



COMPUTER SCIENCE
UNIVERSITY OF MARYLAND

Words & their Meaning: Word Sense Disambiguation

CMSC 470

Marine Carpuat

Today: Word Meaning

2 core issues from an NLP perspective

- **Semantic similarity:** given two words, how similar are they in meaning?
- **Word sense disambiguation:** given a word that has more than one meaning, which one is used in a specific context?

**“Big rig carrying fruit crashes on 210 Freeway,
creates jam”**

<http://articles.latimes.com/2013/may/20/local/la-me-ln-big-rig-crash-20130520>

How do we know that a word (lemma) has distinct senses?

- Linguists often design tests for this purpose
- e.g., **zeugma** combines distinct senses in an uncomfortable way

Which flight serves breakfast?

Which flights serve BWI?

*Which flights serve breakfast and BWI?

Word Senses

- “Word sense” = distinct meaning of a word
- Same word, different senses
 - **Homonyms** (homonymy): unrelated senses; identical orthographic form is coincidental
 - E.g., financial bank vs. river bank
 - **Polysemes** (polysemy): related, but distinct senses
 - E.g., Financial bank vs. blood bank vs. tree bank
 - **Metonyms** (metonymy): “stand in”, technically, a sub-case of polysemy
 - E.g., use “Washington” in place of “the US government”
- Different word, same sense
 - **Synonyms** (synonymy)

WordNet: a lexical database for English

<https://wordnet.princeton.edu/>

- Includes most English nouns, verbs, adjectives, adverbs
- Electronic format makes it amenable to automatic manipulation: used in many NLP applications
- “WordNets” generically refers to similar resources in other languages

Synonymy in WordNet

- WordNet is organized in terms of “synsets”
 - Unordered set of (roughly) synonymous “words” (or multi-word phrases)
- Each synset expresses a distinct meaning/concept

WordNet: Example

Noun

{pipe, tobacco pipe} (a tube with a small bowl at one end; used for smoking tobacco)

{pipe, pipe, piping} (a long tube made of metal or plastic that is used to carry water or oil or gas etc.)

{pipe, tube} (a hollow cylindrical shape)

{pipe} (a tubular wind instrument)

{organ pipe, pipe, pipework} (the flues and stops on a pipe organ)

Verb

{shriek, shrill, pipe up, pipe} (utter a shrill cry)

{pipe} (transport by pipeline) “pipe oil, water, and gas into the desert”

{pipe} (play on a pipe) “pipe a tune”

{pipe} (trim with piping) “pipe the skirt”

WordNet 3.0: Size

Part of speech	Word form	Synsets
Noun	117,798	82,115
Verb	11,529	13,767
Adjective	21,479	18,156
Adverb	4,481	3,621
Total	155,287	117,659

Different inventories can be used to define senses

WordNet Sense	Spanish Translation	Roget Category	Target Word in Context
bass ⁴	lubina	FISH/INSECT	... fish as Pacific salmon and striped bass and...
bass ⁴	lubina	FISH/INSECT	... produce filets of smoked bass or sturgeon...
bass ⁷	bajo	MUSIC	... exciting jazz bass player since Ray Brown...
bass ⁷	bajo	MUSIC	... play bass because he doesn't have to solo...

Different inventories do not always agree on sense distinctions
e.g., translation makes some distinctions but not others

Exercise: how many senses of “drive”?

1. *"Can you drive this four-wheel truck?"*
2. *"We drive to the university every morning"*
3. *"We drive the car to the garage"*
4. *"He drives me mad"*
5. *"She is driven by her passion"*
6. *"Drive a nail into the wall"*
7. *" She is driving away at her doctoral thesis"*
8. *"What are you driving at?"*
9. *"My new truck drives well"*
10. *"She drives for the taxi company in Newark"*
11. *"drive the cows into the barn"*
12. *"We drive the turnpike to work"*
13. *"drive a golf ball"*

Exercise: how many senses of “drive”?

1. *"Can you drive this four-wheel truck?"*
2. *"We drive to the university every morning"*
3. *"We drive the car to the garage"*
4. *"He drives me mad"*
5. *"She is driven by her passion"*
6. *"Drive a nail into the wall"*
7. *" She is driving away at her doctoral thesis"*
8. *"What are you driving at?"*
9. *"My new truck drives well"*
10. *"She drives for the taxi company in Newark"*
11. *"drive the cows into the barn"*
12. *"We drive the turnpike to work"*
13. *"drive a golf ball"*

13 distinct senses
according to WordNet!

Exercise: how many senses of “drive”?

1. *"We drive to the university every morning"* (operate or control a vehicle)
2. *"We drive the car to the garage"* (cause someone or something to move by driving)
3. *"He drives me mad"* (force into or from an action or state, either physically or metaphorically)
4. *"She is driven by her passion"* (to compel or force or urge relentlessly or exert coercive pressure on, or motivate strongly)
5. *"Drive a nail into the wall"* (push, propel, or press with force)
6. *"She is driving away at her doctoral thesis"* (strive and make an effort to reach a goal)
7. *"What are you driving at?"* (move into a desired direction of discourse)
8. *"My new truck drives well"* (have certain properties when driven)
9. *"She drives for the taxi company in Newark"* (work as a driver)
10. *"drive the cows into the barn"* (urge forward)
11. *"We drive the turnpike to work"* (proceed along in a vehicle)
12. *"drive a golf ball"* (strike with a driver, as in teeing off)

13 distinct senses
according to WordNet!

What can we do when humans who annotate senses disagree?

- Disagreement is inevitable when annotating based on human judgments
 - Even with trained annotators
 - There is no “ground truth”
- We cannot measure “correctness” of annotations directly
- Instead, we can measure reliability of annotation
 - Do human annotators make same decisions consistently?
 - Assumption: high reliability implies validity

Quantifying (dis)agreement between human annotators: Cohen's Kappa

- Measures agreement between two annotators while taking into account the possibility of chance agreement

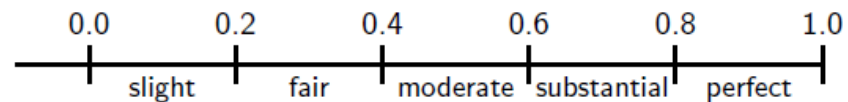
Probability of actual agreement

Probability of expected agreement

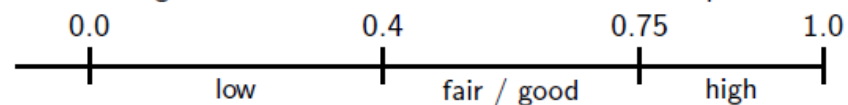
$$K = \frac{\text{Pr}(a) - \text{Pr}(e)}{1 - \text{Pr}(e)}$$

- Scales for interpreting Kappa

Landis & Koch, 1977



Green, 1997



Quantifying (dis)agreement between human annotators: Cohen's Kappa

Consider this confusion matrix for sense annotations by A and B of the same 250 examples

	Sense 1	Sense 2	Sense 3	Total
Sense 1	54	28	3	85
Sense 2	31	18	23	72
Sense 3	0	21	72	93
Total	85	67	98	250

$$K = \frac{\text{Pr}(a) - \text{Pr}(e)}{1 - \text{Pr}(e)}$$

Here $\text{Pr}(a) = 0.576$, $\text{Pr}(e) = 0.339$, $K=0.36$
(agreement is low)

Word Sense Disambiguation

what you should know (so far)

- Word senses distinguish different meanings of same word
- Sense inventories provide definitions of word senses
- Sense distinctions and annotations are based on human judgment
 - no “ground truth”
 - Measure annotation reliability using inter-annotator agreement

Word Sense Disambiguation

- Computational task
 - Given a predefined sense inventory (e.g., WordNet)
 - Goal: automatically select the correct sense of a word
 - Input: a word in context
 - Output: sense of the word
- Motivated by many applications:
 - Information retrieval
 - Machine translation
 - ...

How hard is the problem?

- **Most words in English have only one sense**
 - 62% in Longman's Dictionary of Contemporary English
 - 79% in WordNet
- But the others tend to have several senses
 - Average of 3.83 in LDOCE
 - Average of 2.96 in WordNet
- **Ambiguous words are more frequently used**
 - In the British National Corpus, 84% of instances have more than one sense
- **Some senses are more frequent than others**

Baseline Performance

- Baseline: most frequent sense
 - Equivalent to “take first sense” in WordNet
 - Does surprisingly well!

Freq	Synset	Gloss
338	plant ¹ , works, industrial plant	buildings for carrying on industrial labor
207	plant ² , flora, plant life	a living organism lacking the power of locomotion
2	plant ³	something planted secretly for discovery by another
0	plant ⁴	an actor situated in the audience whose acting is rehearsed but seems spontaneous to the audience

62% accuracy in this case!

Upper Bound Performance

- Upper bound
 - Fine-grained WordNet sense: 75-80% human agreement
 - Coarser-grained inventories: 90% human agreement possible

Simplest WSD algorithm: Lesk's Algorithm

- Intuition: note word overlap between context and dictionary entries
 - **Unsupervised**, but knowledge rich

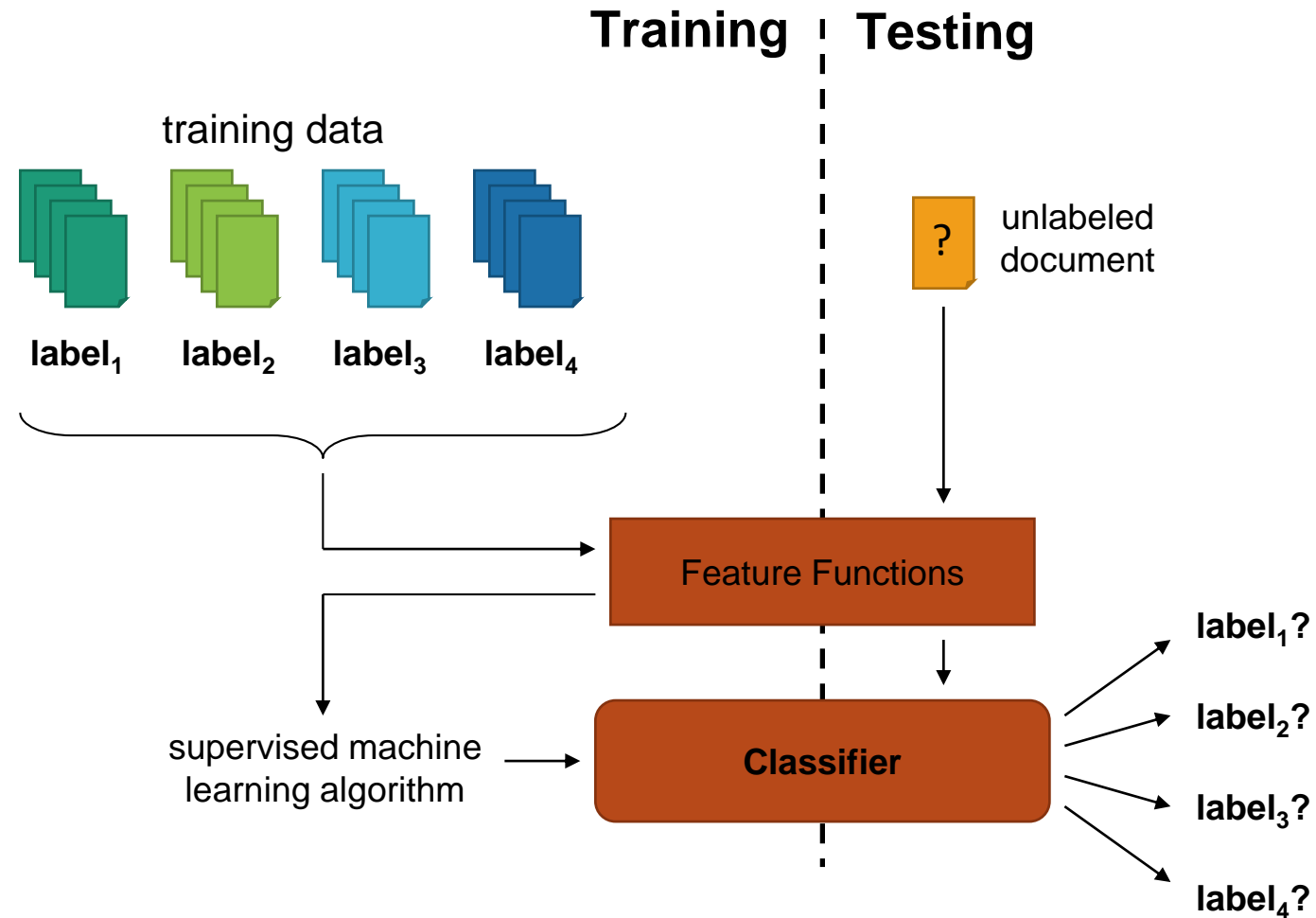
The bank can guarantee deposits will eventually cover future tuition costs because it invests in adjustable-rate mortgage securities.

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”

Lesk's Algorithm

- Simplest implementation:
 - Count overlapping content words between glosses and context
- Lots of variants:
 - Include the examples in dictionary definitions
 - Include hypernyms and hyponyms
 - Give more weight to larger overlaps
 - Give extra weight to infrequent words (e.g., using idf)
 - ...

Alternative: WSD as Supervised Classification



Existing Corpora

- Lexical sample
 - *line-hard-serve* corpus (4k sense-tagged examples)
 - *interest corpus* (2,369 sense-tagged examples)
 - ...
- All-words
 - SemCor (234k words, subset of Brown Corpus)
 - Senseval/SemEval (2081 tagged content words from 5k total words)
 - ...

How are annotated examples used in supervised learning?

- Supervised learning = requires examples annotated with correct prediction
- Used in 2 ways:
 - To find good values for the model (hyper)parameters (**training data**)
 - To evaluate how good the resulting classifier is (**test data**)
- How do we know how good a classifier is?
 - Compare classifier predictions with human annotation
 - On **held out** test examples
 - Evaluation metrics: accuracy, precision, recall

The 2-by-2 contingency table

	correct	not correct
selected	tp	fp
not selected	fn	tn

Precision and recall

- **Precision:** % of selected items that are correct
Recall: % of correct items that are selected

	correct	not correct
selected	tp	fp
not selected	fn	tn

A combined measure: F

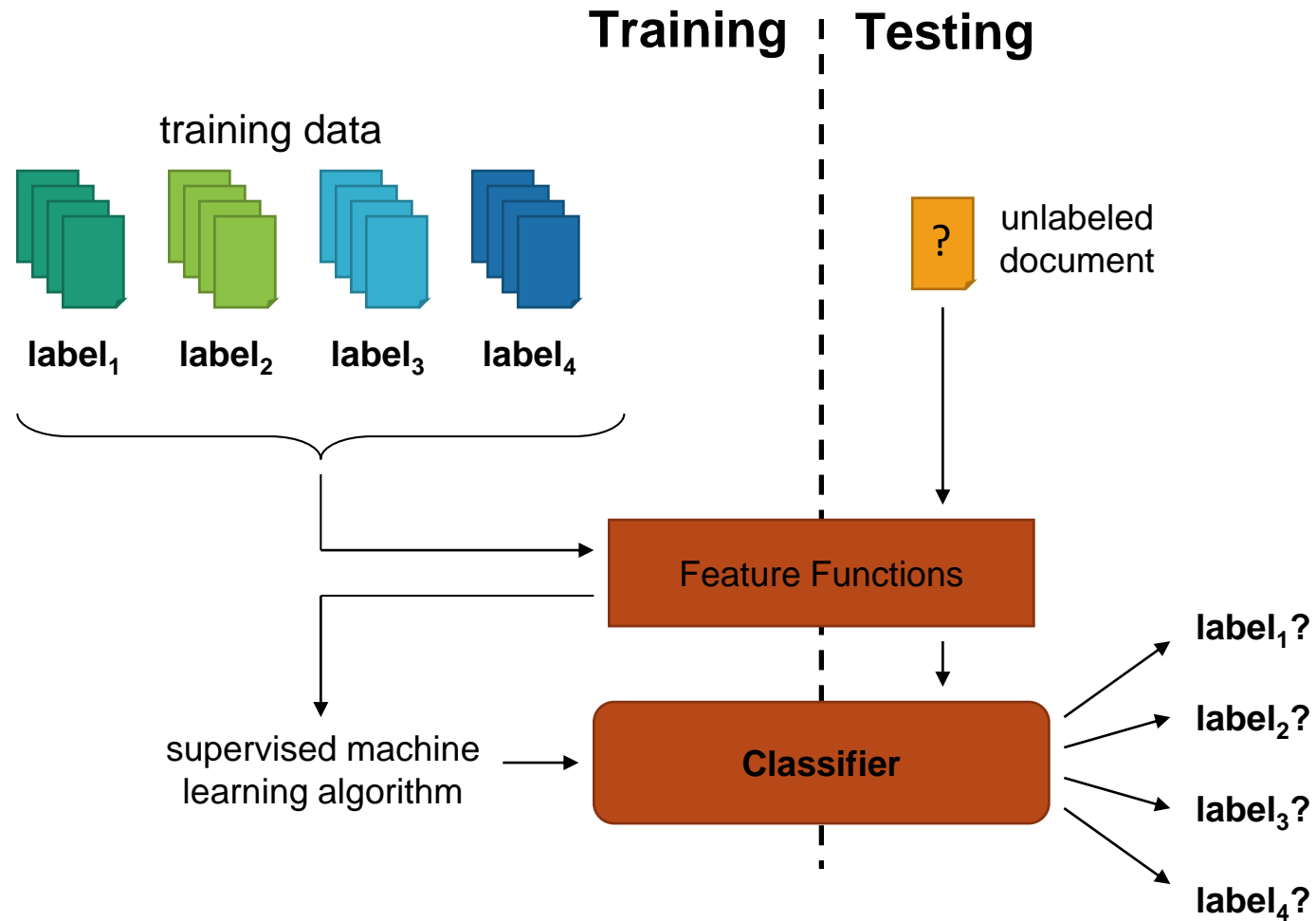
- A combined measure that assesses the P/R tradeoff is F measure (weighted harmonic mean):

$$F = \frac{1}{a \frac{1}{P} + (1-a) \frac{1}{R}} = \frac{(b^2 + 1)PR}{b^2P + R}$$

- People usually use balanced F1 measure
 - i.e., with $\beta = 1$ (that is, $\alpha = \frac{1}{2}$):

$$F = 2PR/(P+R)$$

Multiclass Classification



Is this spam?

From: "Fabian Starr"
<Patrick_Freeman@pamietaniepeerelu.pl>
Subject: Hey! Software for the funny prices!

Get the great discounts on popular software today
for PC and Macintosh

<http://iiled.org/Cj4Lmx>

70-90% Discounts from retail price!!!
All software is instantly available to download - No
Need Wait!

What is the subject of this article?

MEDLINE Article

Available on line at www.ncbi.nlm.nih.gov/pubmed/15482222

SCIENTIFIC DATA
www.nature.com/scientificdata

Brain
Cognition
www.nature.com/scientificdata

Syntactic frame and verb bias in aphasia: Plausibility judgments of underper- subject sentences

Suzanna Galil,¹ Lise Mann,² Gill Ransberger,³ David S. Justilly,² Elizabeth Elder,⁴ Molly Kavage,⁴ and L. Holland Aubrey⁴

¹University of Toronto, Toronto, ON, M5S 1A5
²University of Toronto, Toronto, ON, M5S 1A5
³University of Arizona, Tucson, AZ, 85724
RECEIVED 1 May 2015

Abstract

The study investigates how factors that have been argued to define “canonical form” in sentence comprehension (syntactic structure, semantic role, and frequency of usage) fit the evidence for what the sentence processing mechanism can do to predict difficult-to-analyze or “non-canonical” sentences. Using a plausibility judgment task, we show that a mixed group of aphasic participants performed differently from an executive function control group. For the core of its distribution the population generally favors the underper- subject frame. We show that this effect is modulated by word bias, i.e., the likelihood that a verb appears in a given syntactic structure. Phrases of positive bias verbs were significantly easier than phrases of negative bias verbs. More generally, we show that sentence structure modulates the degree to which the bias of the verb is borne significantly more than attention to what structure and word bias do on their own. These findings suggest that “canonical form” reflects frequency and word bias.

© 2015 The Author(s). All rights reserved.

1. Introduction

The simplicity of “canonical form” or “canonical word order” for normal and aphasic comprehension has often been taken as evidence in the sentence comprehension literature (Ransberger, as has been pointed out by Mann (2012)), the privileged status of canonical form itself needs explanation. In recent definitions of “canonical form” yield widely different predictions. One approach to the definition of canonical word order is that posited in Bates, Prizack, and Wallace (1987) (see also, Bates et al., 1996) that processes with Agent-Action-Object order represent the canonical word order for English. A second approach is based on syntactic “movement” analysis and defines as canonical any word order that diverges from the [NP]_i[Verb]_j[NP]_k configuration assumed for the deep structure of English sentences. Based on the understanding of animacy, King (1988) argues that sentences with unaccusative verbs should be difficult to process for aphasic patients, in particular for patients with “agrammatism,” for reasons that are analogous to the reasons giving rise to the greater difficulty of passive compared to active. Although the precise definition of unaccusativity is contested (see e.g., Levin & Rappaport Hovav, 1995), unaccusative verbs are generally understood to be intransitive verbs whose surface subject represents Underper- arguments. Examples of unaccusative verbs include verbs like melt and float. Under the transformational analysis assumed in King (1988), the surface subjects of unaccusative verbs are linked via movement to their objects in deep structure. Unaccusative verbs therefore inherit the very same difficulty as passive sentences, according to King’s analysis, and should be as hard as passives for aphasic speakers.

A different approach to canonical form has been proposed by Elder et al. (2009) who suggest that canonical form refers to the most frequent syntactic frame for a given verb. Under this view, aphasic problems with producing and understanding passives derive from the fact that, for most intransitive verbs, passives occur less frequently than active. The motivation of this approach, also advanced by Clark (2012), is that comprehension difficulty should vary with the lexical bias of the verbs

SCIENTIFIC DATA | www.nature.com/scientificdata

MeSH Subject Category Hierarchy

- Antagonists and Inhibitors
- Blood Supply
- Chemistry
- Drug Therapy
- Embryology
- Epidemiology
- ...

?

Text Classification

- Assigning subject categories, topics, or genres
- Spam detection
- Authorship identification
- Age/gender identification
- Language Identification
- Sentiment analysis
- ...

Word Sense Disambiguation

what you should know

- Word senses distinguish different meanings of same word
- Sense inventories
- Annotation issues and annotator agreement (Kappa)
- Definition of Word Sense Disambiguation Task
- An unsupervised approach: Lesk algorithm
- Supervised classification:
 - Train vs. test data
 - The most frequent class baseline
- Evaluation metrics: accuracy, precision, recall