

# Lexical Semantics

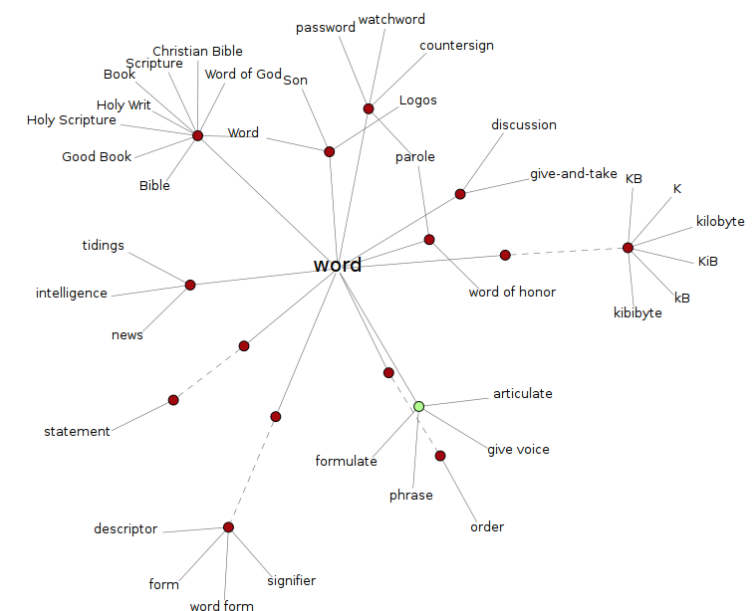
COMP90042

Natural Language Processing

Lecture 9



THE UNIVERSITY OF  
MELBOURNE



# Sentiment Analysis

- Bag of words, kNN classifier. Training data:
  - ▶ “This is a good movie.” → 😊
  - ▶ “This is a great movie.” → 😊
  - ▶ “This is a terrible film.” → 😞
- “This is a wonderful film.” → ?
- Two problems:
  - ▶ The model does not know that “movie” and “film” are synonyms. Since “film” appears only in negative examples the model learns that it is a negative word.
  - ▶ “wonderful” is not in the vocabulary (OOV – Out-Of-Vocabulary).

# Sentiment Analysis

- Comparing words directly will not work. How to make sure we compare word **meanings** instead?
- Solution: add this information explicitly through a **lexical database**.

# Word Semantics

- Lexical semantics (this lecture)
  - ▶ How the meanings of words connect to one another.
  - ▶ Manually constructed resources: lexicons, thesauri, ontologies, etc.
- Distributional semantics (next)
  - ▶ How words relate to each other in the text.
  - ▶ Automatically created resources from corpora.

# What Do Words Mean?

- Referents in the physical or social world
  - ▶ But not usually useful in text analysis
- Their dictionary definition
  - ▶ But dictionary definitions are necessarily circular
  - ▶ Only useful if meaning is already understood

red *n.* the color of blood or a ruby.  
blood *n.* the red liquid that circulates in the heart, arteries and veins of animals.

- Their relationships with other words
  - ▶ Also circular, but more practical

# Word Senses

- A word sense describes one aspect of the meaning of a word

**mouse**<sup>1</sup> : .... a *mouse* controlling a computer system in 1968.

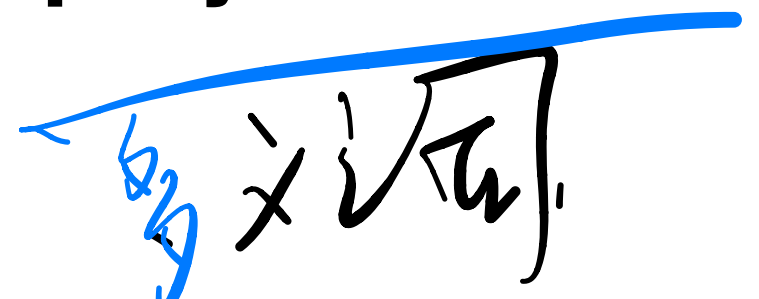
**mouse**<sup>2</sup> : .... a quiet animal like a *mouse*

**bank**<sup>1</sup> : ...a *bank* can hold the investments in a custodial account ...

**bank**<sup>2</sup> : ...as agriculture burgeons on the east *bank*, the river ...

# Word Glosses

- Gloss: textual definition of a sense given by a dictionary
- *Bank*:
  - ▶ financial institution that accepts deposits and channels the money into lending activities
  - ▶ sloping land (especially the slope beside a body of water)
- If a word has multiple senses, it is **polysemous**



# Meaning Through Relations

- Another way to define meaning: by looking at how it relates to other words
- **Synonymy:** near identical meaning
  - ▶ *vomit vs. throw up*
  - ▶ *big vs. large*
- **Antonymy:** opposite meaning
  - ▶ *long vs. short*
  - ▶ *big vs. little*



# Meaning Through Relations (2)

- **Hypernymy:** is-a relation

- ▶ cat is an animal
- ▶ *mango* is a *fruit*

animal is hypernymy  
of cat.

- **Meronymy:** part-whole relation

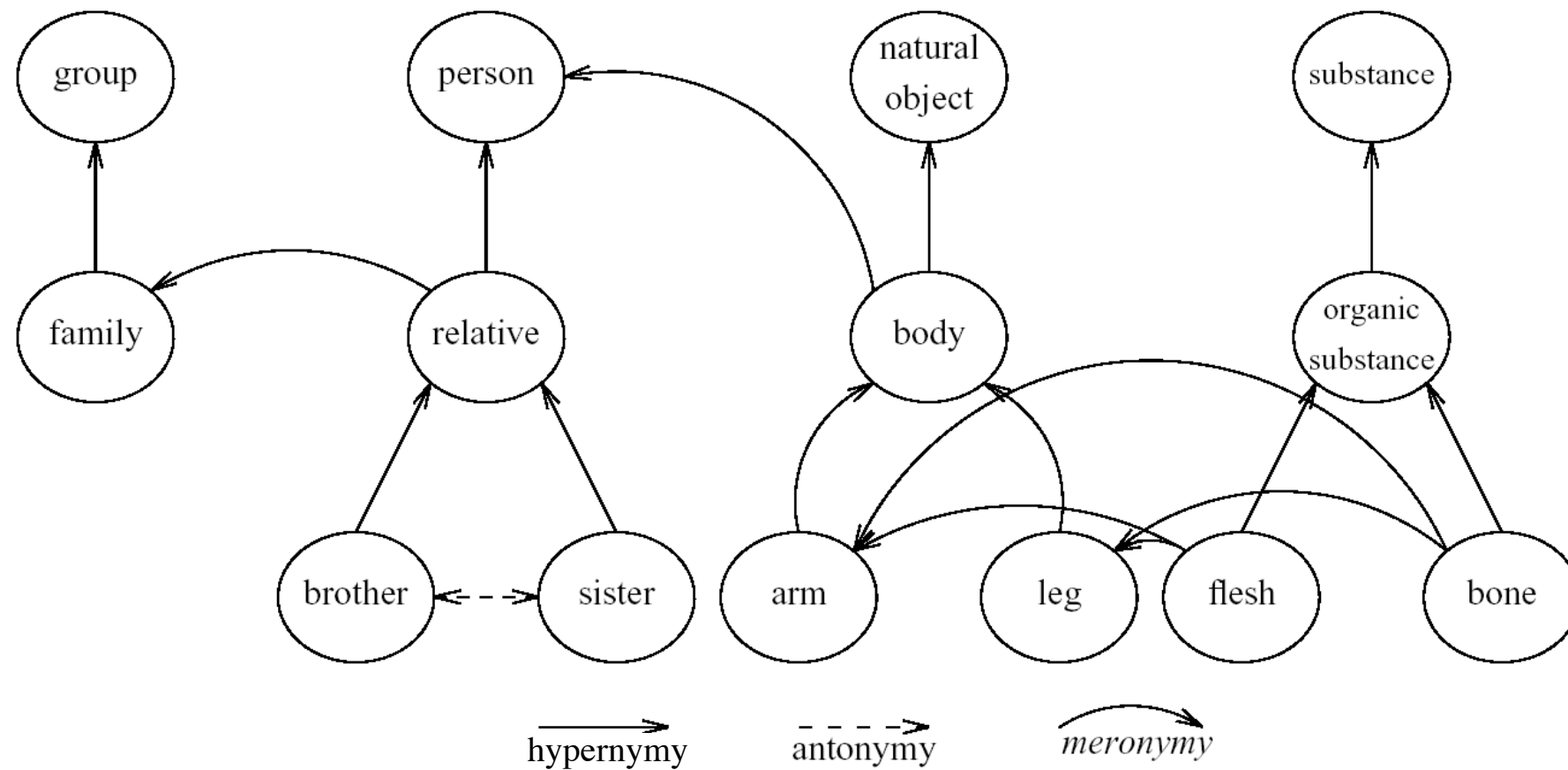
- ▶ *leg* is part of a *chair*
- ▶ *wheel* is part of a *car*

wheel - car

Meronymy.

legs - chair.

# Meaning Through Relations (3)



# WordNet

- A database of lexical relations
- English WordNet includes ~120,000 nouns, ~12,000 verbs, ~21,000 adjectives, ~4,000 adverbs
- On average: noun has 1.23 senses; verbs 2.16
- WordNets available in most major languages ([www.globalwordnet.org](http://www.globalwordnet.org), <https://babelnet.org/>)
- English version freely available (accessible via NLTK)

# WordNet Example

The noun “bass” has 8 senses in WordNet.

1. bass<sup>1</sup> - (the lowest part of the musical range)
2. bass<sup>2</sup>, bass part<sup>1</sup> - (the lowest part in polyphonic music)
3. bass<sup>3</sup>, basso<sup>1</sup> - (an adult male singer with the lowest voice)
4. sea bass<sup>1</sup>, bass<sup>4</sup> - (the lean flesh of a saltwater fish of the family Serranidae)
5. freshwater bass<sup>1</sup>, bass<sup>5</sup> - (any of various North American freshwater fish with lean flesh (especially of the genus Micropterus))
6. bass<sup>6</sup>, bass voice<sup>1</sup>, basso<sup>2</sup> - (the lowest adult male singing voice)
7. bass<sup>7</sup> - (the member with the lowest range of a family of musical instruments)
8. bass<sup>8</sup> - (nontechnical name for any of numerous edible marine and freshwater spiny-finned fishes)

*human perspective.*

# Synsets

*Synonym sets*

- Nodes of WordNet are not words or lemmas, but senses
- They are represented by sets of synonyms, or **synsets**
- *Bass* synsets:
  - ▶ {*bass*<sup>1</sup>, *deep*<sup>6</sup>}
  - ▶ {*bass*<sup>6</sup>, *bass voice*<sup>1</sup>, *basso*<sup>2</sup>}
- Another synset:
  - ▶ {*chump*<sup>1</sup>, *fool*<sup>2</sup>, *gull*<sup>1</sup>, *mark*<sup>9</sup>, *patsy*<sup>1</sup>, *fall guy*<sup>1</sup>, *sucker*<sup>1</sup>, *soft touch*<sup>1</sup>, *mug*<sup>2</sup>}
  - ▶ Gloss: a person who is gullible and easy to take advantage of

# Synsets (2)

>>> nltk.corpus.wordnet.synsets('bank')

[Synset('bank.n.01'), Synset('depository\_financial\_institution.n.01'), Synset('bank.n.03'), Synset('bank.n.04'), Synset('bank.n.05'), Synset('bank.n.06'), Synset('bank.n.07'), Synset('savings\_bank.n.02'), Synset('bank.n.09'), Synset('bank.n.10'), Synset('bank.v.01'), Synset('bank.v.02'), Synset('bank.v.03'), Synset('bank.v.04'), Synset('bank.v.05'), Synset('deposit.v.02'), Synset('bank.v.07'), Synset('trust.v.01')]

>>> nltk.corpus.wordnet.synsets('bank')[0].definition()

u'sloping land (especially the slope beside a body of water)'

>>> nltk.corpus.wordnet.synsets('bank')[1].lemma\_names()

[u'depository\_financial\_institution', u'bank', u'banking\_concern', u'banking\_company']

words participated  
with the word.

# Lexical Relations in WordNet

Relation	Also Called	Definition	Example
Hypernym	Superordinate	From concepts to superordinates	<i>breakfast</i> <sup>1</sup> $\rightarrow$ <i>meal</i> <sup>1</sup>
Hyponym	Subordinate	From concepts to subtypes	<i>meal</i> <sup>1</sup> $\rightarrow$ <i>lunch</i> <sup>1</sup>
Instance Hypernym	Instance	From instances to their concepts	<i>Austen</i> <sup>1</sup> $\rightarrow$ <i>author</i> <sup>1</sup>
Instance Hyponym	Has-Instance	From concepts to their instances	<i>composer</i> <sup>1</sup> $\rightarrow$ <i>Bach</i> <sup>1</sup>
Part Meronym	Has-Part	From wholes to parts	<i>table</i> <sup>2</sup> $\rightarrow$ <i>leg</i> <sup>3</sup>
Part Holonym	Part-Of	From parts to wholes	<i>course</i> <sup>7</sup> $\rightarrow$ <i>meal</i> <sup>1</sup>
Antonym		Semantic opposition between lemmas	<i>leader</i> <sup>1</sup> $\iff$ <i>follower</i> <sup>1</sup>
Derivation		Lemmas w/same morphological root	<i>destruction</i> <sup>1</sup> $\iff$ <i>destroy</i> <sup>1</sup>

# Hypernymy Chain

bass<sup>3</sup>, basso (an adult male singer with the lowest voice)

=> singer, vocalist, vocalizer, vocaliser

=> musician, instrumentalist, player

=> performer, performing artist

=> entertainer

=> person, individual, someone...

=> organism, being

=> living thing, animate thing,

=> whole, unit

=> object, physical object

=> physical entity

=> entity

bass<sup>7</sup> (member with the lowest range of a family of instruments)

=> musical instrument, instrument

=> device

=> instrumentality, instrumentation

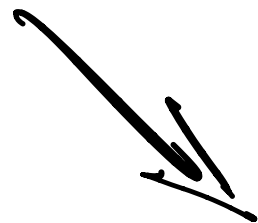
=> artifact, artefact

=> whole, unit

=> object, physical object

=> physical entity

=> entity





# Word Similarity

# Word Similarity

- Synonymy: *film* vs. *movie*
- What about *show* vs. *film*? *opera* vs. *film*?
- Unlike synonymy (which is a binary relation), word similarity is a spectrum
- We can use lexical database (e.g. WordNet) or thesaurus to estimate word similarity

# Word Similarity with Paths

- Given WordNet, find similarity based on path length
- $\text{pathlen}(c_1, c_2) = 1 + \text{edge length}$  in the shortest path between sense  $c_1$  and  $c_2$ 

*in word net.*
- similarity between two senses:
 
$$\text{simpath}(c_1, c_2) = \frac{1}{\text{pathlen}(c_1, c_2)}$$
- similarity between two words
 
$$\text{wordsim}(w_1, w_2) = \max_{c_1 \in \text{senses}(w_1), c_2 \in \text{senses}(w_2)} \text{simpath}(c_1, c_2)$$

*of all senses, what is the pair closest to each other*

*all senses.*

# Examples

$$\text{simpath}(c_1, c_2) = \frac{1}{\text{pathlen}(c_1, c_2)}$$

$$\text{simpath}(\textit{nickel}, \textit{coin}) = 1/2 = 0.5$$

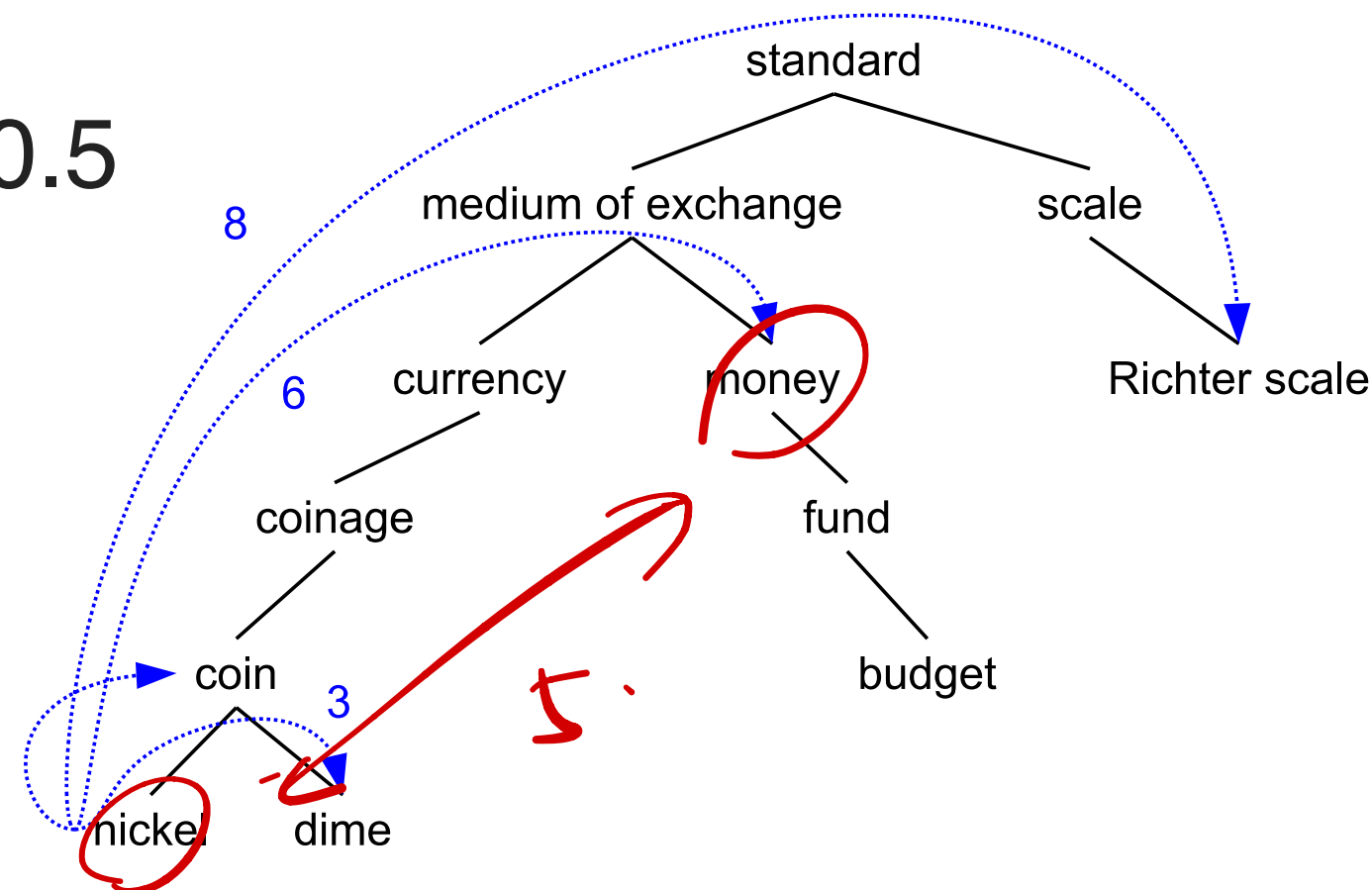
$$\text{simpath}(\textit{nickel}, \textit{currency}) = 1/4 = 0.25$$

$$\text{simpath}(\textit{nickel}, \textit{money}) = 1/6 = 0.17$$

$$\text{simpath}(\textit{nickel}, \textit{Richter scale}) = 1/8 = 0.13$$

nearest 1.

Worthnote.



$$\text{simpaths}(\textit{nickel}, \textit{nickel}) = \frac{1}{1} = 1$$

# Beyond Path Length

- Problem: edges vary widely in actual semantic distance
  - ▶ Much bigger jumps near top of hierarchy
- Solution 1: include depth information (Wu & Palmer)
  - ▶ Use path to find lowest common subsumer (LCS)
  - ▶ Compare using depths

*lowest common ancestor*

$$\text{simwup}(c_1, c_2) = \frac{2 * \text{depth}(\text{LCS}(c_1, c_2))}{\text{depth}(c_1) + \text{depth}(c_2)}$$

*Wu & Palmer*

# Examples

$$\text{simwup}(c_1, c_2) = \frac{2 * \text{depth}(\text{LCS}(c_1, c_2))}{\text{depth}(c_1) + \text{depth}(c_2)}$$

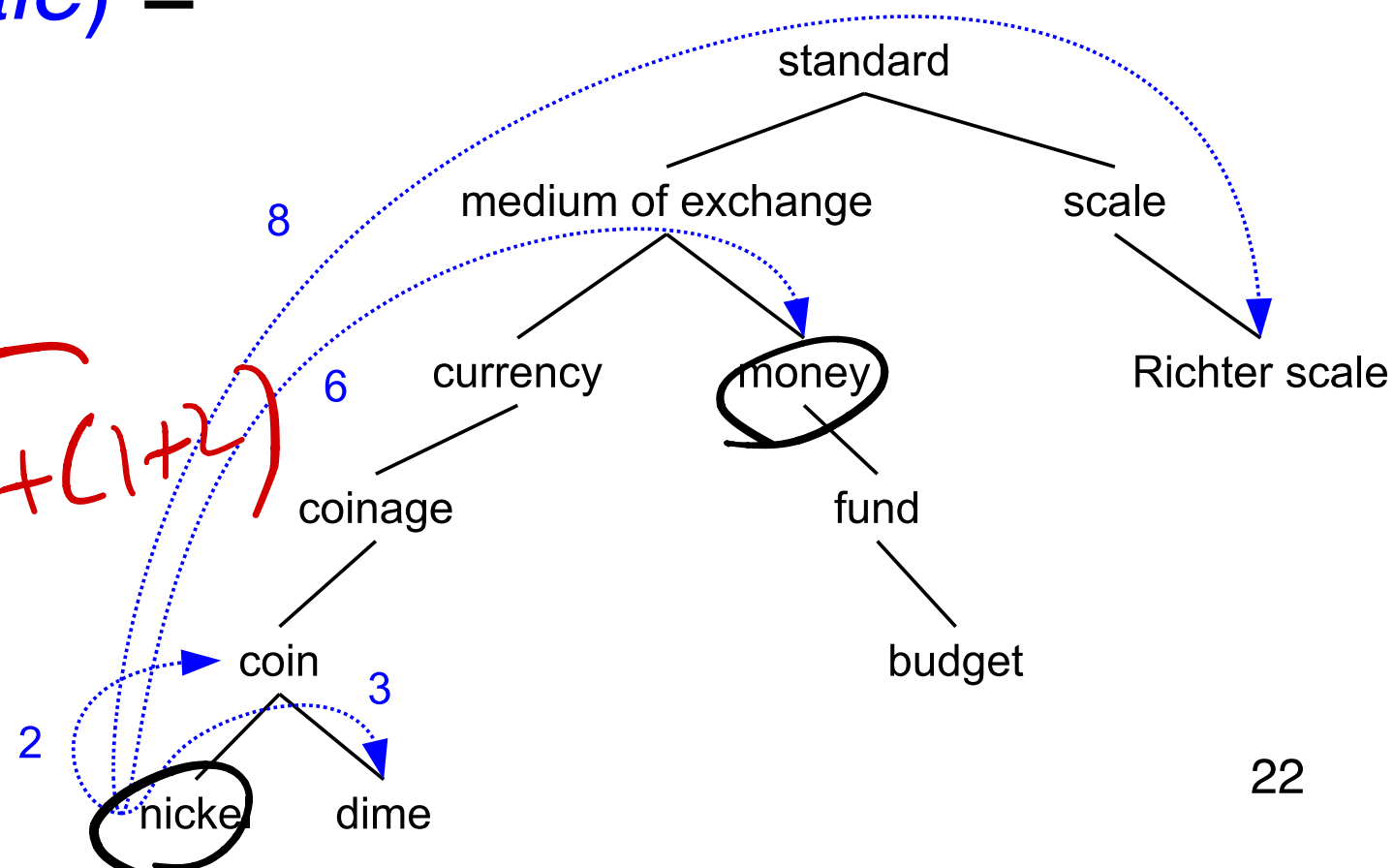
LCS.

$$\text{simwup}(\text{nickel}, \text{money}) = 2 * 2 / (3 + 6) = 0.44$$

$$\text{simwup}(\text{nickel}, \text{Richter scale}) = 2 * 1 / (3 + 6) = 0.22$$

$$\frac{2(1+1)}{(1+2)+(1+5)}$$

$$\frac{2}{(5+1)+(1+2)} = \frac{2}{9}$$



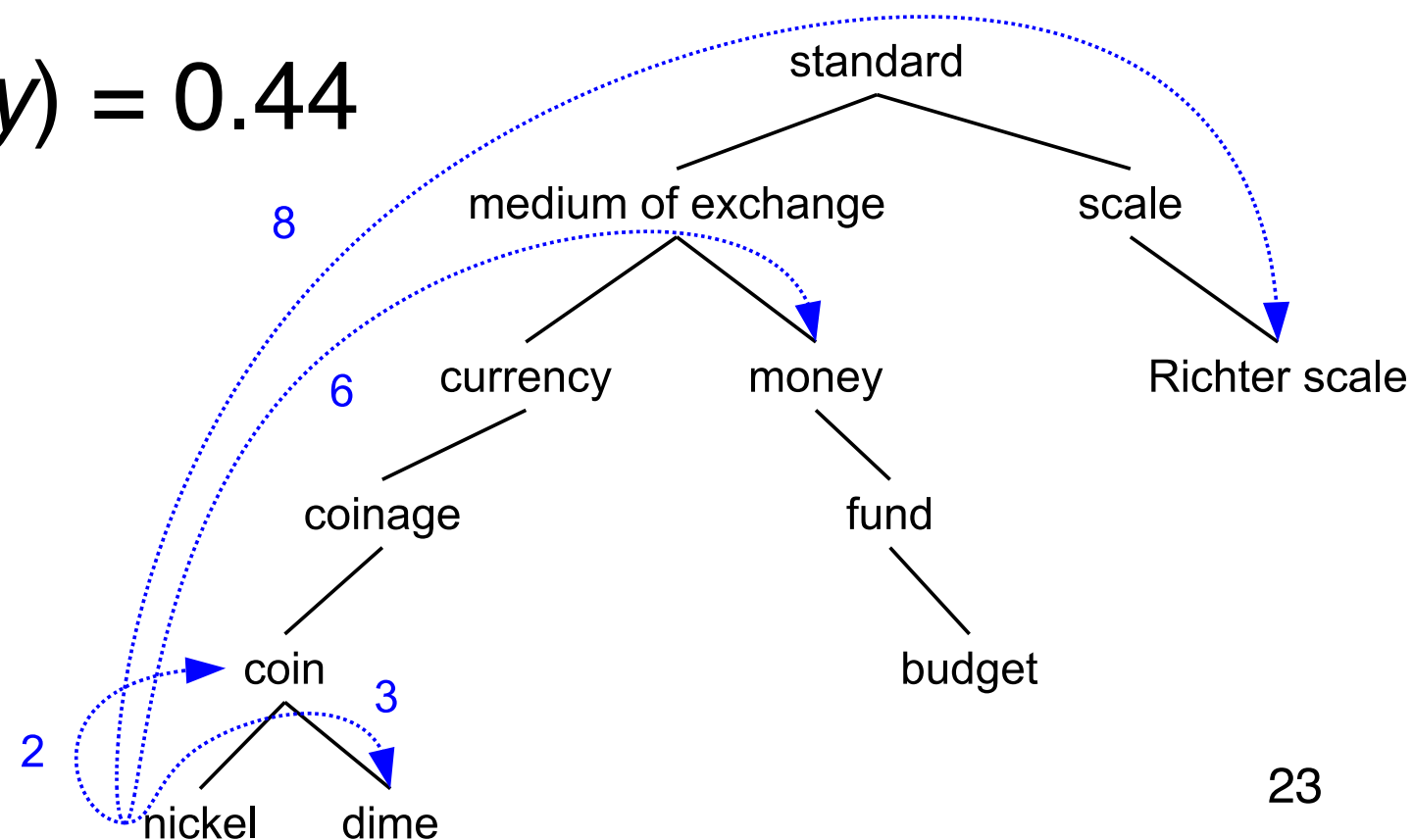
# Abstract Nodes

*penalize words go through abstract words.*

- But count of edges or node depth is still poor semantic distance metric
- Nodes high in the hierarchy is very abstract/general
- How do we make words that connect through very abstract nodes much less similar

►  $\text{simwup}(\text{nickel}, \text{money}) = 0.44$

►  $\text{simwup}(\text{nickel}, \text{Richter scale}) = 0.22$



# Concept Probability

belongs to

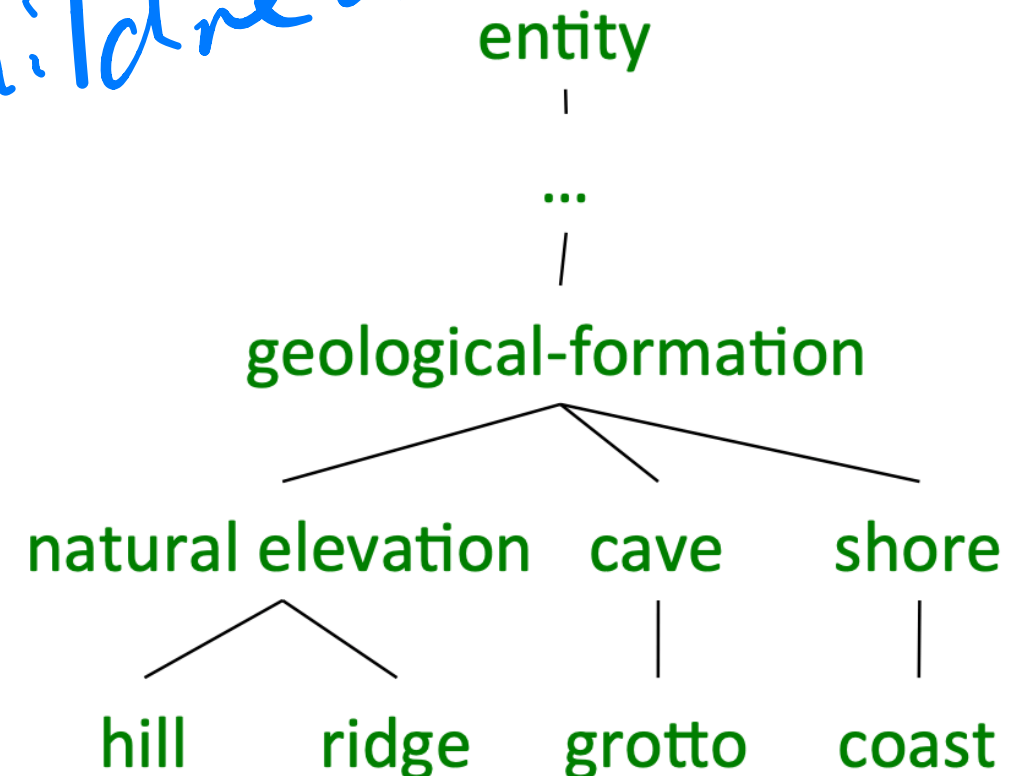
$$P(c) = \frac{\sum_{w \in \text{words}(c)} \text{count}(w)}{N}$$

# words instead of current word.

all corpus counts

- $P(c)$ : probability that a randomly selected word in a corpus is an instance of concept  $c$
- $\text{words}(c)$ : set of all words that are children of  $c$
- $\text{words}(\text{geological-formation}) = \{\text{hill, ridge, grotto, coast, natural elevation, cave, shore}\}$
- $\text{words}(\text{natural elevation}) = \{\text{hill, ridge}\}$

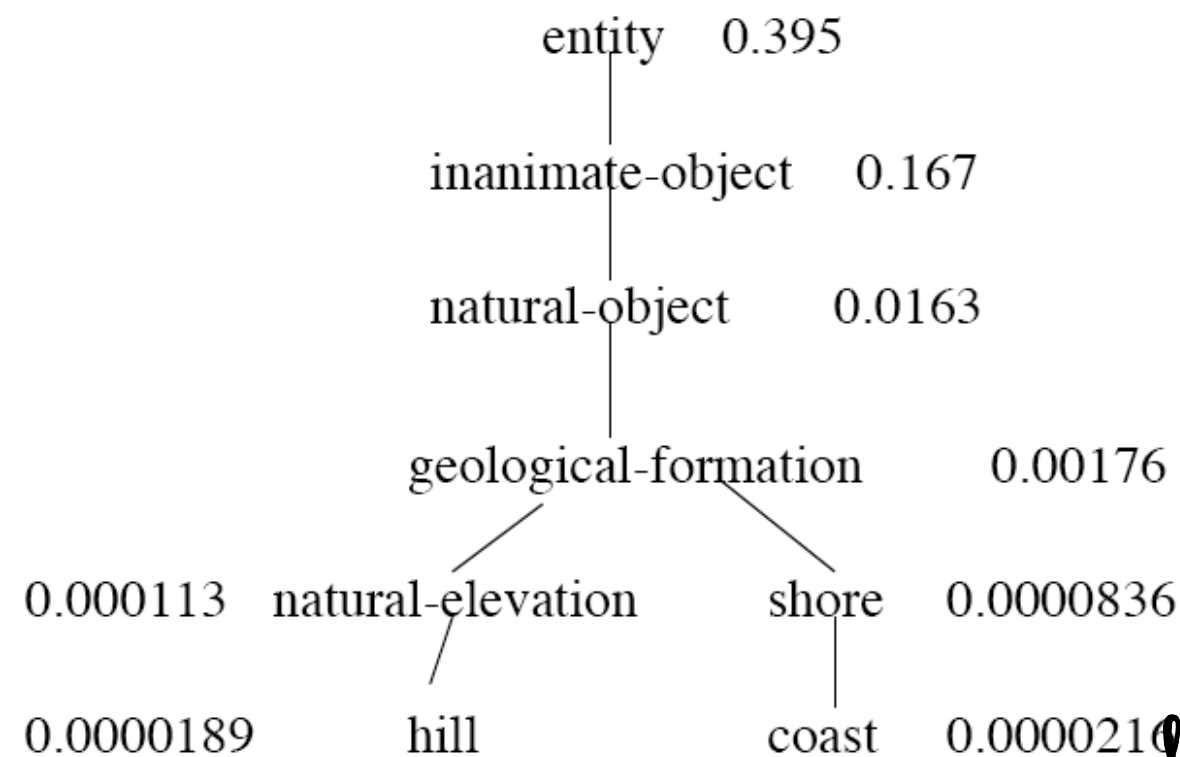
children





# Example

- Abstract nodes higher in the hierarchy has a higher  $P(c)$



Penalize high cross road indirectly.

# Similarity with Information Content

Concept-  
probability

$$IC(c) = -\log P(c)$$

use IC instead of depth in simwup

$$\text{simlin}(c_1, c_2) = \frac{2 \times IC(LCS(c_1, c_2))}{IC(c_1) + IC(c_2)}$$

$$\text{simlin}(\text{hill}, \text{coast}) = \frac{2 \times -\log P(\text{geological-formation})}{-\log P(\text{hill}) - \log P(\text{coast})}$$

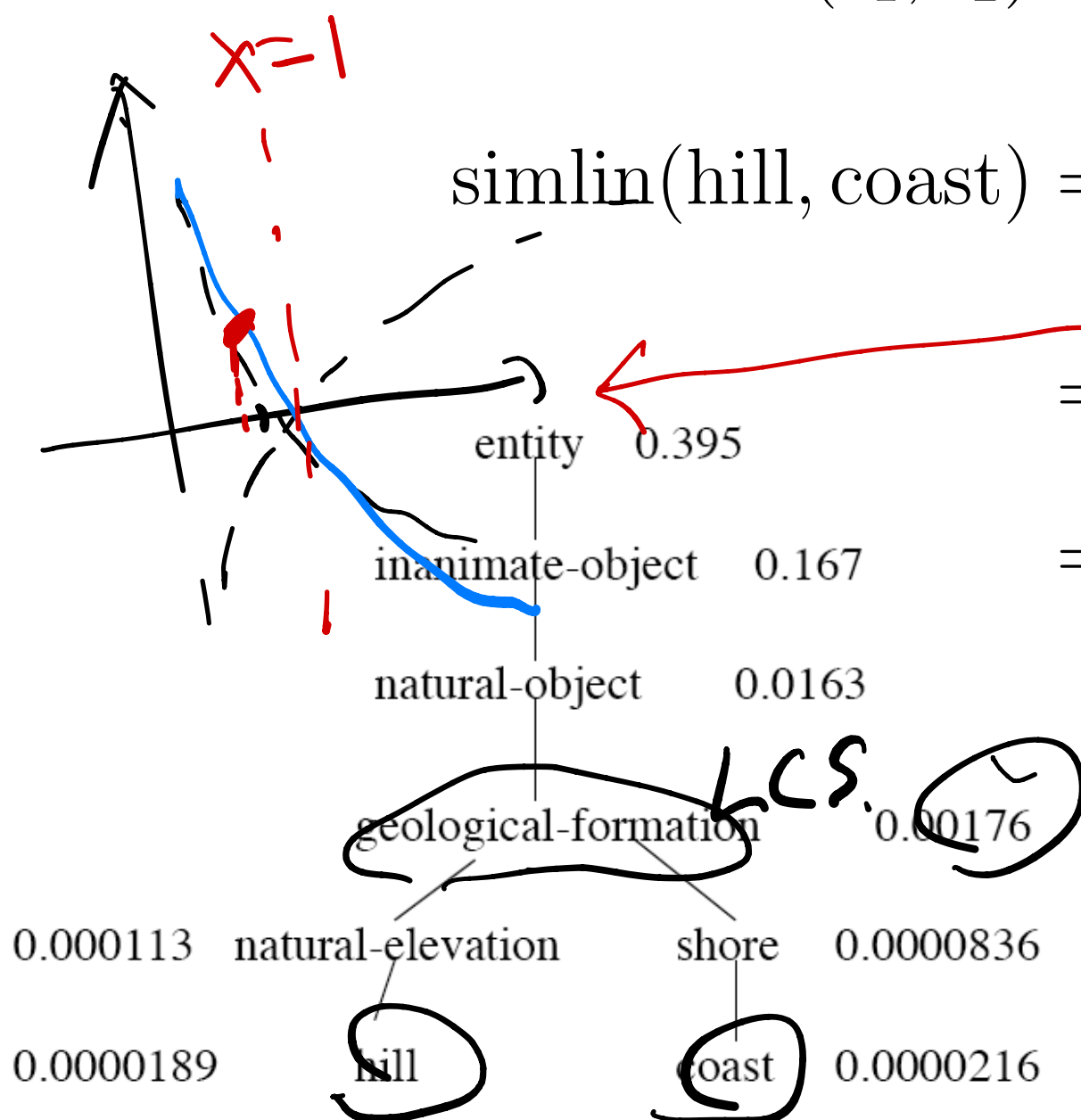
$$= \frac{-2 \log 0.00176}{-\log 0.0000189 - \log 0.0000216}$$

$$= 0.587$$

Catch Low LCS

if LCS node is very high up in the hierarchy  
(say  $P(c) = 0.99$ ), then IC will be very low  
(0.01 in this case)

penalty



# Sentiment Analysis Revisited

- “This is a great movie.” → 😊
- “This is a wonderful film.” → ?
- Comparing words using WordNet paths work well if our classifier is based on word similarities (such as **kNN**)
- But what if we want sense as a general feature representation, so we can employ other classifiers?
- **Solution** **map** words in text to **senses** in WordNet explicitly.

*Logistic Regression*

# Word Sense Disambiguation

- Task: selects the correct sense for words in a sentence
- Baseline:
  - ▶ Assume the most popular sense *predominant sense.*
- Good WSD potentially useful for many tasks in NLP
  - ▶ In practice, often ignored because good WSD too hard
  - ▶ Active research area

# Supervised WSD

- Apply standard machine classifiers
- Feature vectors typically words and syntax around target
  - ▶ But context is ambiguous too!
  - ▶ How big should context window be? (typically very small)
- Requires sense-tagged corpora
  - ▶ E.g. SENSEVAL, SEMCOR (available in NLTK)
  - ▶ Very time consuming to create!

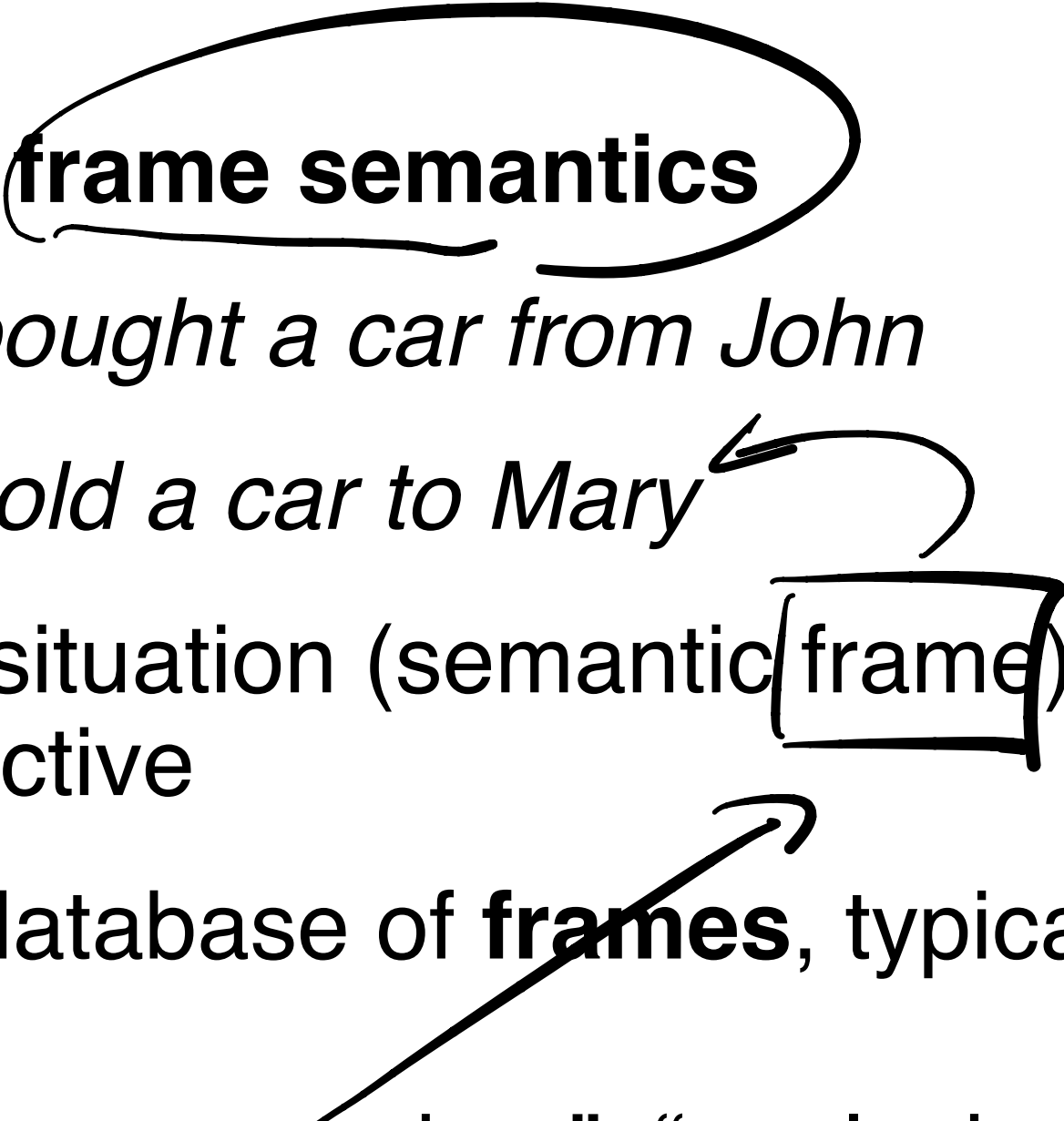
# Less Supervised Approaches

- **Lesk**: Choose sense whose dictionary gloss from WordNet most overlaps with the context *overlap.*
- *The **bank** can guarantee deposits will eventually cover future tuition costs because it invests in adjustable-rate mortgage securities.*
- **bank**<sup>1</sup>: 2 overlapping non-stopwords, *deposits* and *mortgage*
- **bank**<sup>2</sup>: 0

bank <sup>1</sup>	Gloss:	a financial institution that accepts <u>deposits</u> and channels the money into lending activities
	Examples:	“he cashed a check at the bank”, “that bank holds the <u>mortgage</u> on my home”
bank <sup>2</sup>	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”



# Other Databases - FrameNet

- Based on **frame semantics**
    - ▶ *Mary bought a car from John*
    - ▶ *John sold a car to Mary*
    - ▶ Same situation (semantic **frame**), just different perspective
  - A lexical database of **frames**, typically prototypical situations
    - ▶ E.g. “commerce\_buy”, “apply\_heat”
- 

# FrameNet

- Includes lists of *lexical units* that *evoke* the frame
  - ▶ E.g. *cook, fry, bake, boil*, etc.
- Lists of *semantic roles* or *frame elements*
  - ▶ E.g. “the cook”, “the food”, “the container”, “the instrument”
- Semantic *relationships* among frames
  - ▶ “*apply\_heat*” is Causative of “*absorb\_heat*”, is Used by “*cooking\_creation*”



# Moving On To The Corpus

- Manually-tagged lexical resources an important starting point for text analysis
- But much modern work attempts to derive semantic information directly from corpora, without human intervention
- Distributional semantics!

# Reading

- JM3 Ch 19.1-19.3, 19.4.1, 19.5.1