

Word Sense Disambiguation

Shih-Hsiang Lin

Outline

- Introduction
- Methodological Preliminaries
 - Supervised and Unsupervised learning
 - Pseudowords
 - Upper and lower bounds on performance
- Methods for Disambiguating
 - Supervised Disambiguation
 - Dictionary-based
 - Unsupervised Disambiguation

Before we starting.....

- 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)

Introduction

- Because of many words have **several meaning or senses**
 - there is ambiguity about how they are to be interpreted
- The task of **disambiguation** is to **determine** which of **the senses** of an ambiguous word is invoked in a particular use of the word.
- Types of problem
 - Syntactic ambiguity
 - differences in syntactic categories
 - Semantic ambiguity
 - homonymy(同形/音異義)or polysemy(一詞多義)

Methodological Preliminaries

Supervised and Unsupervised learning

- Supervised learning (**classification** 、 **function-fitting**)
 - know the actual status for each piece of data on which we learn
 - each element in a training set is paired with an acceptable response
- Unsupervised learning (**clustering**)
 - we don't know the classification of the data in the training sample
 - adjusts through direct confrontation with new experiences (**self organization**)

Methodological Preliminaries

Pseudowords

- In order to test the performance of algorithms on a natural ambiguous word
 - a large number of occurrence has to be disambiguated by hand ← **time-intensive laborious task**
- Generate artificial evaluation data
 - pseudowords can be created by **conflating two or more natural words**
 - create pseudoword **banana-door** and replaces all occurrence of **banana** and **door** in the corpus
- Easy to create large-scale train/test set

Methodological Preliminaries

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**
 - \aleph We can't expect an automatic procedure to do better
 - Lower bound → assign all contexts to the **most frequent sense**
 - A way to make sense of performance figures
 - A good idea for those which have no standardized evaluation sets for comparing systems

Methods for Disambiguating

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

Notational conventions

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

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)

Bayesian Classification

- 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.

Computing Posterior Probability for Bayes Classification

- We want to assign the ambiguous word w to the sense s' , given context c , where:

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

*Each context word contributes potentially useful information about which sense of the ambiguous word is likely to be used with it

Naive Bayes (Gale et al. 1992)

- An instance of a particular kind of Bayes classifier
- **Naive Bayes assumption**: The attributes (contextual words) used for description are all conditionally independent

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

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

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

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

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      $P(v_j|s_k) = \frac{C(v_j, s_k)}{C(v_j)}$ 
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  $s' = \arg \max_s \text{score}(s_k)$ 
```

Example of Bayesian disambiguation algorithm

Sense	Clues for sense
Medication	prices, prescription, patent, increase, consumer, pharmaceutical
Illegal substance	abuse, paraphernalia, illicit, alcohol, cocaine, traffickers

Clues for two senses of drug used by a Bayesian classifier

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

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

An Information-Theoretic Approach

- In the information theoretic approach try to find a **single contextual feature** that reliably indicates which sense of the ambiguous word is being used

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

Prendre une decision → make a decision | Prendre une mesure → take a measure

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

- The **Flip-Flop** algorithm is used to disambiguate between the different senses of a word using the mutual information as a measure.
- 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
s end
```

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

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

Example of Classification based on Information-Theoretic Approach

- $P = \{t_1, \dots, t_m\} = \{\text{take, make, rise, speak}\}$
 $Q = \{x_1, \dots, x_n\} = \{\text{measure, note, exemple, decision, parole}\}$
- Initial: find random partition P
 - $P1 = \{\text{take, rise}\}$, $P2 = \{\text{make, speak}\}$
- 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

Dictionary-Based Disambiguation

- If we have no information about the sense categorization of a word
 - Relying on the senses in dictionaries and thesauri.
- Sense definitions are extracted from existing sources such as **dictionaries** and **thesauri**(同屬詞典)
- Use distributional properties to improve disambiguation
 - Ambiguous words are only used with one sense in any given discourse and with any given collocate

Disambiguation Based on Sense Definitions (Lesk, 1986)

- A word's dictionary definitions are likely to be good indicators of the senses they define.
- The algorithm:
 - Given a context c for a word w with senses s_1, \dots, s_k
 - Find the bags of words corresponding to each sense s_k in the dictionary (s_k bags of words).
 - Compare with the bag of words formed by combining the context word definitions. Pick the sense which gives maximum overlap with this bag

Example of Disambiguation Based on Sense Definitions

```

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

```

D_1, \dots, D_k the dictionary definitions of the senses S_1, \dots, S_k of the ambiguous word w , represented as the bag of words occurring in the definition.

v_j is the word occurring in the context c of w
 E_{v_j} is the dictionary definition of v_j (union of all the sense definitions of v_j)

	Sense	Definition
Two senses of ash	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

Thesaurus-Based Disambiguation (Walker, 1984)

- The **semantic categories** of the words in a context determine the semantic category of the context as a whole.
 - decide the semantic category of the context
 - then decide which word sense are used
- Each word is assigned one or more subject codes which corresponds to its different meanings
- For each subject code, we count the number of words (from the context) having the same subject code. We select the subject code corresponding to the highest count

Thesaurus-Based Disambiguation (cont.)

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

$t(s_k)$ 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

The score is the number of words that are compatible with the subject code of sense s_k

Problem : A general categorization of words into topics is often inappropriate for a particular domain

Mouse \rightarrow mammal, electronic device

A general topic categorization may also have a problem of coverage

Navratilova \rightarrow sports

Thesaurus-Based Disambiguation

Creating New Categories(Yarowsky, 1992)

- Add new words to a category if they occur more often than chance
- Adapted the algorithm for words that do not occur in the thesaurus but that are very Informative
 - For example *Navratilova* can be added to the sports category

Thesaurus-Based Disambiguation

Creating New Categories (cont.)

```
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      $\text{score}(c_i, t_l) = \log \frac{P(c_i|t_l)}{P(c_i)} P(t_l)$ 
5   end
6 end
7  $\mathbf{t}(c_i) = \{t_l | \text{score}(c_i, t_l) > \alpha\}$ 
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    $T_l = \{c | t_l \in \mathbf{t}(c)\}$ 
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    $\text{score}(s_k) = \log P(\mathbf{t}(s_k)) + \sum_{v_j \text{ in } c} \log P(v_j | \mathbf{t}(s_k))$ 
24 end
25 choose  $s'$  s.t.  $s' = \text{argmax}_{s'} \text{score}(s')$ 
```

Disambiguations Based on Translations (Dagan et al. 91 & 94)

- Words can be disambiguated by looking at how they are translated in other languages
- This method use 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.

Disambiguations Based on Translations (cont.)

- Example: the word “interest” has two translations in German:
 - “Beteiligung” (legal share--50% a interest in the company)
 - “Interesse” (attention, concern--her interest in Mathematics).
- To disambiguate the word “interest”, we identify the sentence 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

Disambiguations Based on Translations (cont.)

```
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    $\text{score}(s_k) = |\{c \in S \mid \exists w' \in T(s_k), v' \in T(v) : R(w', v') \in c\}|$   
4 end  
5 choose  $s' = \text{argmax}_k \text{score}(s_k)$ 
```

- Step1
 - Count the number of times that translations of the two senses of interest occur with translations of show in the second language corpus
- Step2
 - Compare the counts of the two different senses
- Step 3
 - Choose the sense that has the higher counts as a corresponding sense

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.

One Sense per Discourse, One sense per Collocation (cont.)

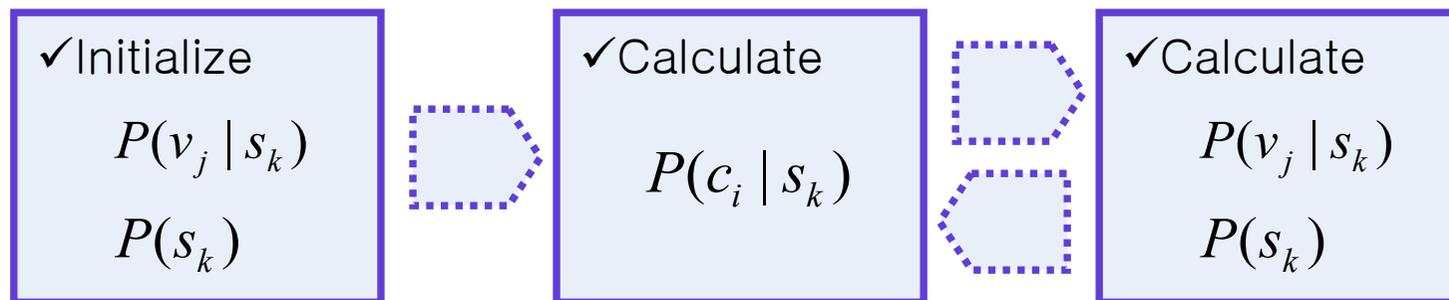
```
1 comment: Initialization
2 for all senses  $s_k$  of  $w$  do
3      $F_k =$  the set of collocations in  $s_k$ 's dictionary definition
4 end
5 for all senses  $s_k$  of  $w$  do
6      $E_k = \emptyset$ 
7 end
8 comment: One sense per collocation
9 while (at least one  $E_k$  changed in the last iteration) do
10     for all senses  $s_k$  of  $w$  do
11          $E_k = \{c_l | \exists f_m : f_m \in c_l \wedge f_m \in F_k\}$ 
12     end
13     for all senses  $s_k$  of  $w$  do
14          $F_k = \{f_m | \forall n \neq k \frac{P(s_k|f_m)}{P(s_n|f_m)} > \alpha\}$ 
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
```

One Sense per Discourse, One sense per Collocation (cont.)

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

Unsupervised Disambiguation

- Sense tagging? Sense discriminate?
- Cluster the contexts of an ambiguous word into a number of groups and discriminate between these groups without labeling them
- The probabilistic model is the Bayesian model but the $P(v_j | s_k)$ are estimated using the EM algorithm



Unsupervised Disambiguation

EM Algorithm

- **Initialize** the parameters μ of model. These 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 : $I(C|\mu) = \log \prod_i \sum_k P(c_j | s_k) P(s_k)$
- while $I(C|\mu)$ is improving repeat:
 - **E-step**: $h_{ik} = P(c_j | s_k) P(s_k) / \sum_k P(c_j | s_k) P(s_k)$ (use Naive bayes to compute $P(c_j | s_k)$)
 - **M-step**: reestimate the parameters $P(v_j | s_k)$ and $P(s_k)$ by MLE:
 $P(v_j | s_k) = \sum_{c_i} h_{jk} / Z_j$ where the sum is over all contexts c_i in which v_j occurs, Z_j a normalizing constant.
 $P(s_k) = \sum_i h_{jk} / \sum_k \sum_i h_{jk} = \sum_i h_{jk} / I$

END