The word “**bank**”, some senses that we found in a dictionary were:
 - bank [1, noun]: the rising ground bordering a lake, river, or sea... (岸)
 - bank [2, verb]: to heap or pile in a bank (築堤防護)
 - bank [3, noun]: an establishment for the custody, loan, or exchange of money (銀行)
 - bank [4, verb]: to deposit money (存錢)
 - bank [5, noun]: a series of objects arranged in a row (排;組)
- Reference: Webster’s Dictionary online <http://www.m-w.com>

Introduction (cont.)

- Two ambiguity in a sentence :
 - Tagging
 - Most part of speech tagging models simply use local context (nearby structure)
 - Word sense disambiguation
 - Word sense disambiguation methods often try to use context words in a broader context
- Ex: You should *butter* your toast.
 - Tagging
 - The word “butter” can be a 名詞 or 動詞
 - Word sense disambiguation
 - The word “butter” can mean 塗奶油 or 說甜言蜜語

7.1 Methodological Preliminaries

- 7.1.1 Supervised and Unsupervised learning
 - Supervised learning (**classification**、**function-fitting**)
 - The actual status for each piece of data on which we train
 - One extrapolates the shape of a function based on some data points.
 - Unsupervised learning (**clustering task**)
 - We don't know the classification of the data in the training sample

7.1 Methodological Preliminaries (cont.)

- 7.1.2 Pseudowords
 - Hand-labeling is a **time intensive** and **laborious task**
 - Test data are hard to come by
 - It is often convenient to generate artificial evaluation data for the comparison and improvement of text processing algorithms
 - Artificial ambiguous words can be created by conflating two or more natural words
 - Ex: banana-door
 - Easy to create large-scale train/test set

7.1 Methodological Preliminaries (cont.)

- 7.1.3 Upper and lower bounds on performance
 - It's meaningless that only consider numerical evaluation
 - Need to consider how difficult the task is
 - Using upper and lower bounds to estimate
 - **Upper bound** → human performance
 - We can't expect an automatic procedure to do better
 - **Lower bound (baseline)** → the simplest possible algorithm
 - Assignment of all contexts to the most frequent sense

Methods for Disambiguation

- 7.2 Supervised Disambiguation
 - Disambiguation based on a labeled training set
- 7.3 Dictionary-based
 - Disambiguation based on lexical resources such as dictionaries and thesauri
- 7.4 Unsupervised Disambiguation
 - Disambiguation based on training on an unlabeled text corpora.

Notational conventions used in this chapter

Symbol	Meaning
w	an ambiguous word
$s_1, \dots, s_k, \dots, s_K$	senses of the ambiguous word w
$c_1, \dots, c_i, \dots, c_I$	contexts of w in a corpus
$v_1, \dots, v_j, \dots, v_J$	words used as contextual features for disambiguation

7.2 Supervised Disambiguation

- Training corpus: Each occurrence of the ambiguous word w is annotated with a semantic label
- Supervised disambiguation is a classification task.
- We will look at:
 - **Bayesian classification** (Gale et al. 1992).
 - **Information-theoretic approach** (Brown et al. 1991)

7.2 Supervised Disambiguation (cont.)

- 7.2.1 Bayesian classification (Gale et al.1992)
 - The approach treats the context of occurrence as a bag of words without structure, but it integrates information from many words in the context window. (feature)
 - Bayes Decision rule
 - Decide s' if $P(s' | c) > P(s_k | c)$ for $s_k \neq s'$
 - Bayes decision rule is optimal because it minimizes the probability of error
 - Choose the class (or sense) with the highest conditional probability and hence the smallest error rate.

7.2 Supervised Disambiguation (cont.)

- 7.2.1 Bayesian classification (Gale et al.1992)
 - Computing **Posterior Probability** for Bayes Classification
 - We want to assign the ambiguous word w to the sense s' , given context c , where:

$$\begin{aligned} s' &= \arg \max P(s_k | c) \\ &= \arg \max \frac{P(c | s_k)}{P(c)} P(s_k) \\ &= \arg \max P(c | s_k) P(s_k) \\ &= \arg \max [\log \underline{P(c | s_k)} + \log P(s_k)] \end{aligned} \quad \left. \begin{array}{l} \text{Bay's Rule} \\ \text{log} \end{array} \right\}$$

7.2 Supervised Disambiguation (cont.)

- 7.2.1 Bayesian classification (Gale et al. 1992)

- Naive Bayes assumption (Gale et al. 1992)

- An instance of a particular kind of Bayes classifier

$$P(c | s_k) = P(\{v_j | v_j \text{ in } c\} | s_k) = \prod_{v_j \text{ in } c} P(v_j | s_k)$$

- Consequences of this assumption:
 - 1. Bag of words model: the structure and linear ordering of words within the context is ignored.
 - 2. The presence of one word in the bag is independent of another

7.2 Supervised Disambiguation (cont.)

- 7.2.1 Bayesian classification (Gale et al.1992)

- Decision Rule for Naive Bayes

- Decide s' if

$$s' = \arg \max_{s_k} [\log P(s_k) + \sum_{v_j \text{ in } c} \log P(v_j | s_k)]$$

- $P(v_j | s_k)$ and $P(s_k)$ are computed via Maximum-Likelihood Estimation, perhaps with appropriate smoothing, from the labeled training corpus

$$P(v_j | s_k) = \frac{C(v_j, s_k)}{\sum_t C(v_t, s_k)} , P(s_k) = \frac{C(s_k)}{C(w)}$$

7.2 Supervised Disambiguation (cont.)

- 7.2.1 Bayesian classification (Gale et al.1992)
 - Bayesian disambiguation algorithm

```
1 comment: Training
2 for all senses  $s_k$  of w do
3     for all words  $v_j$  in the vocabulary do
4          $P(v_j|s_k) = \frac{C(v_j, s_k)}{C(v_i)}$ 
5     end
6 end
7 for all senses  $s_k$  of w do
8      $P(s_k) = \frac{C(s_k)}{C(w)}$ 
9 end
10 comment: Disambiguation
11 for all senses  $s_k$  of w do
12      $\text{score}(s_k) = \log P(s_k)$ 
13     for all words  $v_j$  in the context window c do
14          $\text{score}(s_k) = \text{score}(s_k) + \log P(v_j|s_k)$ 
15     end
16 end
17 choose  $s' = \arg \max_{s_k} \text{score}(s_k)$ 
```

7.2 Supervised Disambiguation (cont.)

- **7.2.1 Bayesian classification (Gale et al.1992)**

- Example of Bayesian disambiguation algorithm
 - $w = \text{drug}$

Sense ($s_1..s_k$)	Clues for sense ($v_1..v_j$)
Medication	prices, prescription, patent, increase, consumer, pharmaceutical
Illegal substance	abuse, paraphernalia, illicit, alcohol, cocaine, traffickers

$$P(\text{prices}|\text{'medication'}) > P(\text{price}|\text{'illicit substance'})$$

$$P(v_j | s_k) = \frac{C(v_j, s_k)}{\sum_t C(v_t, s_k)}$$

- Bayes Classifier uses information from all words in the context window by using an **independence assumption**
 - Unrealistic independence assumption

7.2 Supervised Disambiguation (cont.)

- **7.2.2 An information-theoretic approach (Brown et al. 1991)**

- The approach looks at only one informative feature in the context, which may be sensitive to text structure. But this feature is carefully selected from a large number of potential “informants.”

Ambiguous word	Indicator	Examples: value → sense	
prendre	object	measure	→ to take
		decision	→ to make
vouloir	tense	present	→ to want
		conditional	→ to like
cent	word to the left	per	→ %
		number	→ c.[money]

Highly informative indicators for three ambiguous French words

French

English

- Prendre une mesure → take a measure
- Prendre une décision → make a decision

7.2 Supervised Disambiguation (cont.)

- **7.2.2 An information-theoretic approach (Brown et al. 1991)**

- **Flip-Flop Algorithm** (Brown et al., 1991)

- The algorithm is used to disambiguate between the different senses of a word using the **mutual information** as a measure.
 - The algorithm is an efficient **linear-time** algorithm for computing the best partition of values for a particular indicator.
 - Categorize the informant (contextual word) as to which sense it indicates.

```
1 find random partition  $P = \{P_1, P_2\}$  of  $\{t_1, \dots, t_m\}$ 
2 while (improving) do
3     find partition  $Q = \{Q_1, Q_2\}$  of  $\{x_1, \dots, x_n\}$ 
4         that maximizes  $I(P; Q)$ 
5     find partition  $P = \{P_1, P_2\}$  of  $\{t_1, \dots, t_m\}$ 
6         that maximizes  $I(P; Q)$ 
7
8 end
```

t_1, \dots, t_m be the translation of the ambiguous word
 x_1, \dots, x_n the possible values of the indicator

$$I(X; Y) = \sum_{x \in X} \sum_{y \in Y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}$$

7.2 Supervised Disambiguation (cont.)

- **7.2.2 An information-theoretic approach (Brown et al. 1991)**
 - Flip-Flop Algorithm
 - Example
 - $P = \{t_1, \dots, t_m\} = \{\text{take, make, rise, speak}\}$
 - $Q = \{x_1, \dots, x_n\} = \{\text{mesure, note, exemple, decision, parole}\}$
 - Step1
 - Initial: find random partition P
 - » $P1 = \{\text{take, rise}\}, P2 = \{\text{make, speak}\}$
 - Step2
 - Find partition Q of the indicator values would give us maximum $I(P; Q)$
 - » $Q1 = \{\text{measure, note, exemple}\}, Q2 = \{\text{decision, parole}\}$
 - Repartition P and also maximum $I(P; Q)$
 - » $P1 = \{\text{take}\}, P2 = \{\text{make, rise, speak}\}$
 - If improving repeat step2

7.3 Dictionary-based Disambiguation

- If we have no information about the sense categorization of specific instance of a word, we can fall back on a general characterization of the senses.
- Sense definitions are extracted from existing sources such as **dictionaries** and **thesaurus** (類語詞典，表示一串有著相似、相關或相反意思的單字)
- The different types of informational method have been used:
 - 7.3.1 Disambiguation based on sense definitions
 - 7.3.2 Thesaurus-based disambiguation
 - 7.3.3 Disambiguation based on translations in a second-language corpus
 - 7.3.4 One sense per discourse, one sense per collocation

7.3 Dictionary-based Disambiguation (cont.)

- 7.3.1 Disambiguation based on sense definitions (Lesk, 1986)
 - A word's dictionary definitions are likely to be good indicators of the senses they define.
 - Lesk's dictionary-based disambiguation algorithm
 - Ambiguous word w
 - Senses of $w : S_1 \dots S_k$ (bags of words)
 - Dictionary definition of senses : $D_1 \dots D_k$
 - $E_{v_j} : \text{the set of words occurring in the dictionary definition } (D_1 \dots D_k) \text{ of word } v_j$ (bags of words)

```
1 comment: Given: context c
2 for all senses  $s_k$  of  $w$  do
3   score( $s_k$ ) = overlap( $D_k, \bigcup_{v_j \in c} E_{v_j}$ )
4 end
5 choose  $s'$  s.t.  $s' = \operatorname{argmax}_{S_k} \text{score}(s_k)$ 
```

7.3 Dictionary-based Disambiguation (cont.)

- **7.3.1 Disambiguation based on sense definitions (Lesk, 1986)**
 - Lesk's dictionary-based disambiguation algorithm

- Ex: Two senses of *ash*

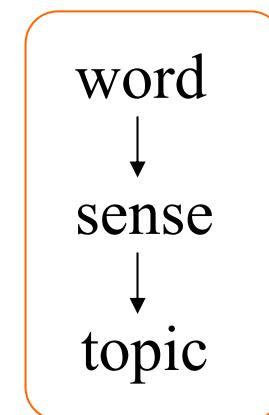
Sense	Definition
S_1 tree	a tree of the olive family
S_2 burned stuff	The solid residue left when combustible material is burned

Scores		Context
S_1	S_2	
0	1	This cigar burns slowly and creates a stiff ash
1	0	The ash is one of the last trees to come into leaf

7.3 Dictionary-based Disambiguation (cont.)

- 7.3.2 Thesaurus-based disambiguation
 - Simple thesaurus-based algorithm (Walker, 1987)
 - Each word is assigned one or more subject codes in the dictionary
 - If the word is assigned several subject codes, then we assume that they corresponds to the different senses of the word.
 - $t(s_k)$ is the subject code of sense s_k
 - $\delta(t(s_k), v_j) = 1$ iff $t(s_k)$ is one of the subject codes of v_j and 0 otherwise

```
1 comment: Given: context c  
2 for all senses  $s_k$  of w do  
3   score( $s_k$ ) =  $\sum_{v_j \text{ in } c} \delta(t(s_k), v_j)$   
4 end  
5 choose  $s'$  s.t.  $s' = \operatorname{argmax}_{s_k} \text{score}(s_k)$ 
```



7.3 Dictionary-based Disambiguation (cont.)

- 7.3.2 Thesaurus-based disambiguation
 - Simple thesaurus-based algorithm (**Walker, 1987**)
 - Problem
 - A general categorization of words into topics is often inappropriate for a particular domain
 - » Mouse → mammal, electronic device
 - » When in a computer manual
 - A general topic categorization may also have a problem of coverage
 - » Navratilova (網球運動員) → sports
 - » When “Navratilova” is not found in the thesaurus.....

7.3 Dictionary-based Disambiguation (cont.)

- 7.3.2 Thesaurus-based disambiguation
 - Adaptation thesaurus-based algorithm (Yarowsky, 1987)
 - Adapted the algorithm for words that do not occur in the thesaurus but that are very informative
 - Using Bayes classifier for both adaptation and disambiguation

```
1 comment: Categorize contexts based on categorization of words
2 for all contexts  $c_i$  in the corpus do
3   for all thesaurus categories  $t_l$  do
4     score( $c_i, t_l$ ) =  $\log \frac{P(c_i|t_l)}{P(c_i)} P(t_l)$ 
5   end
6 end
7  $t(c_i) = \{t_l | \text{score}(c_i, t_l) > \alpha\}$ 
8 comment: Categorize words based on categorization of contexts
9 for all words  $v_j$  in the vocabulary do
10    $V_j = \{c | v_j \text{ in } c\}$ 
11 end
12 for all topics  $t_l$  do
13    $T_l = \{c | t_l \in t(c)\}$ 
14 end
15 for all words  $v_j$ , all topics  $t_l$  do
16    $P(v_j|t_l) = |V_j \cap T_l| / \sum_j |V_j \cap T_l|$ 
17 end
18 for all topics  $t_l$  do
19    $P(t_l) = (\sum_j |V_j \cap T_l|) / (\sum_l \sum_j |V_j \cap T_l|)$ 
20 end
21 comment: Disambiguation
22 for all senses  $s_k$  of  $w$  occurring in  $c$  do
23   score( $s_k$ ) =  $\log P(t(s_k)) + \sum_{v_j \text{ in } c} \log P(v_j|t(s_k))$ 
24 end
25 choose  $s'$  s.t.  $s' = \text{argmax}_{s_k} \text{score}(s_k)$ 
```

7.3 Dictionary-based Disambiguation (cont.)

- 7.3.3 Disambiguation based on translations in a second-language corpus (Dagan et al. 1991; Dagan and Itai 1994)
 - Words can be disambiguated by looking at how they are translated in other languages
 - This method uses of word correspondences in a bilingual dictionary
 - First Language
 - The one for which we want to disambiguation
 - Second Language
 - Target language in the bilingual dictionary
 - For example, if we want to disambiguate English based on German corpus, then English is the 1st language, and the German is the 2nd language.

7.3 Dictionary-based Disambiguation (cont.)

- **7.3.3 Disambiguation based on translations in a second-language corpus (Dagan et al. 1991; Dagan and Itai 1994)**
 - Ex : w = *interest*
 - To disambiguate the word “interest”, we identify the **phrase** it occurs in, search a German corpus for instances of the phrase, and assign the meaning associated with the German use of the word in that phrase

	Sense 1	Sense 2
Definition	legal share (利息)	attention, concern (興趣)
Translation	<i>Beteiligung</i>	<i>Interesse</i>
English collocation	acquire an interest	show interest
Translation	<i>Beteiligung erwerben</i>	<i>Interesse zeigen</i>

7.3 Dictionary-based Disambiguation (cont.)

- **7.3.3 Disambiguation based on translations in a second-language corpus (Dagan et al. 1991; Dagan and Itai 1994)**
 - Disambiguation based on a second-language corpus
 - S is the second-language corpus
 - $T(s_k)$ is the set of possible translations of sense s_k
 - $T(v)$ is the set of possible translations of v

1 **comment:** Given: a context c in which w occurs in relation $R(w, v)$

2 **for** all senses s_k of w **do**

3 $\text{score}(s_k) = |\{c \in S | \exists w' \in T(s_k), v' \in T(v) : R(w', v') \in c\}|$

4 **end**

 轉換成第二語言的sense 與轉換過的sense相關的文字

5 choose $s' = \text{argmax}_{s_k} \text{score}(s_k)$

$R(\text{Interesse}, \text{zeigen})$ would be higher than count of $R(\text{Beteiligung}, \text{zeigen})$

7.3 Dictionary-based Disambiguation (cont.)

- 7.3.4 One sense per discourse, one sense per collocation (Yarowsky, 1995)
 - There are constraints between different occurrences of an ambiguous word within a corpus that can be exploited for disambiguation
 - **One sense per discourse**
 - The sense of a target word is highly consistent within any given document
 - **One sense per collocation**
 - **Nearby words** provide strong and consistent **clues** to the sense of a target word, conditional on relative distance, order and syntactic relationship

7.3 Dictionary-based Disambiguation (cont.)

- 7.3.4 One sense per discourse, one sense per collocation (Yarowsky, 1995)
 - Look one sense per discourse

Discourse	Initial label	Context
d_1	living	the existence of plant and animal life
	living	classified as either <i>plant</i> or animal
	?	Although bacterial and plant cells are enclosed
d_2	living	contains a varied plant and animal life
	living	the most common <i>plant</i> life
	living	slight within Arctic plant species
	factory	are protected by plant parts remaining from

? will be “living”

7.3 Dictionary-based Disambiguation (cont.)

- 7.3.4 One sense per discourse, one sense per collocation (Yarowsky, 1995)
 - **One sense per collocation** : Most sensed are strongly correlated with certain contextual **features** like other words in the same

```
1 comment: Initialization
2 for all senses  $s_k$  of w do
3    $F_k$  = the set of collocations in  $s_k$ 's dictionary definition
4 end
5 for all senses  $s_k$  of w do
6    $E_k$  =  $\emptyset$ 
7 end
8 comment: One sense per collocation
9 while (at least one  $E_k$  changed in the last iteration) do
10   for all senses  $s_k$  of w do
11      $E_k$  =  $\{c_i \mid \exists f_m : f_m \in c_i \wedge f_m \in F_k\}$ 
12   end
13   for all senses  $s_k$  of w do
14      $F_k$  =  $\{f_m \mid \forall n \neq k \frac{P(s_k | f_m)}{P(s_n | f_m)} > \alpha\}$ 
15   end
16 end
17 comment: One sense per discourse
18 for all documents  $d_m$  do
19   determine the majority sense  $s_k$  of w in  $d_m$ 
20   assign all occurrences of w in  $d_m$  to  $s_k$ 
21 end
```

F_k contains characteristic collocations.
 E_k is the set of contexts of the ambiguous word w that are currently assigned to s_k .

Collocational features are ranked according to the ratio (similar with information-theoretic method 7.2.2)

This is a surprisingly good performance given that the algorithm does not need a labeled set of training examples.

7.4 Unsupervised Disambiguation

- Cluster the contexts of an ambiguous word into a number of groups
- Discriminate between these groups without labeling them
 - Probabilistic model is same the same with section 7.2.1
 - Word : w
 - Senses : $s_1 \dots s_k$
 - Estimate $P(v_j|s_k)$,
 - In contrast to Gale et al.'s Bayes classifier, parameter estimation in unsupervised disambiguation is not based on a labeled training set.
 - Instead, we start with a random initialization of the parameters $P(v_j|s_k)$. The $P(v_j|s_k)$ are then reestimated by the EM algorithm.

7.4 Unsupervised Disambiguation (cont.)

- EM Algorithm
 - Learning a word sense clustering.
 - K : number of desired senses
 - $c_1, \dots, c_i, \dots, c_I$ are the contexts of the ambiguous word in the corpus
 - $v_1, \dots, v_j, \dots, v_J$ are the words being used as disambiguating features
 - 1. Initialize
 - Initialize the parameters of the model μ randomly. The parameters are $P(v_j | s_k)$ and $P(s_k)$, $j = 1, 2, \dots, J$, $k = 1, 2, \dots, K$.
 - Compute the log likelihood of corpus C given the model μ as the product of the probabilities $P(c_i)$ of the individual contexts c_i (where $P(c_i) = \sum_k P(c_i | s_k) P(s_k)$) :

$$l(C|\mu) = \log \prod_{i=1}^I \sum_{k=1}^K P(c_i | s_k) P(s_k) = \sum_{i=1}^I \log \sum_{k=1}^K P(c_i | s_k) P(s_k)$$

7.4 Unsupervised Disambiguation (cont.)

- EM Algorithm
 - 2. While $I(C|\mu)$ is improving repeat:
 - **E-step:** estimate h_{ik} , the posterior probability that s_k generated c_i , as follows:

$$h_{ik} = \frac{P(s_k)P(c_i|s_k)}{\sum_{k=1}^K P(s_k)P(c_i|s_k)}$$

To compute $P(c_i|s_k)$, we make the by now familiar Naïve Bayes assumption:

$$P(c_i|s_k) = \prod_j P(v_j|s_k)^{C(v_j \text{ in } c_i)}$$

- **M-step:** Re-estimate the parameters $P(v_j|s_k)$ and $P(s_k)$ by way of maximum likelihood estimation:

$$P(v_j|s_k) = \frac{\sum_i C(v_j \text{ in } c_i) \cdot h_{ik}}{\sum_j \sum_i C(v_j \text{ in } c_i) \cdot h_{ik}}$$

Recompute the probabilities of the senses as follows:

$$P(s_k) = \frac{\sum_{i=1}^I h_{ik}}{\sum_{k=1}^K \sum_{i=1}^I h_{ik}} = \frac{\sum_{i=1}^I h_{ik}}{I}$$

Thanks for your attention !