Word Sense Disambiguation

Ling571 Deep Processing Techniques for NLP February 23, 2011

Word Sense Disambiguation

- Robust Approaches
 - Supervised Learning Approaches
 - Naïve Bayes
 - Dictionary-based Approaches
 - Bootstrapping Approaches
 - One sense per discourse/collocation
 - Similarity-based approaches
 - Thesaurus-based techniques
 - Unsupervised approaches
 - Why they work, Why they don't

Disambiguation Features

- Key: What are the features?
 - Part of speech
 - Of word and neighbors
 - Morphologically simplified form
 - Words in neighborhood
 - Question: How big a neighborhood?
 - Is there a single optimal size? Why?
 - (Possibly shallow) Syntactic analysis
 - E.g. predicate-argument relations, modification, phrases
 - Collocation vs co-occurrence features
 - Collocation: words in specific relation: p-a, 1 word +/-
 - Co-occurrence: bag of words..

WSD Evaluation

- Ideally, end-to-end evaluation with WSD component
 - Demonstrate real impact of technique in system
 - Difficult, expensive, still application specific
- Typically, intrinsic, sense-based
 - Accuracy, precision, recall
 - SENSEVAL/SEMEVAL: all words, lexical sample
- Baseline:
 - Most frequent sense, Lesk
- Topline:
 - Human inter-rater agreement: 75-80% fine; 90% coarse

- Supervised learning approach
 - Input: feature vector X label

Supervised learning approach
Input: feature vector X label

Best sense = most probable sense given f

$$\hat{s} = \underset{s \in S}{\operatorname{argmax}} P(s \mid \vec{f})$$

- Supervised learning approach
 Input: feature vector X label
- Best sense = most probable sense given f

$$\hat{s} = \underset{s \in S}{\operatorname{arg\,max}} P(s \mid \vec{f})$$
$$\hat{s} = \underset{s \in S}{\operatorname{arg\,max}} \frac{P(\vec{f} \mid s)P(s)}{P(\vec{f})}$$

Issue:

Data sparseness: full feature vector rarely seen

Issue:

Data sparseness: full feature vector rarely seen

- "Naïve" assumption:
 - Features independent given sense

$$P(\vec{f} \mid s) \approx \prod_{j=1}^{n} P(f_j \mid s)$$
$$\hat{s} = \underset{s \in S}{\operatorname{argmax}} P(s) \prod_{j=1}^{n} P(f_j \mid s)$$

Training NB Classifier $\hat{s} = \underset{s \in S}{\operatorname{argmax}} P(s) \prod_{j=1}^{n} P(f_j \mid s)$ • Estimate P(s): • Prior





• Estimate $P(f_i|s)$



Training NB Classifier

$$\hat{s} = \underset{s \in S}{\operatorname{argmax}} P(s) \prod_{j=1}^{n} P(f_j \mid s)$$

• Estimate P(s):
• Prior
 $P(s_i) = \frac{count(s_i, w_j)}{count(w_j)}$
• Estimate P(f_j | s) $P(f_j \mid s) = \frac{count(f_j, s)}{count(s)}$

Issues:

Underflow => log prob



- Estimate P(f_j|s) $P(f_j|s) = \frac{count(f_j,s)}{count(s)}$
- Issues:
 - Underflow => log prob
 - Sparseness => smoothing

- (Simplified) Lesk algorithm
 - "How to tell a pine cone from an ice cream cone"

- (Simplified) Lesk algorithm
 - "How to tell a pine cone from an ice cream cone"
 - Compute context 'signature' of word to disambiguate
 - Words in surrounding sentence(s)

- (Simplified) Lesk algorithm
 - "How to tell a pine cone from an ice cream cone"
 - Compute context 'signature' of word to disambiguate
 - Words in surrounding sentence(s)
 - Compare overlap w.r.t. dictionary entries for senses

- (Simplified) Lesk algorithm
 - "How to tell a pine cone from an ice cream cone"
 - Compute context 'signature' of word to disambiguate
 - Words in surrounding sentence(s)
 - Compare overlap w.r.t. dictionary entries for senses
 - Select sense with highest (non-stopword) overlap

Applying Lesk

• The bank can guarantee deposits will eventually cover future tuition costs because it invests in mortgage securities.

| bank ¹ | Gloss: | a financial institution that accepts deposits and channels the |
|-------------------|-----------|--|
| | | money into lending activities |
| | Examples: | "he cashed a check at the bank", "that bank holds the mortgage |
| | | on my home" |
| bank ² | Gloss: | sloping land (especially the slope beside a body of water) |
| | Examples: | "they pulled the canoe up on the bank", "he sat on the bank of |
| | | the river and watched the currents" |

Applying Lesk

• The bank can guarantee deposits will eventually cover future tuition costs because it invests in mortgage securities.

| bank ¹ | Gloss: | a financial institution that accepts deposits and channels the |
|-------------------|-----------|--|
| | | money into lending activities |
| | Examples: | "he cashed a check at the bank", "that bank holds the mortgage |
| | | on my home" |
| bank ² | Gloss: | sloping land (especially the slope beside a body of water) |
| | Examples: | "they pulled the canoe up on the bank", "he sat on the bank of |
| | | the river and watched the currents" |



Applying Lesk

• The bank can guarantee deposits will eventually cover future tuition costs because it invests in mortgage securities.

| bank ¹ | Gloss: | a financial institution that accepts deposits and channels the |
|-------------------|-----------|--|
| | | money into lending activities |
| | Examples: | "he cashed a check at the bank", "that bank holds the mortgage |
| | | on my home" |
| bank ² | Gloss: | sloping land (especially the slope beside a body of water) |
| | Examples: | "they pulled the canoe up on the bank", "he sat on the bank of |
| | | the river and watched the currents" |

- Bank¹: 2
- Bank²: 0

- Overlap score:
 - All words equally weighted (excluding stopwords)

- Overlap score:
 - All words equally weighted (excluding stopwords)
- Not all words equally informative

- Overlap score:
 - All words equally weighted (excluding stopwords)
- Not all words equally informative
 - Overlap with unusual/specific words better
 - Overlap with common/non-specific words less good

- Overlap score:
 - All words equally weighted (excluding stopwords)
- Not all words equally informative
 - Overlap with unusual/specific words better
 - Overlap with common/non-specific words less good
- Employ corpus weighting:
 - IDF: inverse document frequency
 - $Idf_i = log (Ndoc/nd_i)$

Minimally Supervised WSD

- Yarowsky's algorithm (1995)
- Bootstrapping approach:
 - Use small labeled seedset to iteratively train

Minimally Supervised WSD

- Yarowsky's algorithm (1995)
- Bootstrapping approach:
 - Use small labeled seedset to iteratively train
- Builds on 2 key insights:
 - One Sense Per Discourse
 - Word appearing multiple times in text has same sense
 - Corpus of 37232 bass instances: always single sense

Minimally Supervised WSD

- Yarowsky's algorithm (1995)
- Bootstrapping approach:
 - Use small labeled seedset to iteratively train
- Builds on 2 key insights:
 - One Sense Per Discourse
 - Word appearing multiple times in text has same sense
 - Corpus of 37232 bass instances: always single sense
 - One Sense Per Collocation
 - Local phrases select single sense
 - Fish -> Bass¹
 - Play -> Bass²

- Training Decision Lists
 - 1. Pick Seed Instances & Tag

- Training Decision Lists
 - 1. Pick Seed Instances & Tag
 - 2. Find Collocations: Word Left, Word Right, Word <u>+</u>K

- Training Decision Lists
 - 1. Pick Seed Instances & Tag
 - 2. Find Collocations: Word Left, Word Right, Word +K
 - (A) Calculate Informativeness on Tagged Set,
 Order: abs(log P(Sense₁ | Collocation) P(Sense₂ | Collocation)

- Training Decision Lists
 - 1. Pick Seed Instances & Tag
 - 2. Find Collocations: Word Left, Word Right, Word +K
 - (A) Calculate Informativeness on Tagged Set,
 Order: abs(log P(Sense, |Collocation))
 - (B) Tag New Instances with Rules

- Training Decision Lists
 - 1. Pick Seed Instances & Tag
 - 2. Find Collocations: Word Left, Word Right, Word +K
 - (A) Calculate Informativeness on Tagged Set,
 Order: abs(log P(Sense, |Collocation))
 - (B) Tag New Instances with Rules
 - (C) Apply 1 Sense/Discourse

• (D)

- Training Decision Lists
 - 1. Pick Seed Instances & Tag
 - 2. Find Collocations: Word Left, Word Right, Word +K
 - (A) Calculate Informativeness on Tagged Set,
 Order: abs(log P(Sense, |Collocation))
 - (B) Tag New Instances with Rules
 - (C) Apply 1 Sense/Discourse
 - (D) If Still Unlabeled, Go To 2

- Training Decision Lists
 - 1. Pick Seed Instances & Tag
 - 2. Find Collocations: Word Left, Word Right, Word +K
 - (A) Calculate Informativeness on Tagged Set,
 Order: abs(log P(Sense₁ | Collocation) P(Sense₂ | Collocation)
 - (B) Tag New Instances with Rules
 - (C) Apply 1 Sense/Discourse
 - (D) If Still Unlabeled, Go To 2
 - 3. Apply 1 Sense/Discourse

- Training Decision Lists
 - 1. Pick Seed Instances & Tag
 - 2. Find Collocations: Word Left, Word Right, Word <u>+</u>K
 - (A) Calculate Informativeness on Tagged Set,
 - Order: $abs(log \frac{P(Sense_1 | Collocation)}{P(Sense_2 | Collocation)})$
 - (B) Tag New Instances with Rules
 - (C) Apply 1 Sense/Discourse
 - (D) If Still Unlabeled, Go To 2
 - 3. Apply 1 Sense/Discourse

Disambiguation: First Rule Matched

| Į | Initia | decision list for plant (abbrevia | ated) |
|---|--------|-------------------------------------|-----------------|
| [| LogL | Collocation | Sense |
| l | 8.10 | plant life | $\Rightarrow A$ |
| ۱ | 7.58 | manufacturing plant | ⇒ B |
| l | 7.39 | life (within $\pm 2-10$ words) | ⇒ A |
| l | 7.20 | manufacturing (in ±2-10 words) | ⇒ B |
| l | 6.27 | animal (within $\pm 2-10$ words) | ⇒ A |
| l | 4.70 | equipment (within $\pm 2-10$ words) | ⇒ B |
| Į | 4.39 | employee (within $\pm 2-10$ words) | ⇒ B |
| l | 4.30 | assembly plant | ⇒ B |
| | 4.10 | plant closure | ⇒ B |
| ۱ | 3.52 | plant species | ⇒ A |
| Ì | 3.48 | automate (within $\pm 2-10$ words) | ⇒ B |
| | 3.45 | microscopic plant | \Rightarrow A |
| í | | | |

Iterative Updating



There are more kinds of plants and animals in the rainforests than anywhere else on Earth. Over half of the millions of known species of plants and animals live in the rainforest. Many are found nowhere else. There are even plants and animals in the rainforest that we have not yet discovered. **Biological Example**

The Paulus company was founded in 1938. Since those days the product range has been the subject of constant expansions and is brought up continuously to correspond with the state of the art. We' re engineering, manufacturing and commissioning worldwide ready-to-run plants packed with our comprehensive knowhow. Our Product Range includes pneumatic conveying systems for carbon, carbide, sand, lime andmany others. We use reagent injection in molten metal for the... Industrial Example

Label the First Use of "Plant"

<u>Sense Choice With</u> Collocational Decision Lists

- Create Initial Decision List
 Rules Ordered by abs(log P(Sense, |Collocation))
- Check nearby Word Groups (Collocations)
 - Biology: "Animal" in <u>+</u> 2-10 words
 - Industry: "Manufacturing" in <u>+</u> 2-10 words
- Result: Correct Selection
 - 95% on Pair-wise tasks

• Synonymy:

• Synonymy:

• True propositional substitutability is rare, slippery

- Synonymy:
 - True propositional substitutability is rare, slippery
- Word similarity (semantic distance):
 - Looser notion, more flexible

- Synonymy:
 - True propositional substitutability is rare, slippery
- Word similarity (semantic distance):
 - Looser notion, more flexible
 - Appropriate to applications:
 - IR, summarization, MT, essay scoring

- Synonymy:
 - True propositional substitutability is rare, slippery
- Word similarity (semantic distance):
 - Looser notion, more flexible
 - Appropriate to applications:
 - IR, summarization, MT, essay scoring
 - Don't need binary +/- synonym decision

- Synonymy:
 - True propositional substitutability is rare, slippery
- Word similarity (semantic distance):
 - Looser notion, more flexible
 - Appropriate to applications:
 - IR, summarization, MT, essay scoring
 - Don't need binary +/- synonym decision
 - Want terms/documents that have high similarity

- Synonymy:
 - True propositional substitutability is rare, slippery
- Word similarity (semantic distance):
 - Looser notion, more flexible
 - Appropriate to applications:
 - IR, summarization, MT, essay scoring
 - Don't need binary +/- synonym decision
 - Want terms/documents that have high similarity
 - Differ from relatedness

- Synonymy:
 - True propositional substitutability is rare, slippery
- Word similarity (semantic distance):
 - Looser notion, more flexible
 - Appropriate to applications:
 - IR, summarization, MT, essay scoring
 - Don't need binary +/- synonym decision
 - Want terms/documents that have high similarity
 - Differ from relatedness
- Approaches:
 - Thesaurus-based
 - Distributional

• Key idea:

• Shorter path length in thesaurus, smaller semantic dist.

- Key idea:
 - Shorter path length in thesaurus, smaller semantic dist.
 - Words similar to parents, siblings in tree
 - Further away, less similar

- Key idea:
 - Shorter path length in thesaurus, smaller semantic dist.
 - Words similar to parents, siblings in tree
 - Further away, less similar
- Pathlength=# edges in shortest route in graph b/t nodes

- Key idea:
 - Shorter path length in thesaurus, smaller semantic dist.
 - Words similar to parents, siblings in tree
 - Further away, less similar
- Pathlength=# edges in shortest route in graph b/t nodes
 - Sim_{path}= -log pathlen(c₁,c₂) [Leacock & Chodorow]

- Key idea:
 - Shorter path length in thesaurus, smaller semantic dist.
 - Words similar to parents, siblings in tree
 - Further away, less similar
- Pathlength=# edges in shortest route in graph b/t nodes
 - Sim_{path}= -log pathlen(c₁,c₂) [Leacock & Chodorow]
- Problem 1:

- Key idea:
 - Shorter path length in thesaurus, smaller semantic dist.
 - Words similar to parents, siblings in tree
 - Further away, less similar
- Pathlength=# edges in shortest route in graph b/t nodes
 - Sim_{path}= -log pathlen(c₁,c₂) [Leacock & Chodorow]
- Problem 1:
 - Rarely know which sense, and thus which node

- Key idea:
 - Shorter path length in thesaurus, smaller semantic dist.
 - Words similar to parents, siblings in tree
 - Further away, less similar
- Pathlength=# edges in shortest route in graph b/t nodes
 - Sim_{path}= -log pathlen(c₁,c₂) [Leacock & Chodorow]
- Problem 1:
 - Rarely know which sense, and thus which node
- Solution: assume most similar senses estimate

- Key idea:
 - Shorter path length in thesaurus, smaller semantic dist.
 - Words similar to parents, siblings in tree
 - Further away, less similar
- Pathlength=# edges in shortest route in graph b/t nodes
 - Sim_{path}= -log pathlen(c₁,c₂) [Leacock & Chodorow]
- Problem 1:
 - Rarely know which sense, and thus which node
- Solution: assume most similar senses estimate
 - Wordsim(w_1, w_2) = max sim(c_1, c_2)

Path Length

• Path length problem:



Path Length

- Path length problem:
 - Links in WordNet not uniform
 - Distance 5: Nickel->Money and Nickel->Standard



• Solution 1:

• Solution 1:

• Build position-specific similarity measure

- Solution 1:
 - Build position-specific similarity measure
 - Not general
- Solution 2:

- Solution 1:
 - Build position-specific similarity measure
 - Not general
- Solution 2:
 - Add corpus information: information-content measure
 - P(c) : probability that a word is instance of concept c

- Solution 1:
 - Build position-specific similarity measure
 - Not general
- Solution 2:
 - Add corpus information: information-content measure
 - P(c) : probability that a word is instance of concept c
 - Words(c) : words subsumed by concept c; N: words in corpus

$$P(c) = \frac{\sum_{w \in words(c)} count(w)}{N}$$

- Information content of node:
 - IC(c) = -log P(c)

- Information content of node:
 - IC(c) = -log P(c)
- Least common subsumer (LCS):
 - Lowest node in hierarchy subsuming 2 nodes

- Information content of node:
 - IC(c) = -log P(c)
- Least common subsumer (LCS):
 - Lowest node in hierarchy subsuming 2 nodes
- Similarity measure:
 - $sim_{RESNIK}(c_1,c_2) = \log P(LCS(c_1,c_2))$

- Information content of node:
 - IC(c) = -log P(c)
- Least common subsumer (LCS):
 - Lowest node in hierarchy subsuming 2 nodes
- Similarity measure:
 - $sim_{RESNIK}(c_1, c_2) = \log P(LCS(c_1, c_2))$
- Issue:

- Information content of node:
 - IC(c) = -log P(c)
- Least common subsumer (LCS):
 - Lowest node in hierarchy subsuming 2 nodes
- Similarity measure:
 - $sim_{RESNIK}(c_1,c_2) = \log P(LCS(c_1,c_2))$
- Issue:
 - Not content, but difference between node & LCS

- Information content of node:
 - IC(c) = -log P(c)
- Least common subsumer (LCS):
 - Lowest node in hierarchy subsuming 2 nodes
- Similarity measure:
 - $sim_{RESNIK}(c_1,c_2) = \log P(LCS(c_1,c_2))$
- Issue:
 - Not content, but difference between node & LCS $sim_{Lin}(c_1, c_2) = \frac{2 \times \log P(LCS(c_1, c_2))}{\log P(c_1) + \log P(c_2)}$