

Word Senses and WordNet

(Thesaurus-based)

Word Senses and Relations

- Homonymy, Polysemy, Synonymy, and more
- Online Resources
- Thesaurus methods for word similarity

Lexical Semantics

- Focus on word meanings:
 - **Relations of meaning among words**
 - Similarities & differences of meaning in similar context
 - Internal meaning structure of words
 - Basic internal units combine for meaning

Word Definitions

- What's a word?
 - Definitions so far: Types, tokens, stems, roots, inflected forms, etc...
 - **Lexeme**: An entry in a lexicon consisting of a pairing of a form with a single meaning representation
 - **Lexicon**: A collection of lexemes

I. Possible Word Relations

- Homonymy
- Polysemy
- Synonymy
- Antonymy
- Hypernymy
- Hyponymy
- Meronymy

Homonymy

- Lexemes share a form
 - Phonological, orthographic or both
 - But have unrelated, distinct meanings
- Clear examples
 - **Bat** (wooden stick-like thing) vs. **bat** (flying scary mammal thing)
 - **Bank** (financial institution) versus **bank** (riverside)
- Can be homophones, homographs:
 - Homophones:
 - Write/right, piece/peace, to/too/two
 - Homographs:
 - Desert/desert
 - Bass/bass

Issues for NLP Applications

- Text-to-Speech
 - Same orthographic form but different phonological form
 - **bass** vs. **bass**
- Information retrieval
 - Different meanings same orthographic form
 - QUERY: **bat care**
- Machine Translation (English -> Spanish)
 - bat: **murciélago** (animal) or **bate** (for baseball)

Polysemy

- The **bank** is constructed from red brick
I withdrew the money from the **bank**
 - Are these the same sense? Different?
- Or consider the following WSJ example
 - **While some banks furnish sperm only to married women, others are less restrictive**
 - Which sense of bank is this?
 - Is it distinct from the river bank sense?
 - The savings bank sense?

Polysemy

- A single lexeme with multiple **related** meanings (**bank** the building, **bank** the financial institution)
- Most non-rare words have multiple meanings
 - Number of meanings related to word frequency
 - Verbs tend more to polysemy
 - Distinguishing polysemy from homonymy isn't always easy (or necessary)

Metonymy or Systematic Polysemy: A systematic relationship between senses

- Lots of types of polysemy are systematic
 - School, university, hospital
 - All can mean the institution or the building.
- A systematic relationship:
 - **Building** ↔ **Organization**
- Other such kinds of systematic polysemy:

Author (Jane Austen wrote Emma)

↔ **Works of Author** (I love Jane Austen)

Tree (Plums have beautiful blossoms)

↔ **Fruit** (I ate a preserved plum)

Metaphor vs. Metonymy

- **Metaphor**: two different meaning domains are related
 - **Citibank claimed it was misrepresented.**
 - Corporation as person
- **Metonymy**: use of one aspect of a concept to refer to other aspects of entity or to entity itself
 - **The Citibank is on the corner of Main and State.**
 - Building stands for organization

How Do We Identify Words with Multiple Senses?

- ATIS examples
 - **Which flights *serve* breakfast?**
 - **Does America West *serve* Philadelphia?**
- The “zeugma” test: conjoin two potentially similar/dissimilar senses
 - **?Does United *serve* breakfast and San Jose?**
 - **Does United *serve* breakfast and lunch?**

Synonymy

- Word that have the same meaning in some or all contexts.
 - **filbert / hazelnut**
 - **couch / sofa**
 - **big / large**
 - **automobile / car**
 - **vomit / throw up**
 - **Water / H₂O**
- Two lexemes are synonyms if they can be successfully substituted for each other in all situations
 - If so they have the same **propositional meaning**

Few Examples of Perfect Synonymy

- Even if many aspects of meaning are identical
 - Still may not preserve the acceptability based on notions of politeness, slang, register, genre, etc.
- E.g, **water** and **H₂O**, **coffee** and **java**

Terminology

- Lemmas and wordforms
 - A **lexeme** is an abstract pairing of meaning and form
 - A **lemma** or citation form is the grammatical form that is used to represent a lexeme.
 - **Carpet** is the lemma for **carpets**
 - Specific surface forms **carpets**, **sung** are called wordforms
- The lemma **bank** has two senses:
 - **Instead, a bank can hold the investments in a custodial account in the client's name.**
 - **But as agriculture burgeons on the east bank, the river will shrink even more.**
- A sense is a discrete representation of one aspect of the meaning of a word

Synonymy Relates Senses not Words

- Consider *big* and *large*
- Are they synonyms?
 - How **big** is that plane?
 - Would I be flying on a **large** or a small plane?
- How about:
 - Miss Nelson, for instance, became a kind of **big** sister to Benjamin.
 - ?Miss Nelson, for instance, became a kind of **large** sister to Benjamin.
- Why?
 - *big* has a sense that means being older, or grown up
 - *large* lacks this sense

Antonyms

- Senses that are *opposites* with respect to one feature of their meaning
- Otherwise, they are very similar
 - dark / light
 - short / long
 - hot / cold
 - up / down
 - in / out
- More formally: antonyms can
 - Define a binary opposition or an attribute at opposite ends of a scale (*long/short, fast/slow*)
 - Be reversives: *rise/fall, up/down*

Hyponyms

- A sense is a **hyponym** of another if the first sense is more specific, denoting a subclass of the other
 - *car* is a hyponym of *vehicle*
 - *dog* is a hyponym of *animal*
 - *mango* is a hyponym of *fruit*
- Conversely
 - *vehicle* is a hypernym/superordinate of *car*
 - *animal* is a hypernym of *dog*
 - *fruit* is a hypernym of *mango*

superordinate	vehicle	fruit	furniture	mammal
hyponym	car	mango	chair	dog

Hypernymy Defined

- Extensional
 - The class denoted by the **superordinate**
 - Extensionally includes class denoted by the **hyponym**
- Entailment
 - *A sense A is a hyponym of sense B if being an A entails being a B*
- Hyponymy is usually *transitive*
 - (A hypo B and B hypo C entails A hypo C)
- Another name: the **IS-A hierarchy**
 - A **IS-A** B (or A **ISA** B)
 - B **subsumes** A

Meronymy

- The part-whole relation
 - A *leg* is part of a *chair*; a *wheel* is part of a *car*.
- *Wheel* is a **meronym** of *car*, and *car* is a **holonym** of *wheel*.

II. WordNet

- A hierarchically organized lexical database
- On-line thesaurus + aspects of a dictionary
 - Versions for other languages are under development

Category	Unique Forms
Noun	117,798
Verb	11,529
Adjective	22,479
Adverb	4,481

Hyponyms and Instances

- WordNet has both **classes** and **instances**.
- An **instance** is an individual, a proper noun that is a unique entity
 - San Francisco is an **instance** of city
 - But city is a class
 - city is a **hyponym** of municipality...location...

WordNet Entries

The noun “bass” has 8 senses in WordNet.

1. bass¹ - (the lowest part of the musical range)
2. bass², bass part¹ - (the lowest part in polyphonic music)
3. bass³, basso¹ - (an adult male singer with the lowest voice)
4. sea bass¹, bass⁴ - (the lean flesh of a saltwater fish of the family Serranidae)
5. freshwater bass¹, bass⁵ - (any of various North American freshwater fish with lean flesh (especially of the genus *Micropterus*))
6. bass⁶, bass voice¹, basso² - (the lowest adult male singing voice)
7. bass⁷ - (the member with the lowest range of a family of musical instruments)
8. bass⁸ - (nontechnical name for any of numerous edible marine and freshwater spiny-finned fishes)

The adjective “bass” has 1 sense in WordNet.

1. bass¹, deep⁶ - (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 Noun Relations

Relation	Also called	Definition	Example
Hypernym	Superordinate	From concepts to superordinates	<i>breakfast</i> ¹ → <i>meal</i> ¹
Hyponym	Subordinate	From concepts to subtypes	<i>meal</i> ¹ → <i>lunch</i> ¹
Member Meronym	Has-Member	From groups to their members	<i>faculty</i> ² → <i>professor</i> ¹
Has-Instance		From concepts to instances of the concept	<i>composer</i> ¹ → <i>Bach</i> ¹
Instance		From instances to their concepts	<i>Austen</i> ¹ → <i>author</i> ¹
Member Holonym	Member-Of	From members to their groups	<i>copilot</i> ¹ → <i>crew</i> ¹
Part Meronym	Has-Part	From wholes to parts	<i>table</i> ² → <i>leg</i> ³
Part Holonym	Part-Of	From parts to wholes	<i>course</i> ⁷ → <i>meal</i> ¹
Antonym		Opposites	<i>leader</i> ¹ → <i>follower</i> ¹

WordNet Verb Relations

Relation	Definition	Example
Hypernym	From events to superordinate events	<i>fly^y → travel^p</i>
Troponym	From a verb (event) to a specific manner elaboration of that verb	<i>walk¹ → stroll¹</i>
Entails	From verbs (events) to the verbs (events) they entail	<i>snore¹ → sleep¹</i>
Antonym	Opposites	<i>increase¹ ↔ decrease¹</i>

WordNet Hierarchies

```

Sense 3
bass, 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
                                            => entity
                                                => entity
                                                    => entity
                                                        => entity
                                                            => entity
                                                                => entity
                                                                    => entity

```

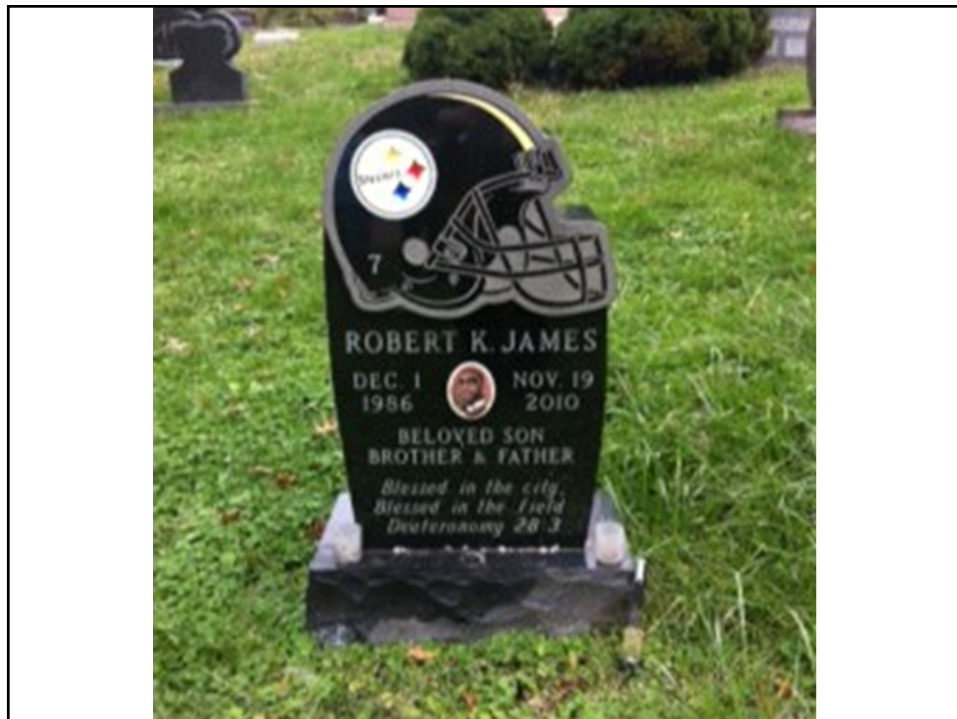
```

Sense 7
bass --
(the member with the lowest range of a family of
musical instruments)
=> musical instrument, instrument
    => device
        => instrumentality, instrumentation
            => artifact, artefact
                => whole, unit
                    => object, physical object
                        => physical entity
                            => entity

```

How is 'Sense' Defined in WordNet?

- The set of near-synonyms for a WordNet sense is called a **synset** (synonym set); their version of a sense or a concept
- Example: **chump** as a noun to mean 'a person who is gullible and easy to take advantage of'
{chump¹, fool², gull¹, mark⁹, patsy¹, fall guy¹, sucker¹, soft touch¹, mug²}
- Each of these senses share this same gloss (but not every sense)
- For WordNet, the meaning of this sense of **chump** is this list.



S: (n) **field** (a piece of land cleared of trees and usually enclosed)

S: (n) battlefield, battleground, field of battle, field of honor, **field** (a region where a battle is being (or has been) fought)

S: (n) **field** (somewhere (away from a studio or office or library or laboratory) where practical work is done or data is collected)

S: (n) discipline, subject, subject area, subject field, **field**, field of study, study, bailiwick (a branch of knowledge)

S: (n) **field**, field of force, force field (the space around a radiating body within which its electromagnetic oscillations can exert force on another similar body not in contact with it)

S: (n) **field**, field of operation, line of business (a particular kind of commercial enterprise)

S: (n) sphere, domain, area, orbit, **field**, arena (a particular environment or walk of life)

S: (n) playing field, athletic field, playing area, **field** (a piece of land prepared for playing a game)

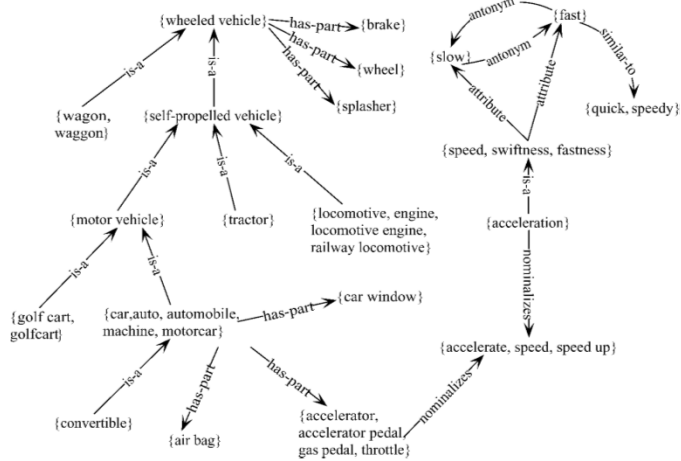


Time flies.

(thanks to Dr. Wiebe for Allegheny Cemetery photos)

S: (v) fly, wing (travel through the air; be airborne) "Man cannot fly"
 S: (v) fly (move quickly or suddenly) "He flew about the place"
 S: (v) fly, aviate, pilot (operate an airplane) "The pilot flew to Cuba"
 S: (v) fly (transport by aeroplane) "We fly flowers from the Caribbean to North America"
 S: (v) fly (cause to fly or float) "fly a kite"
 S: (v) fly (be dispersed or disseminated) "Rumors and accusations are flying"
 S: (v) fly (change quickly from one emotional state to another) "fly into a rage"
 S: (v) fly, fell, vanish (pass away rapidly) "Time flies like an arrow"; "Time fleeing beneath him"

WordNet: Viewed as a graph



“Supersenses”

The top level hypernyms in the hierarchy

A word’s supersense can be a useful coarse-grained representation of word meaning for NLP tasks (counts from Schneider and Smith 2013’s Streusel corpus)

Noun				Verb	
GROUP	1469 <i>place</i>	BODY	87 <i>hair</i>	STATIVE	2922 <i>is</i>
PERSON	1202 <i>people</i>	STATE	56 <i>pain</i>	COGNITION	1093 <i>know</i>
ARTIFACT	971 <i>car</i>	NATURAL OBJ.	54 <i>flower</i>	COMMUNIC.*	974 <i>recommend</i>
COGNITION	771 <i>way</i>	RELATION	35 <i>portion</i>	SOCIAL	944 <i>use</i>
FOOD	766 <i>food</i>	SUBSTANCE	34 <i>oil</i>	MOTION	602 <i>go</i>
ACT	700 <i>service</i>	FEELING	34 <i>discomfort</i>	POSSESSION	309 <i>pay</i>
LOCATION	638 <i>area</i>	PROCESS	28 <i>process</i>	CHANGE	274 <i>fix</i>
TIME	530 <i>day</i>	MOTIVE	25 <i>reason</i>	EMOTION	249 <i>love</i>
EVENT	431 <i>experience</i>	PHENOMENON	23 <i>result</i>	PERCEPTION	143 <i>see</i>
COMMUNIC.*	417 <i>review</i>	SHAPE	6 <i>square</i>	CONSUMPTION	93 <i>have</i>
POSSESSION	339 <i>price</i>	PLANT	5 <i>tree</i>	BODY	82 <i>get...done</i>
ATTRIBUTE	205 <i>quality</i>	OTHER	2 <i>stuff</i>	CREATION	64 <i>cook</i>
QUANTITY	102 <i>amount</i>			CONTACT	46 <i>put</i>
ANIMAL	88 <i>dog</i>			COMPETITION	11 <i>win</i>
				WEATHER	0 —

33

WordNet 3.1

- Where it is:
 - <http://wordnetweb.princeton.edu/perl/webwn>
- Libraries
 - Python: WordNet from NLTK
 - <http://www.nltk.org/search.html?q=wordnet>
 - Java:
 - JWNL, extJWNL on sourceforge

Other (domain specific) thesauri

MeSH: Medical Subject Headings thesaurus from the National Library of Medicine

- **MeSH (Medical Subject Headings)**
 - 177,000 entry terms that correspond to 26,142 biomedical “headings”

- **Hemoglobins**

Synset

Entry Terms: Eryhem, Ferrous Hemoglobin, Hemoglobin

Definition: The oxygen-carrying proteins of ERYTHROCYTES. They are found in all vertebrates and some invertebrates. The number of globin subunits in the hemoglobin quaternary structure differs between species. Structures range from monomeric to a variety of multimeric arrangements

The MeSH Hierarchy

<ul style="list-style-type: none"> 1. + Anatomy [A] 2. + Organisms [B] 3. + Diseases [C] 4. - Chemicals and Drugs [D] <ul style="list-style-type: none"> o Inorganic Chemicals [D01] + o Organic Chemicals [D02] + o Heterocyclic Compounds [D03] + o Polycyclic Compounds [D04] + o Macromolecular Substances [D05] + o Hormones, Hormone Substitutes, and o Enzymes and Coenzymes [D08] + o Carbohydrates [D09] + o Lipids [D10] + o Amino Acids, Peptides, and Proteins o Nucleic Acids, Nucleotides, and Nucl o Complex Mixtures [D20] + o Biological Factors [D23] + o Biomedical and Dental Materials [D25] + o Pharmaceutical Preparations [D26] + o Chemical Actions and Uses [D27] + 5. + Analytical, Diagnostic and Therapeutic Techniques and Equipment [E] 6. + Psychiatry and Psychology [F] 7. + Phenomena and Processes [G] 	<p style="text-align: center;">Amino Acids, Peptides, and Proteins [D12]</p> <p style="text-align: center;">Proteins [D12.776]</p> <p style="text-align: center;">Blood Proteins [D12.776.124]</p> <p style="text-align: center;">Acute-Phase Proteins [D12.776.124.050] +</p> <p style="text-align: center;">Anion Exchange Protein 1, Erythrocyte [D12.776.124.078]</p> <p style="text-align: center;">Ankyrins [D12.776.124.080]</p> <p style="text-align: center;">beta 2-Glycoprotein I [D12.776.124.117]</p> <p style="text-align: center;">Blood Coagulation Factors [D12.776.124.125] +</p> <p style="text-align: center;">Cholesterol Ester Transfer Proteins [D12.776.124.197]</p> <p style="text-align: center;">Fibrin [D12.776.124.270] +</p> <p style="text-align: center;">Glycophorin [D12.776.124.300]</p> <p style="text-align: center;">Hemocyanin [D12.776.124.337]</p> <p style="text-align: center;">▶ Hemoglobins [D12.776.124.400]</p> <p style="text-align: center;">Carboxyhemoglobin [D12.776.124.400.141]</p> <p style="text-align: center;">Erythrocytorins [D12.776.124.400.220]</p>
--	--

37

Uses of the MeSH Ontology

- Provide synonyms (“entry terms”)
 - E.g., glucose and dextrose
- Provide hypernyms (from the hierarchy)
 - E.g., glucose ISA monosaccharide
- Indexing in MEDLINE/PubMED database
 - NLM’s bibliographic database:
 - 20 million journal articles
 - Each article hand-assigned 10-20 MeSH terms

III. Word Similarity

- **Synonymy**: a binary relation
 - Two words are either synonymous or not
- **Similarity** (or **distance**): a looser metric
 - Two words are more similar if they share more features of meaning
- Similarity is properly a relation between **senses**
 - The word “bank” is not similar to the word “slope”
 - Bank¹ is similar to fund³
 - Bank² is similar to slope⁵
- But we’ll compute similarity over both words and senses

Why word similarity

- A practical component in lots of NLP tasks
 - Question answering
 - Natural language generation
 - Automatic essay grading
 - Plagiarism detection
- A theoretical component in many linguistic and cognitive tasks
 - Historical semantics
 - Models of human word learning
 - Morphology and grammar induction

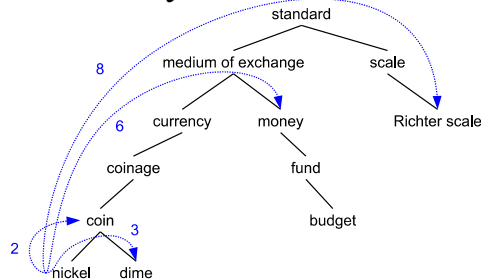
Word similarity and word relatedness

- We often distinguish **word similarity** from **word relatedness**
 - **Similar words**: near-synonyms
 - **Related words**: can be related any way
 - car, bicycle: **similar**
 - car, gasoline: **related**, not similar

Two classes of similarity algorithms

- Thesaurus-based algorithms
 - Are words “nearby” in hypernym hierarchy?
 - Do words have similar glosses (definitions)?
- Distributional algorithms
 - Do words have similar distributional contexts?
 - Distributional (Vector) semantics (prior classes)

Path based similarity



- Two concepts (senses/synsets) are similar if they are near each other in the thesaurus hierarchy
 - =have a short path between them
 - concepts have path 1 to themselves

Refinements to path-based similarity

- $\text{pathlen}(c_1, c_2) = 1 + \text{number of edges in the shortest path in the hypernym graph between sense nodes } c_1 \text{ and } c_2$

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

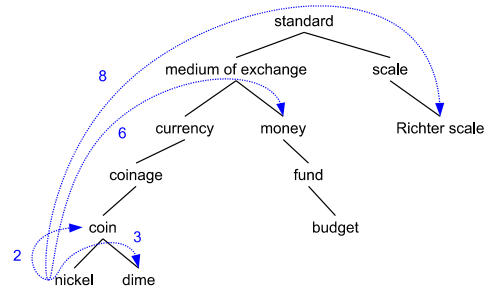
- $\text{wordsim}(w_1, w_2) = \max_{c_1 \in \text{senses}(w_1), c_2 \in \text{senses}(w_2)} \text{sim}(c_1, c_2)$

$$c_1 \in \text{senses}(w_1), c_2 \in \text{senses}(w_2)$$

Example: path-based similarity

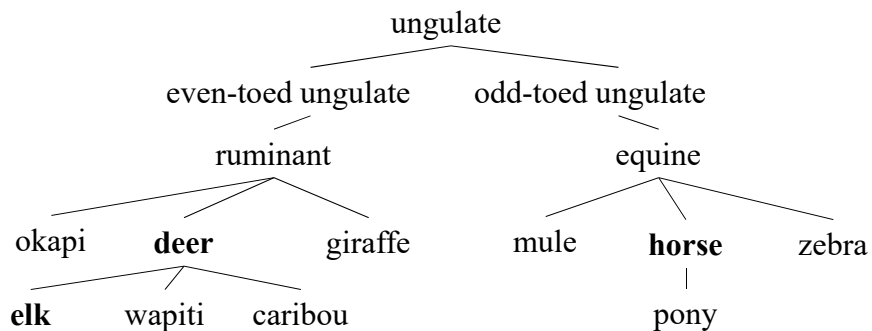
$$\text{simpath}(c_1, c_2) = 1/\text{pathlen}(c_1, c_2)$$

- $\text{simpath}(\text{nickel}, \text{coin}) = 1/2 = .5$
- $\text{simpath}(\text{fund}, \text{budget}) = 1/2 = .5$
- $\text{simpath}(\text{nickel}, \text{currency}) = 1/4 = .25$
- $\text{simpath}(\text{nickel}, \text{money}) = 1/6 = .17$
- $\text{simpath}(\text{coinage}, \text{Richter scale}) = 1/6 = .17$



Which pair of words exhibits the greatest similarity?

1. Deer-elk
2. Deer-horse



Problem with basic path-based similarity

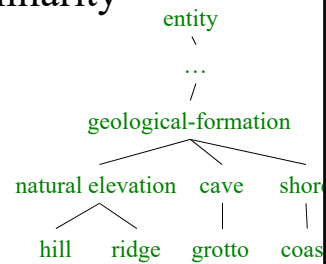
- Assumes each link represents a uniform distance
 - But *nickel to money* seems to us to be closer than *nickel to standard*
 - Nodes high in the hierarchy are very abstract
- We instead want a metric that
 - Represents the cost of each edge independently
 - Words connected only through abstract nodes
 - are less similar

Information content similarity metrics

Resnik 1995

- Let's define $P(c)$ as:
 - The probability that a randomly selected word in a corpus is an instance of concept c
 - Formally: there is a distinct random variable, ranging over words, associated with each concept in the hierarchy
 - for a given concept, each observed noun is either
 - a member of that concept with probability $P(c)$
 - not a member of that concept with probability $1-P(c)$
 - All words are members of the root node (Entity)
 - $P(\text{root})=1$
 - The lower a node in hierarchy, the lower its probability

Information content similarity



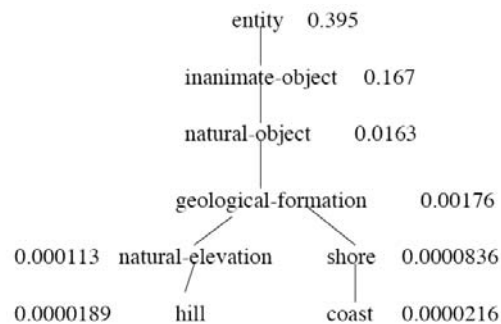
- Train by counting in a corpus
 - Each instance of `hill` counts toward frequency of *natural elevation*, *geological formation*, *entity*, etc
 - Let $\text{words}(c)$ be the set of all words that are children of node c
 - $\text{words}(\text{"geo-formation"}) = \{\text{hill, ridge, grotto, coast, cave, shore, natural elevation}\}$
 - $\text{words}(\text{"natural elevation"}) = \{\text{hill, ridge}\}$

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

Information content similarity

- WordNet hierarchy augmented with probabilities

D. Lin. 1998. An Information-Theoretic Definition of Similarity. ICML 1998



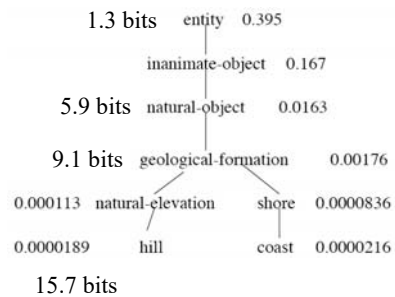
Information content and probability

- The **self-information** of an event, also called its **surprisal**:
 - how surprised we are to know it; how much we learn by knowing it.
 - The more surprising something is, the more it tells us when it happens
 - We'll measure self-information in **bits**.
$$I(w) = -\log_2 P(w)$$
- I flip a coin; $P(\text{heads}) = 0.5$
- How many bits of information do I learn by flipping it?
 - $I(\text{heads}) = -\log_2(0.5) = -\log_2(1/2) = \log_2(2) = 1$ bit
- I flip a biased coin: $P(\text{heads}) = 0.8$ I don't learn as much
 - $I(\text{heads}) = -\log_2(0.8) = -\log_2(0.8) = .32$ bits

51

Information content: definitions

- **Information content:**
 $IC(c) = -\log P(c)$
- **Most informative subsumer**
 (Lowest common subsumer)
 $LCS(c_1, c_2) =$
 The most informative (lowest)
 node in the hierarchy subsuming
 both c_1 and c_2



Using information content for similarity: the Resnik method

Philip Resnik. 1995. Using Information Content to Evaluate Semantic Similarity in a Taxonomy. IJCAI 1995.

Philip Resnik. 1999. Semantic Similarity in a Taxonomy: An Information-Based Measure and its Application to Problems of Ambiguity in Natural Language. JAIR, 11, 95-130.

- The similarity between two words is related to their common information
- The more two words have in common, the more similar they are
- Resnik: measure common information as:
 - The information content of the most informative (lowest) subsumer (MIS/LCS) of the two nodes
 - $\text{sim}_{\text{resnik}}(c_1, c_2) = -\log P(\text{LCS}(c_1, c_2))$

Dekang Lin method

Dekang Lin. 1998. An Information-Theoretic Definition of Similarity. ICM

- Intuition: Similarity between A and B is not just what they have in common
- The more **differences** between A and B, the less similar they are:
 - Commonality: the more A and B have in common, the more similar they are
 - Difference: the more differences between A and B, the less similar
- Commonality: $\text{IC}(\text{common}(A, B))$
- Difference: $\text{IC}(\text{description}(A, B)) - \text{IC}(\text{common}(A, B))$

Dekang Lin similarity theorem

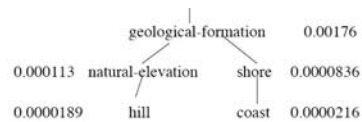
- The similarity between A and B is measured by the ratio between the amount of information needed to state the commonality of A and B and the information needed to fully describe what A and B are

$$sim_{Lin}(A, B) \propto \frac{IC(common(A, B))}{IC(description(A, B))}$$

- Lin (altering Resnik) defines $IC(common(A, B))$ as 2 x information of the LCS

$$sim_{Lin}(c_1, c_2) = \frac{2 \log P(LCS(c_1, c_2))}{\log P(c_1) + \log P(c_2)}$$

Lin similarity



$$sim_{Lin}(A, B) = \frac{2 \log P(LCS(c_1, c_2))}{\log P(c_1) + \log P(c_2)}$$

$$\begin{aligned}
 sim_{Lin}(hill, coast) &= \frac{2 \log P(\text{geological-formation})}{\log P(hill) + \log P(coast)} \\
 &= \frac{2 \ln 0.00176}{\ln 0.0000189 + \ln 0.0000216} \\
 &= .59
 \end{aligned}$$

The (extended) Lesk Algorithm

- A thesaurus-based measure that looks at **glosses**
- Two concepts are similar if their glosses contain similar words
 - *Drawing paper*: **paper** that is **pecially prepared** for use in drafting
 - *Decal*: the art of transferring designs from **pecially prepared paper** to a wood or glass or metal surface
- For each n -word phrase that's in both glosses
 - Add a score of n^2
 - **Paper** and **pecially prepared** for $1 + 2^2 = 5$
 - Compute overlap also for other relations
 - glosses of hypernyms and hyponyms

Summary: thesaurus-based similarity

$$\text{sim}_{\text{path}}(c_1, c_2) = \frac{1}{\text{pathlen}(c_1, c_2)}$$

$$\text{sim}_{\text{resnik}}(c_1, c_2) = -\log P(\text{LCS}(c_1, c_2)) \quad \text{sim}_{\text{lin}}(c_1, c_2) = \frac{2 \log P(\text{LCS}(c_1, c_2))}{\log P(c_1) + \log P(c_2)}$$

$$\text{sim}_{\text{jiangconrath}}(c_1, c_2) = \frac{1}{\log P(c_1) + \log P(c_2) - 2 \log P(\text{LCS}(c_1, c_2))}$$

$$\text{sim}_{e\text{Lesk}}(c_1, c_2) = \sum_{r, q \in \text{RELS}} \text{overlap}(\text{gloss}(r(c_1)), \text{gloss}(q(c_2)))$$

Libraries for computing thesaurus-based similarity

- NLTK
 - [http://nltk.github.com/api/nltk.corpus.reader.html?highlight=similarity - nltk.corpus.reader.WordNetCorpusReader.res_similarity](http://nltk.github.com/api/nltk.corpus.reader.html?highlight=similarity-nltk.corpus.reader.WordNetCorpusReader.res_similarity)
- WordNet::Similarity
 - <http://wn-similarity.sourceforge.net/>

59

Evaluating similarity

- Extrinsic (task-based, end-to-end) Evaluation:
 - Question Answering
 - Essay grading
- Intrinsic Evaluation:
 - Correlation between algorithm and human word similarity ratings
 - Wordsim353: 353 noun pairs rated 0-10.
 $sim(plane, car) = 5.77$
 - Taking TOEFL multiple-choice vocabulary tests
 - Levied is closest in meaning to:
imposed, believed, requested, correlated

Problems with thesaurus-based meaning

- Not every language has a thesaurus
- Even if we have a thesaurus, *recall* problems
 - Many words are missing
 - Most phrases are missing
 - Some connections between senses are missing
 - Adjectives and verbs have less structured hyponymy

Hyponymy and Other Relations

- Could we discover new relationships and add them to a taxonomy?
- Why – unknown word problem (at one time Microsoft or IBM, but not Google)

Hearst Approach

- Based on hand-built patterns
- E.g. *NP-0 such as NP-1* implies *hyponym (NP-1, NP-0)*
- Corpus-based pattern extraction (Snow, Jurafsky, Ng 2005)

Summary

- Lexical Semantics
 - Homonymy, Polysemy, Synonymy, etc.
- Computational resource for lexical semantics
 - WordNet, etc.
- Word Similarity
 - Thesaurus methods