

### Word Sense Disambiguation

CS-585

Natural Language Processing

Sonjia Waxmonsky

### Word Sense Disambiguation

- Many words have multiple meanings
  - E.g, river bank, financial bank
- Problem: Assign proper sense to each ambiguous word in text

- Applications:
  - Machine translation
  - Information retrieval
  - Semantic interpretation of text

### Sense Tagging

 Idea: Treat sense disambiguation like POS tagging, just with "semantic tags"

### **Distributional Similarity**

- The problems differ:
  - POS tags depend on specific structural cues (mostly neighboring tags)
  - Senses depend on semantic context less structured, longer distance dependency

### Wordnet and Synsets

- Wordnet is a manually-compiled machine-readable dictionary for English (and a few other languages), maintained by Princeton University
  - https://wordnet.princeton.edu/
- It can be used programmatically to look up word "synsets" (senses related to a set of words)

shoe.n.01: footwear shaped to fit the foot (below
the ankle) with a flexible upper of leather or
plastic and a sole and heel of heavier material

shoe

shoe.n.02: (card games) a case from which playing
cards are dealt one at a time

horseshoe.n.02: U-shaped plate nailed to underside of horse's hoof

### **Approaches**

- Dictionary-Based Learning
  - Learn to distinguish senses from dictionary entries
- Supervised learning
  - Learn from a pretagged corpus
- Unsupervised Learning
  - Automatically cluster word occurrences into different senses
  - "Clustering": Partioning of datapoints into related groups or clusters

### WSD Evaluation

- Train and test on pretagged texts is difficult to obtain
- Pseudowords Artificial data: 'merge' two words to form an 'ambiguous' word with two 'senses'
- Example, replace all occurrences of "door" and of "banana" with "doorbanana" and see if the system figures out which is which

"The jasmine, almond, **doorbanana**, cork and coco-nut palm are among the trees"

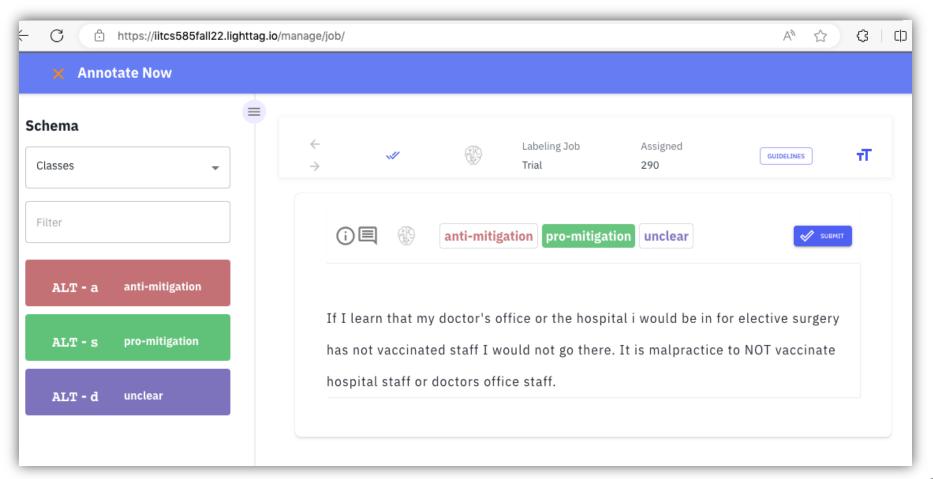
### Performance Bounds

- How good is (say) 80%?
- Evaluate performance relative to lower and upper bounds:
  - Baseline performance: how well does the simplest "reasonable" algorithm do?
    - Majority class → What if we always assign most common label? Serves as <u>lower bound</u>
  - Human performance: what percentage of the time do people agree on classification?
    - Lower agreement → Harder problem
    - Serves as <u>upper bound</u> for machine learning

# ANNOTATION & INTER-RATER RELIABILTY (AGREEMENT)

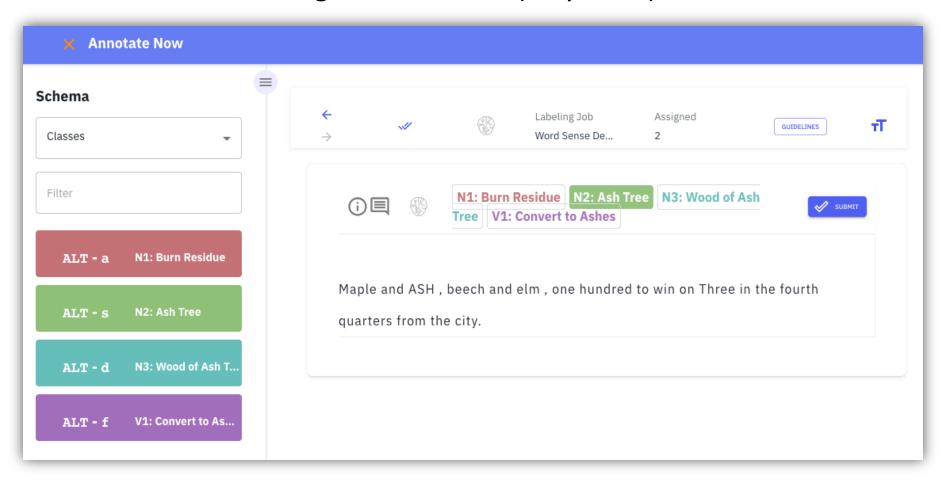
### **Annotation for NLP**

#### Annotation for Text Classification (3-class task)



### **Annotation for NLP**

Word Sense Disambiguation: "Ash" (4 Synsets)



### Measure how often humans agree on annotations

- If they don't often agree, then the task is ill-defined
- Agreement probability P(agree)
   Number of times raters agree / Number of ratings
  - But if 90% of things are annotated as X, then agreement could be high by chance
- Cohen's Kappa

$$\frac{P_{agree} - P_{chance}}{1 - P_{chance}}$$

Cohen's Kappa

$$\frac{P_{agree} - P_{chance}}{1 - P_{chance}}$$

- $P_{agree}$ : Observed agreement rate between annotators (or annotator/system)
- $P_{chance}$ : Expected agreement rate between two annotators assigning labels randomly, but using the true class distribution

- For a binary classification task with equiprobable outcomes,  $P_{chance}$  is 0.5. We'd expect raters using the two classes with equal frequency to agree half the time.
- So in this case, if  $P_{agree} = 0.7$ , then

$$\kappa = \frac{P_{agree} - P_{chance}}{1 - P_{chance}}$$

$$= \frac{0.7 - 0.5}{1 - 0.5}$$

$$= 0.4$$

- For a distribution with N classes,  $P_{chance} = \sum_{i=1}^{N} P_i^2$
- For example, for labels distributed according to (0.1,0.3,0.4,0.2):

	A (0.1)	B (0.3)	C (0.4)	D (0.2)
A (0.1)	0.01	0.03	0.04	0.02
B (0.3)	0.03	0.09	0.12	0.06
C (0.4)	0.04	0.12	0.16	0.08
D (0.2)	0.02	0.06	0.08	0.04

$$P_{chance} = 0.01 + 0.09 + 0.16 + 0.04 = 0.3$$

- For labels distributed according to (0.1,0.3,0.4,0.2),  $P_{chance} = 0.3$
- So if  $P_{agree} = 0.7$ ,

$$\kappa = \frac{P_{agree} - P_{chance}}{1 - P_{chance}}$$

$$= \frac{0.7 - 0.3}{1 - 0.3}$$

$$\approx 0.57$$

Question: Can Cohen's Kappa be negative?

### **DICTIONARY-BASED LEARNING**

### Dictionary-Based Disambiguation

- Idea: Choose between senses of a word given in a dictionary based on the words in the definitions
- ash:
  - s<sub>1</sub>: a tree of the olive family
  - s<sub>2</sub>: the solid residue left when combustible material is burned

# Algorithm (Lesk 1986)

- Define  $D_i(w)$  as the bag of words in the ith definition for w
- Define E(w) as  $\bigcup_i D_i(w)$
- For all senses s<sub>k</sub> of w, do:

$$Score(s_k) = similarity \left(D_k(w), \left[\bigcup_{v \in c} E(v)\right]\right)$$

Choose

$$s = \underset{s_k}{\operatorname{argmax}} Score(s_k)$$

# Similarity Metrics

 $similarity(X,Y) = \begin{cases} &\text{Dice coefficient } \frac{|X \cap Y|}{|X| + |Y|} \\ &\text{Jaccard coefficient } \frac{|X \cap Y|}{|X \cup Y|} \end{cases}$  Overlap coefficient  $\frac{|X \cap Y|}{\min(|X|,|Y|)}$ 

#### ash:

s<sub>1</sub>: a tree of the olive family

s<sub>2</sub>: the solid residue left when combustible material is burned

The fire had left behind nothing but a pile of ash,

The ash<sub>1</sub> can be recognized by its serrated leaves

After being struck by **lightning** the **maple** was reduced to ash,

#### ash:

s<sub>1</sub>: a tree of the olive family

s<sub>2</sub>: the solid residue left when combustible material is burned

The **fire** had left behind nothing but a pile of ash<sub>2</sub>

After be ash fire: be recognized by its serrated leaves

1. combustion or burning, in which substances uced to combine chemically with oxygen from the air 2. the shooting of projectiles from weapons

#### ash:

s<sub>1</sub>: a tree of the olive family

s<sub>2</sub>: the solid residue left when combustible material is burned

The **fire** had left behind nothing but a pile of ash<sub>2</sub>

The ash<sub>1</sub> can be recognized by its serrated leaves

### After be ileaf: truck by lightning the maple was reduced to

- ash<sub>?</sub>
- 1. a flattened structure of a higher plant or **tree**, typically green and blade-like
- 2. a thing that resembles a leaf in being flat and thin

```
ash:
```

s<sub>1</sub>: a tree of the olive family

s<sub>2</sub>: the solid residue left when combustible material is burned

#### lightning:

1. the occurrence of a natural electrical discharge of The fire habetween a cloud and the ground, often causing The as a combustion can be recognized by its serrated leaves 2. very fast

After being struck by lightning the maple was reduced to ash,

#### ash:

s<sub>1</sub>: a tree of the olive family

s<sub>2</sub>: the solid residue left when combustible material is burned

#### maple:

1. a **tree** or shrub with lobed leaves, winged fruits, The **fir** ha and colorful autumn foliage to a pile of ash<sub>2</sub>

The as 2. maple syrup or maple sugar serrated leaves

After being struck by lightning the maple was reduced to ash,

### Some Improvements

- Lesk obtained results of 50-70% accuracy
- Possible improvements:
  - Run iteratively, each time only using definitions of "appropriate" senses for context words
  - Expand each word to a set of synonyms

### **SUPERVISED LEARNING**

# Supervised Learning

- Each ambiguous word token  $w_i$  in the training is tagged with a sense from Senses $(w_i) = s_1, ..., s_k$
- ullet Each word token occurs in a context  $c_i$ 
  - (usually defined as a window around the word occurrence – up to ~100 words long)
- ullet Each context contains a set of words used as features  $v_{ii}$

# **Bayesian Classification**

- Bayes decision rule:
  - Classify  $s(w_i) = \operatorname{argmax}_s P(s \mid c_i)$
- Minimizes probability of error
- How to compute? Use Bayes' Theorem:

$$P(s_k|c) = \frac{P(c|s_k)P(s_k)}{P(c)}$$

# Bayes' Classifier (cont.)

 Note that P(c) is constant for all senses, therefore:

$$s(w_i) = \operatorname{argmax}_s P(s|c)$$
  
=  $\operatorname{argmax}_s \frac{P(c|s)}{P(c)} P(s)$   
=  $\operatorname{argmax}_s P(c|s) P(s)$ 

$$s(w_i) = \operatorname{argmax}_s (\log P(c \mid s) + \log P(s))$$

### Naïve Bayes

#### Assume:

- Features are conditionally independent, given the example class
- Feature order doesn't matter
- (bag of words model repetition counts)

$$P(c|s) = P(\{v_j : v_j \in c\}|s)$$

$$= \prod_{v_j \in c} P(v_j|s)$$
Naïve Bayes
Assumption
$$\log P(c|s) = \sum_{v_j \in c} \log P(v_j|s)$$

### Naïve Bayes Training

- For all senses  $s_k$  of w, do:
  - For all words  $v_i$  in the vocabulary, do:

$$P(v_j|s_k) = \frac{Count(v_j, s_k)}{Count(s_k)}$$

• For all senses  $s_k$  of w, do:

$$P(s_k) = \frac{Count(s_k)}{Count(w)}$$

# Naïve Bayes Classification

For all senses s<sub>k</sub> of w<sub>i</sub>, do:

$$Score(s_k) = \log P(s_k)$$

 For all words v<sub>j</sub> in the context window c<sub>i</sub>, do:

$$Score(s_k) += \log P(v_i|s_k)$$

Choose

$$s(w_i) = \underset{s_k}{\operatorname{argmax}} Score(s_k)$$

# Significant Features

Senses of "drug" (Gale et al. 1992):

'medication' prices, prescription, patent,

increase, consumer,

pharmaceutical

'illegal substance'

abuse, paraphernalia, illicit, alcohol, cocaine, traffickers

### **UNSUPERVISED LEARNING**

# Why Unsupervised Learning?

### Some issues with supervised learning:

- Domain-dependence: In computer manuals, "mouse" will not be evidence for topic "mammal"
- Coverage: "Michael Jordan" will not likely be in a wordnet, but is an excellent indicator for topic "sports"

# Tuning for a Specific Corpus

Use a naïve-Bayes formulation:

$$P(s|c) = \frac{P(s) \prod_{v \in c} P(v|s)}{\prod_{v \in c} P(v)}$$

- Initialize probabilities as uniform
- Re-estimate P(s) and  $P(v_j \mid s)$  for each sense s and each word  $v_j$  by evaluating all contexts in the corpus, assuming the context has sense s if  $P(s \mid c) > \theta$  (where  $\theta$  is a predefined threshold)
- Disambiguate by choosing the highest

### LEVERAGING BILINGUAL DATA

# Using a Bilingual Corpus

Use correlations between phrases in two languages to disambiguate

E.g, interest = 'legal share' (acquire an interest)

'attention' (show interest)

In German Beteiligung erwerben

Interesse zeigen

Depending on where the translations of related words occur, determine which sense applies

# Scoring

- Given a context c in which a syntactic relation R(w, v) holds between w and a context word v:
  - Score of sense  $s_k$  is the number of contexts c' in the second language such that  $R(w', v') \in c'$  where w' is a translation of  $s_k$  and v' is a translation of v.
  - Choose highest-scoring sense

# Using a Bilingual Corpus

#### Challenges:

- In related languages, senses may share a translation
- No occurrences found for some senses

#### Translation List for channel:

- 1. canal, canal de transmisión : a path over which signals can pass
- canal, conducto: a passage for water (or other fluids)
- [No Spanish sense]: groove
- canal, estrecho: a relatively narrow body of water linking two larger bodies; "the ship went aground in the channel"
- 5. [No Spanish sense] : line, communication channel
- canal, vía: bodily passage or tube conveying a secretion or other substance
- [No Spanish sense]: television channel

# RECENT RESEARCH IN WORD SENSE DISAMBIGUATION

### Words-in-Context Task

#### Words-in-Context dataset and task

- Evaluate context-sensitive word embeddings
- In SuperGLUE shared task more difficult NLU tasks
- Evaluate pairs of sentences (different or same sense) rather than assign label to single sentence/context

Label	Target	Context-1	Context-2
F	bed	There's a lot of trash on the <u>bed</u> of the river	I keep a glass of water next to my <u>bed</u> when I sleep
F	land	The pilot managed to land the airplane safely	The enemy <u>landed</u> several of our aircrafts
F	justify	Justify the margins	The end <u>justifies</u> the means
Т	beat	We <u>beat</u> the competition	Agassi <u>beat</u> Becker in the tennis championship
Т	air	<u>Air</u> pollution	Open a window and let in some <u>air</u>
Т	window	The expanded $\underline{\text{window}}$ will give us time to catch the thieves	You have a two-hour $\underline{\text{window}}$ of clear weather to finish working on the lawn

### Words-in-Context Task

WiC: the Word-in-Context Dataset for Evaluating Context-Sensitive Meaning Representation (Pilehvar & Camacho-Collados, NAACL 2019)

"Around 22% of the pairs in the test set had at least one of their target words not covered by these models. For such outof-vocabulary cases, we used BERT's default tokenizer...."

What is unique about BERT's default tokenizer?

### **EWISE: Applying Word Embeddings**

### Zero-shot Word Sense Disambiguation using Sense Definition Embeddings (Kumar et al. 2019)

WordNet + BiLSTM→ Sense Embeddings

