

Linear models for classification:

$$\hat{y} = \arg \max_y \theta^T \mathbf{f}(\mathbf{x}, y) \quad (1)$$

Linear models for classification:

$$\hat{y} = \arg \max_y \boldsymbol{\theta}^T \mathbf{f}(\mathbf{x}, y) \quad (1)$$

- Feature function representation

Linear models for classification:

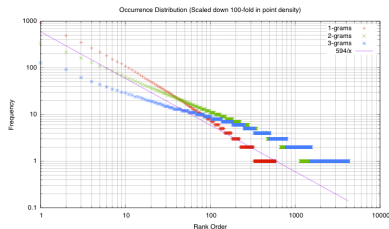
$$\hat{y} = \arg \max_y \boldsymbol{\theta}^T \mathbf{f}(\mathbf{x}, y) \quad (1)$$

- Feature function representation
- Weights

# A question for you

Remember Zipf's law?

$$freq \propto \frac{1}{rank}$$

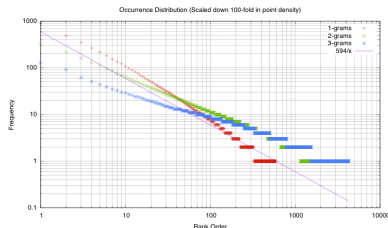


# A question for you

Remember Zipf's law?

$$freq \propto \frac{1}{rank}$$

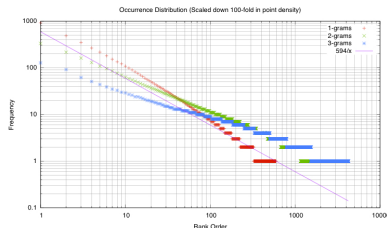
- Suppose you have a corpus with  $N$  word tokens
- $K$  of these tokens appear exactly once.  
(These are called *hapax legomena*.)



# A question for you

Remember Zipf's law?

$$\text{freq} \propto \frac{1}{\text{rank}}$$



- Suppose you have a corpus with  $N$  word tokens
- $K$  of these tokens appear exactly once.  
(These are called *hapax legomena*.)
- Now suppose you get  $2N$  tokens from the same corpus.  
How many words appear exactly once in the new corpus?
  - 1 roughly  $2K$
  - 2 more than  $K$ , but less than  $2K$
  - 3 roughly  $K$
  - 4 less than  $K$  but more than  $\frac{K}{2}$
  - 5 roughly  $\frac{K}{2}$

# IPython Notebook

- <http://ipython.org/notebook.html>
- Browser-based IDE for Python
- Integrates code, text, LaTeX, ...

Linear models for classification:

$$\hat{y} = \arg \max_y \theta^T \mathbf{f}(\mathbf{x}, y) \quad (1)$$



Linear models for classification:

$$\hat{y} = \arg \max_y \boldsymbol{\theta}^T \mathbf{f}(\mathbf{x}, y) \quad (1)$$

- Feature function representation

Linear models for classification:

$$\hat{y} = \arg \max_y \boldsymbol{\theta}^T \mathbf{f}(\mathbf{x}, y) \quad (1)$$

- Feature function representation
- Weights

# Feature functions

Suppose  $y \in \mathcal{Y} = \{\text{pos}, \text{neg}\}$ . Then,

$$\mathbf{f}(\mathbf{x}, y = \text{pos}) = [\mathbf{x}^T, \mathbf{0}^T]^T$$

$$\mathbf{f}(\mathbf{x}, y = \text{neg}) = [\mathbf{0}^T, \mathbf{x}^T]^T$$

# Feature functions

Suppose  $y \in \mathcal{Y} = \{\text{pos}, \text{neg}, \text{neut}\}$ . Then,

$$\mathbf{f}(\mathbf{x}, y = \text{pos}) = [\mathbf{x}^T, \mathbf{0}^T, \mathbf{0}^T]^T$$

$$\mathbf{f}(\mathbf{x}, y = \text{neg}) = [\mathbf{0}^T, \mathbf{x}^T, \mathbf{0}^T]^T$$

$$\mathbf{f}(\mathbf{x}, y = \text{neut}) = [\mathbf{0}^T, \mathbf{0}^T, \mathbf{x}^T]^T$$

# Feature functions

Suppose  $y \in \mathcal{Y} = \{\text{pos}, \text{neg}, \text{neut}\}$ . Then,

$$\mathbf{f}(\mathbf{x}, y = \text{pos}) = [\mathbf{x}^T, \mathbf{0}^T, \mathbf{0}^T]^T$$

$$\mathbf{f}(\mathbf{x}, y = \text{neg}) = [\mathbf{0}^T, \mathbf{x}^T, \mathbf{0}^T]^T$$

$$\mathbf{f}(\mathbf{x}, y = \text{neut}) = [\mathbf{0}^T, \mathbf{0}^T, \mathbf{x}^T]^T$$

The feature vector is composed of individual feature functions, e.g.:

$$\begin{aligned} f_{176}(\mathbf{x}, y) &:= x_{176} \times \delta(y = \text{pos}) \\ &= \delta(\text{great} \in \mathbf{w} \wedge y = \text{pos}) \end{aligned}$$

$$f_{177}(\mathbf{x}, y) := x_{177} \times \delta(y = \text{pos})$$

$$f_{10176}(\mathbf{x}, y) := x_{176} \times \delta(y = \text{neg}) \dots$$

# Feature functions

Suppose  $y \in \mathcal{Y} = \{\text{pos}, \text{neg}, \text{neut}\}$ . Then,

$$\mathbf{f}(\mathbf{x}, y = \text{pos}) = [\mathbf{x}^T, 1, \mathbf{0}^T, \mathbf{0}^T]^T$$

$$\mathbf{f}(\mathbf{x}, y = \text{neg}) = [\mathbf{0}^T, \mathbf{x}^T, 1, \mathbf{0}^T]^T$$

$$\mathbf{f}(\mathbf{x}, y = \text{neut}) = [\mathbf{0}^T, \mathbf{0}^T, \mathbf{x}^T, 1]^T$$

The feature vector is composed of individual feature functions, e.g.:

$$\begin{aligned} f_{176}(\mathbf{x}, y) &:= x_{176} \times \delta(y = \text{pos}) \\ &= \delta(\text{great} \in \mathbf{w} \wedge y = \text{pos}) \end{aligned}$$

$$f_{177}(\mathbf{x}, y) := x_{177} \times \delta(y = \text{pos})$$

$$f_{10176}(\mathbf{x}, y) := x_{176} \times \delta(y = \text{neg}) \dots$$

We usually add an “offset” feature at the end of each vector.

# Weights: Naive Bayes

$$\begin{aligned}\theta^T \mathbf{f}(\mathbf{x}, y) &:= \log P(\mathbf{x}, y; \phi, \mu) \\ &= \log P(\mathbf{x}|y; \phi) P(y; \mu) \\ &= \log P(\mathbf{x}|y; \phi) + \log P(y; \mu)\end{aligned}$$

# Weights: Naive Bayes

$$\begin{aligned}\theta^T \mathbf{f}(\mathbf{x}, y) &:= \log P(\mathbf{x}, y; \phi, \mu) \\ &= \log P(\mathbf{x}|y; \phi) P(y; \mu) \\ &= \log P(\mathbf{x}|y; \phi) + \log P(y; \mu) \\ &= \log \text{Multinomial}(\mathbf{x}; \phi_y) + \log \text{Cat}(y; \mu)\end{aligned}$$



# Weights: Naive Bayes

$$\begin{aligned}\theta^T \mathbf{f}(\mathbf{x}, y) &:= \log P(\mathbf{x}, y; \phi, \mu) \\ &= \log P(\mathbf{x}|y; \phi) P(y; \mu) \\ &= \log P(\mathbf{x}|y; \phi) + \log P(y; \mu) \\ &= \log \text{Multinomial}(\mathbf{x}; \phi_y) + \log \text{Cat}(y; \mu) \\ &= \log \frac{(\sum_n x_n)!}{\prod_n x_n!} + \log \prod_n \phi_{y,n}^{x_n} + \log \mu_y\end{aligned}$$

# Weights: Naive Bayes

$$\begin{aligned}\theta^T \mathbf{f}(\mathbf{x}, y) &:= \log P(\mathbf{x}, y; \phi, \mu) \\ &= \log P(\mathbf{x}|y; \phi) P(y; \mu) \\ &= \log P(\mathbf{x}|y; \phi) + \log P(y; \mu) \\ &= \log \text{Multinomial}(\mathbf{x}; \phi_y) + \log \text{Cat}(y; \mu) \\ &= \log \frac{(\sum_n x_n)!}{\prod_n x_n!} + \log \prod_n \phi_{y,n}^{x_n} + \log \mu_y \\ &\propto \sum_n x_n \log \phi_{y,n} + \log \mu_y\end{aligned}$$

# Weights: Naive Bayes

$$\begin{aligned}\theta^T \mathbf{f}(\mathbf{x}, y) &:= \log P(\mathbf{x}, y; \phi, \mu) \\ &= \log P(\mathbf{x}|y; \phi) P(y; \mu) \\ &= \log P(\mathbf{x}|y; \phi) + \log P(y; \mu) \\ &= \log \text{Multinomial}(\mathbf{x}; \phi_y) + \log \text{Cat}(y; \mu) \\ &= \log \frac{(\sum_n x_n)!}{\prod_n x_n!} + \log \prod_n \phi_{y,n}^{x_n} + \log \mu_y \\ &\propto \sum_n x_n \log \phi_{y,n} + \log \mu_y \\ &= \theta^T \mathbf{f}(\mathbf{x}, y)\end{aligned}$$

where

$$\begin{aligned}\theta &= [\log \phi_1^T, \log \mu_1, \log \phi_2^T, \log \mu_2, \dots]^T \\ \mathbf{f}(\mathbf{x}, y) &= [\mathbf{0}, \dots, \mathbf{0}, \mathbf{x}^T, 1, \mathbf{0}, \dots, \mathbf{0}]^T\end{aligned}$$

# Today

## Naive Bayes

- Recap maximum likelihood estimation
- Smoothing, and bias-variance tradeoff
- Practical details of machine learning
- Features, and the naivety of Naive Bayes

# Today

## Naive Bayes

- Recap maximum likelihood estimation
- Smoothing, and bias-variance tradeoff
- Practical details of machine learning
- Features, and the naivety of Naive Bayes

## Perceptron

- Error-driven classification
- Averaged perceptron
- Mira (maybe)

# Today

## Naive Bayes

- Recap maximum likelihood estimation
- Smoothing, and bias-variance tradeoff
- Practical details of machine learning
- Features, and the naivety of Naive Bayes

## Perceptron

- Error-driven classification
- Averaged perceptron
- Mira (maybe)

## Word sense disambiguation

- Definition of word senses
- Formulation as a classification problem

Remember these headlines?

- Iraqi head seeks arms
- Prostitutes appeal to Pope
- Drunk gets nine years in violin case

Remember these headlines?

- Iraqi head seeks arms
- Prostitutes appeal to Pope
- Drunk gets nine years in violin case

They are ambiguous because words have multiple senses.

- head: BODY-PART, LEADER
- arms: BODY-PART, WEAPON



Remember these headlines?

- Iraqi head seeks arms
- Prostitutes appeal to Pope
- Drunk gets nine years in violin case

They are ambiguous because words have multiple senses.

- head: BODY-PART, LEADER
- arms: BODY-PART, WEAPON

Can you see what is ambiguous about the other examples?

# Word sense disambiguation

Word Sense Disambiguation (WSD) is the problem of identifying the intended sense of each word token.

- Part of a larger field of research called **lexical semantics**
- Part-of-speech ambiguity (**i'm heading out of town**) is usually considered to be a different problem.
- For WSD, words include their POS tag (e.g., **heading/V**)
- Technically, we want to differentiate senses of each *lemma*.  
A *lemma* is a linguistic term for a group of inflected forms: **arm**, **arms**; **serve**, **served**, **serves**, **serving**.

# How many word senses?

Words (lemmas) may have *many* more than two senses.

For example, **serve**:

- [FUNCTION]: The tree stump served as a table
- [ENABLE]: His evasive replies only served to heighten suspicion
- [DISH]: We serve only the rawest fish here
- [ENLIST]: She served her country in the marines
- [JAIL]: He served six years in Alcatraz
- [TENNIS]: Nobody can return his double-reverse spin serve
- [LEGAL]: They were served with subpoenas
- more?

# How many word senses?

How can we test that these senses are really different?

We can construct a **zeugma**, which combines antagonistic senses in an uncomfortable way:

- Which flight serves breakfast?
- Which flights serve Tuscon?

# How many word senses?

How can we test that these senses are really different?


We can construct a **zeugma**, which combines antagonistic senses in an uncomfortable way:

- Which flight serves breakfast?
- Which flights serve Tuscon?
- \*Which flights serve breakfast and Tuscon?


The asterisk is a linguistic notation for utterances which would not be judged to be grammatical by fluent speakers of a language.

# The WSD task: Output

- What should the output of WSD be?  
What are the possible senses for each word?
- We could just look in the dictionary.

**<sup>1</sup>plunge**  *verb* \ˈplʌŋj\

**plunged** | **plung·ing**

**Definition of PLUNGE** 

*transitive verb*

- 1** : to cause to penetrate or enter quickly and forcibly into something <*plunged the dagger*>
- 2** : to cause to enter a state or course of action usually suddenly, unexpectedly, or violently <*plunged the nation into economic depression*>

*intransitive verb*

- 1** : to thrust or cast oneself into or as if into water
- 2 a** : to become pitched or thrown headlong or violently forward and downward; *also* : to move oneself in such a manner <*plunged off the embankment*>  
**b** : to act with reckless haste : enter suddenly or unexpectedly <*plunges into project after project*>  
**c** : to bet or gamble heavily and recklessly
- 3** : to descend or dip suddenly <*the stock's value plunged*>

WSD research is dominated by a computational resource called WORDNET. (<http://wordnet.princeton.edu>)

## WordNet Search - 3.1

[- WordNet home page](#) - [Glossary](#) - [Help](#)

Word to search for:

Display Options:

Key: "S:" = Show Synset (semantic) relations, "W:" = Show Word (lexical) relations  
Display options for sense: (gloss) "an example sentence"

### Noun

- **S: (n) bass** (the lowest part of the musical range)
- **S: (n) bass, bass part** (the lowest part in polyphonic music)
- **S: (n) bass, basso** (an adult male singer with the lowest voice)
- **S: (n) sea bass, bass** (the lean flesh of a saltwater fish of the family Serranidae)
  - [direct hyponym](#) / [full hyponym](#)
  - [direct hypernym](#) / [inherited hypernym](#) / [sister term](#)
    - **S: (n) saltwater fish** (flesh of fish from the sea used as food)
  - [part holonym](#)
- **S: (n) freshwater bass, bass** (any of various North American freshwater fish with lean flesh (especially of the genus *Micropterus*))
- **S: (n) bass, bass voice, basso** (the lowest adult male singing voice)
- **S: (n) bass** (the member with the lowest range of a family of musical instruments)
- **S: (n) bass** (nontechnical name for any of numerous edible marine and freshwater spiny-finned fishes)

### Adjective

- **S: (adj) bass, deep** (having or denoting a low vocal or instrumental range) "a deep voice"; "a bass voice is lower than a baritone voice"; "a bass clarinet"

- WordNet consists of roughly 100K *synsets*, groups of words or phrases with an identical meaning. (e.g., {CHUMP<sup>1</sup>, FOOL<sup>2</sup>, SUCKER<sup>1</sup>, MARK<sup>9</sup>})  
A lemma is polysemous if it participates in multiple synsets.



- WordNet consists of roughly 100K *synsets*, groups of words or phrases with an identical meaning. (e.g., {CHUMP<sup>1</sup>, FOOL<sup>2</sup>, SUCKER<sup>1</sup>, MARK<sup>9</sup>})  
A lemma is polysemous if it participates in multiple synsets.
- WordNet also describes many other lexical relationships:
  - antonymy (x means the opposite of y)
  - hyponymy (x is a hyponym of y if x is-a y)
  - ...

Some statistics of English Wordnet 3:

POS	polysemy
NOUN	1.24
VERB	2.17
ADJECTIVE	1.40
ADVERB	1.25

# WordNet Senses: Pros and cons

- WordNet played a big role in helping WSD move from toy systems to large-scale quantitative evaluations.

# WordNet Senses: Pros and cons

- WordNet played a big role in helping WSD move from toy systems to large-scale quantitative evaluations.
- WordNet's sense granularity may be too fine [IW06].  
Humans agree on 75-80% of examples using WordNet senses.

# WordNet Senses: Pros and cons

- WordNet played a big role in helping WSD move from toy systems to large-scale quantitative evaluations.
- WordNet's sense granularity may be too fine [IW06].  
Humans agree on 75-80% of examples using WordNet senses.
- Are word senses real?  
The premise that word senses can be differentiated in a task-neutral way has been criticized as linguistically naïve [Kil97].

# WordNet Senses: Pros and cons

- WordNet played a big role in helping WSD move from toy systems to to large-scale quantitative evaluations.
- WordNet's sense granularity may be too fine [IW06].  
Humans agree on 75-80% of examples using WordNet senses.
- Are word senses real?  
The premise that word senses can be differentiated in a task-neutral way has been criticized as linguistically naïve [Kil97].
- WordNets are heavyweight.
  - expensive to develop for new languages
  - become outdated as language changes  
(consider: I'm **dead** **tired**, **sick** as a positive adjective, etc)
  - Would WordNet have good coverage for Twitter?

# Translation Sets as Word Senses

- An alternative is to use translation to differentiate word senses.
- E.g., since **bill** is translated as **pico** or **cuenta** in spanish, there are clearly two senses.
- But if there is no language with different spellings of the purported senses, then they are not meaningfully different.
- Most WSD research has focused on WordNet, so we will too.

- **Synthetic** data: different words are conflated (**banana-phone**), the system must identify the original word.



- **Synthetic** data: different words are conflated (**banana-phone**), the system must identify the original word.
- **Lexical sample**: disambiguate a few target words (e.g., “plant” etc).  
First large-scale WSD evaluation, SENSEVAL-1 (1998).

- **Synthetic** data: different words are conflated (**banana-phone**), the system must identify the original word.
- **Lexical sample**: disambiguate a few target words (e.g., “plant” etc). First large-scale WSD evaluation, SENSEVAL-1 (1998).
- **All-words** WSD: a sense must be identified for every token.
  - A **semantic concordance** is a corpus in which each open-class word (nouns, verbs, adjectives, and adverbs) is tagged with its word sense from the target dictionary or thesaurus.
  - SEMCOR is a semantic concordance built from 234K tokens of the Brown corpus.

As of Sunday<sub>n</sub><sup>1</sup> night<sub>n</sub><sup>1</sup> there was<sub>v</sub><sup>4</sup> no word<sub>n</sub><sup>2</sup> ...

# How can we solve WSD?

- How can we tell living **plants** from manufacturing **plants**?

# How can we solve WSD?

- How can we tell living **plants** from manufacturing **plants**?
- **Context**

# How can we solve WSD?

- How can we tell living **plants** from manufacturing **plants**?
- **Context**
  - Town officials are hoping to attract new manufacturing plants through weakened environmental regulations.
  - The endangered plant plays an important role in the local ecosystem.

# How can we solve WSD?

- How can we tell living **plants** from manufacturing **plants**?
- **Context**
  - Town officials are hoping to attract new manufacturing plants through weakened environmental regulations.
  - The endangered plant plays an important role in the local ecosystem.
- Approaches:
  - Knowledge-based
  - Supervised
  - Semi-supervised
  - Unsupervised

# The Lesk Algorithm

- For each sentence  $s_i$  and target word  $w_{ij}$ 
  - Set  $maxOverlap \leftarrow 0$ ,  $bestSense \leftarrow \emptyset$
  - For each possible sense  $t$ 
    - Compute word overlap between  $s_i$  and definition  $w_{ij}[t]$
    - If greater than  $maxOverlap$ , then update  $maxOverlap$  and  $bestSense$ .

# The Lesk Algorithm

- For each sentence  $s_i$  and target word  $w_{ij}$ 
  - Set  $maxOverlap \leftarrow 0$ ,  $bestSense \leftarrow \emptyset$
  - For each possible sense  $t$ 
    - Compute word overlap between  $s_i$  and definition  $w_{ij}[t]$
    - If greater than  $maxOverlap$ , then update  $maxOverlap$  and  $bestSense$ .

Example text: I stopped by the **bank** to try to get a loan

Example definitions:

- Bank<sup>1</sup>: financial institution which borrows and loans money
- Bank<sup>2</sup>: body of land adjacent to a river

The first sense is preferred because the word “loan” appears in both the definition and the query sentence.



# Selectional restrictions

Some verbs have strong selectional restrictions about their arguments:

- They closed the bank<sup>1</sup> after discovering its malfeasance.
- They rested on the bank<sup>2</sup> of the Seine.
- Closed can only take an argument which is an ORGANIZATION.
- Rested can only take an argument which is a PHYSICAL-OBJECT.

Some ontologies categorize common nouns in terms of such properties.

# Supervised WSD

- With labeled data, we can treat WSD as a standard supervised learning problem.
- Some features
  - Bag-of-words
  - Positional (collocation) features
  - Patterns
  - Syntax
  - Document features

# Bag-of-words features

Bag-of-words models are a very typical approach. For example,

$$f(y, \text{bank, I went to the bank to deposit my paycheck}) = \\ \{\langle \text{went}, y \rangle : 1, \langle \text{deposit}, y \rangle : 1, \langle \text{paycheck}, y \rangle : 1\}$$

# Bag-of-words features

Bag-of-words models are a very typical approach. For example,

$$f(y, \text{bank}, \text{I went to the bank to deposit my paycheck}) = \{ \langle \text{went}, y \rangle : 1, \langle \text{deposit}, y \rangle : 1, \langle \text{paycheck}, y \rangle : 1 \}$$

Some examples:

- **bank**[FINANCIAL]:

*a an and are ATM Bonnie card charges check Clyde criminals  
deposit famous for get I much My new overdraft really robbers  
the they think to too two went were*

- **bank**[RIVER]:

*a an and big campus cant catfish East got grandfather great has  
his I in is Minnesota Mississippi muddy My of on planted pole  
pretty right River The the there University walk Wets*

# Positional (collocation) features

- An extension of bag-of-words models is to encode the position of each context word, e.g.,

$$f(y, \text{bank}, \text{I went to the bank to deposit my paycheck}) = \\ \{ \langle i - 3, \text{went}, y \rangle : 1, \langle i + 2, \text{deposit}, y \rangle : 1, \langle i + 4, \text{paycheck}, y \rangle : 1 \}$$

- J&M (optional textbook) call these collocation features; the POS tag of each word can also be included.

# Pattern features

Pattern features extend the idea of positional features with explicit, regex-like patterns:

- bank account
- bank of COUNTRY.

Such features are often used in combination with non-linear classifiers such as decision lists.

# Syntactic features

- Rather than look at local neighbors, we can give special priority to the heads of phrases.
- For example, in

*I deposited my paycheck when I got to the bank*

the most revealing features are **deposit** and **paycheck**.

- **deposit** is the head of the main verb phrase for the sentence, and **paycheck** is the direct object.
- This is a clue that they are more relevant than the words immediately surrounding **bank**.

# Document-level features

According to the “one-sense-per-discourse” heuristic, a document about financial institutions is very unlikely to use the word **bank** in the **river bank** sense.

Word	Senses	Accuracy	Applicblty
plant	living/factory	99.8 %	72.8 %
tank	vehicle/contr	99.6 %	50.5 %
poach	steal/boil	100.0 %	44.4 %
palm	tree/hand	99.8 %	38.5 %
axes	grid/tools	100.0 %	35.5 %
sake	benefit/drink	100.0 %	33.7 %
bass	fish/music	100.0 %	58.8 %
space	volume/outer	99.2 %	67.7 %
motion	legal/physical	99.9 %	49.8 %
crane	bird/machine	100.0 %	49.1 %
<b>Average</b>		<b>99.8 %</b>	<b>50.1 %</b>

(Yarowsky 1995)



# Is Word Sense Disambiguation Important?

- Early machine translation researchers were really worried about WSD.
  - bill[BIRD JAW] → pico
  - bill[INVOICE] → cuenta

# Is Word Sense Disambiguation Important?

- Early machine translation researchers were really worried about WSD.
  - bill[BIRD JAW] → pico
  - bill[INVOICE] → cuenta
- Bar-Hillel, an expert-turned-skeptic, poses this problem:

*"Little John was looking for his toy box. Finally he found it. The box was in the pen." Is pen a writing instrument or a place where children play?*

# Is Word Sense Disambiguation Important?


- Early machine translation researchers were really worried about WSD.
  - bill[BIRD JAW] → pico
  - bill[INVOICE] → cuenta
- Bar-Hillel, an expert-turned-skeptic, poses this problem:

*“Little John was looking for his toy box. Finally he found it. The box was in the pen.” Is pen a writing instrument or a place where children play?*
- The suggestion is this example requires deep knowledge and inference (a box is bigger than a pen[WRITING], but not bigger than a pen[ENCLOSURE]).
- Bar-Hillel was so discouraged that he gave up on MT!

# The Role of WSD Today

- WSD was also thought to be important for information retrieval: **bass experts, help with cures**, etc.
- Many thought the NLP pipeline required a WSD module.  
preprocessing → POS tagging → WSD → application

---


<sup>1</sup>The survey argues that WSD will become relevant as performance improves. 

# The Role of WSD Today

- WSD was also thought to be important for information retrieval: *bass experts, help with cures*, etc.
- Many thought the NLP pipeline required a WSD module.  
preprocessing → POS tagging → WSD → application
- However, years of research on WSD have produced little evidence that it helps downstream applications. A recent survey of WSD notes:

*Unfortunately, to date explicit WSD has not yet demonstrated real benefits in human language technology applications (Navigli 2009).<sup>1</sup>*

---

<sup>1</sup>The survey argues that WSD will become relevant as performance improves. 

# The Role of WSD Today

- There is some evidence that WSD helps translation [CNC07, CW07]
- But in many tasks, higher-order n-grams encode much the same information as WSD.
  - If we have the bigram **bank teller** as a feature, we don't need to disambiguate **bank**.
  - Phrase-based machine translation uses a similar idea.

# Homework 2

- Download the SemCor data.
- Compare the word sense annotations with WordNet online.
- Explain why alternative senses were not chosen.
- Do word sense annotations for one sentence of text from an (English language) blog that you like.



Yee Seng Chan, Hwee Tou Ng, and David Chiang.

Word sense disambiguation improves statistical machine translation.  
*In Proceedings of the Association for Computational Linguistics (ACL)*, Prague, 2007.



Marine Carpuat and Dekai Wu.

Improving statistical machine translation using word sense disambiguation.

*In Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 61–72, 2007.



N. Ide and Y. Wilks.

Making sense about sense, 2006.



Adam Kilgarriff.

I don't believe in word senses.  
*CoRR*, cmp-lg/9712006, 1997.