



Word Meaning and Similarity

Word Similarity:
Thesaurus Methods



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

Dan Jurafsky



Why word similarity

- Information retrieval
- Question answering
- Machine translation
- Natural language generation
- Language modeling
- Automatic essay grading
- Plagiarism detection
- Document clustering



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

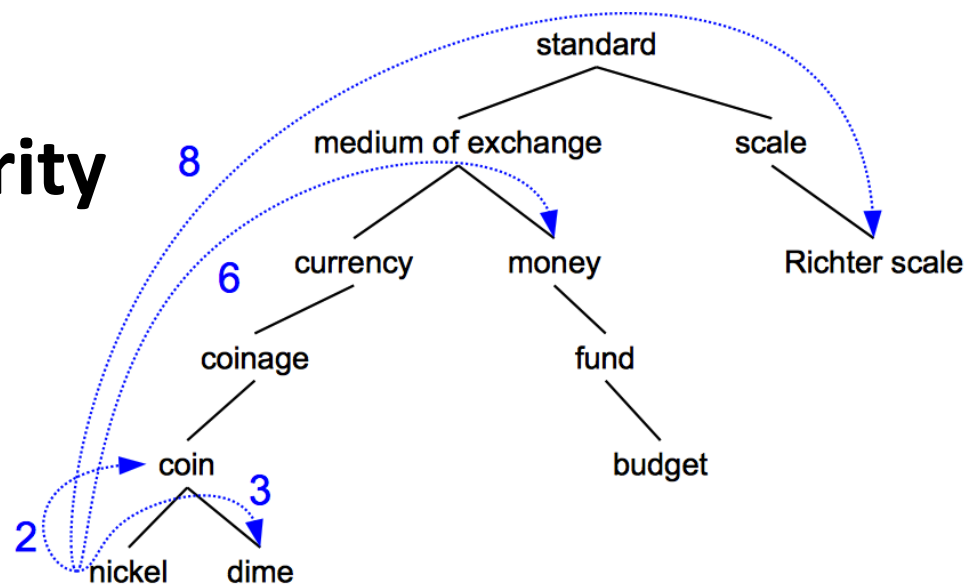


Two classes of similarity algorithms

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



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

- $\text{pathlen}(c_1, c_2) = 1 + \text{number of edges in the shortest path in the hypernym graph between sense nodes } c_1 \text{ and } c_2$
- ranges from 0 to 1 (identity)
- $\text{simpath}(c_1, c_2) = \frac{1}{\text{pathlen}(c_1, c_2)}$
- $\text{wordsim}(w_1, w_2) = \max_{c_1 \in \text{senses}(w_1), c_2 \in \text{senses}(w_2)} \text{sim}(c_1, c_2)$



Example: path-based similarity

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

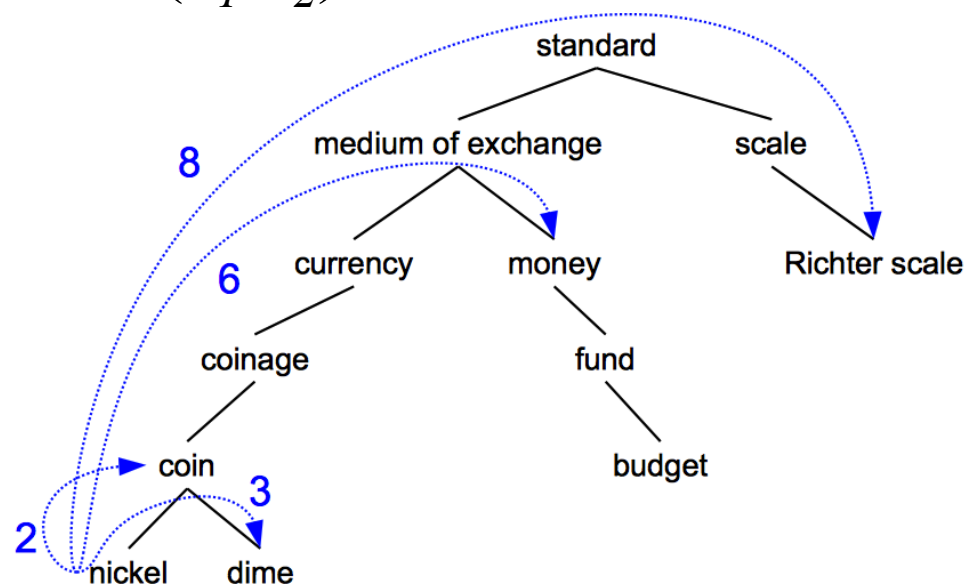
$$\text{simpath}(\text{nickel}, \text{coin}) = 1/2 = .5$$

$$\text{simpath}(\text{fund}, \text{budget}) = 1/2 = .5$$

$$\text{simpath}(\text{nickel}, \text{currency}) = 1/4 = .25$$

$$\text{simpath}(\text{nickel}, \text{money}) = 1/6 = .17$$

$$\text{simpath}(\text{coinage}, \text{Richter scale}) = 1/6 = .17$$





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. Using information content to evaluate semantic similarity in a taxonomy. IJCAI

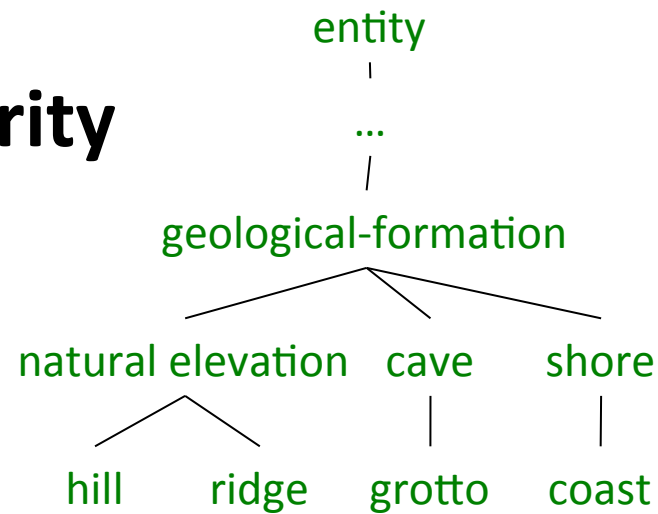
- Let's define $P(c)$ as:
 - The probability that a randomly selected word in a corpus is an instance of concept c
 - Formally: there is a distinct random variable, ranging over words, associated with each concept in the hierarchy
 - for a given concept, each observed noun is either
 - a member of that 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
 - $\text{words}(\text{"geo-formation"}) = \{\text{hill}, \text{ridge}, \text{grotto}, \text{coast}, \text{cave}, \text{shore}, \text{natural elevation}\}$
 - $\text{words}(\text{"natural elevation"}) = \{\text{hill}, \text{ridge}\}$



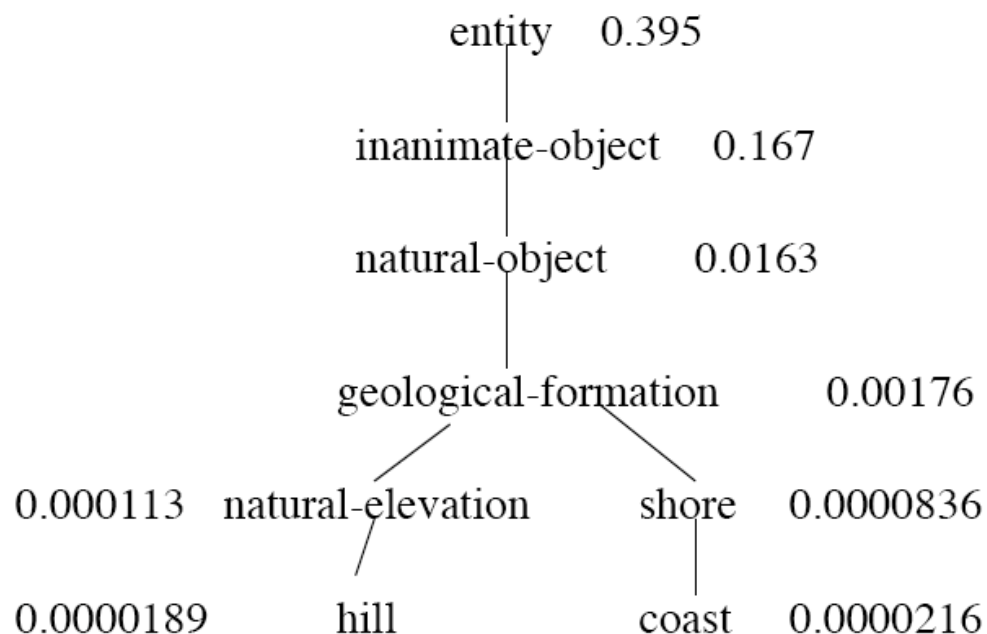
$$P(c) = \frac{\sum_{w \in \text{words}(c)} \text{count}(w)}{N}$$



Information content similarity

- WordNet hierarchy augmented with probabilities $P(c)$

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



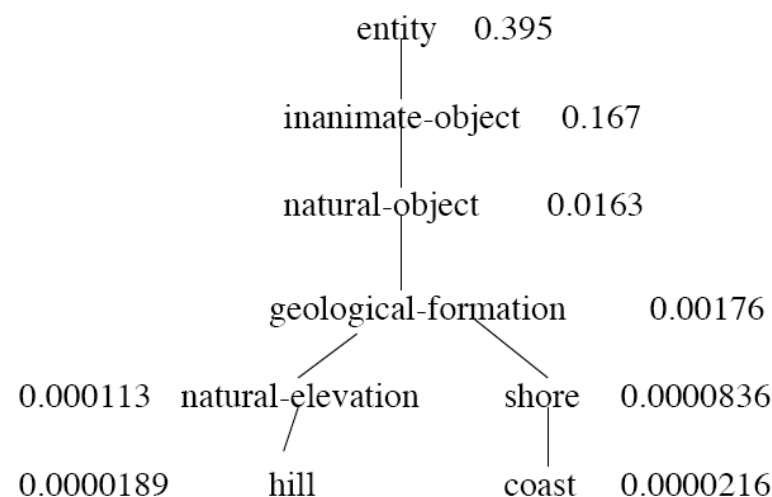


Information content: definitions

- 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



Dan Jurafsky



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. ICML

- 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: $IC(\text{common}(A,B))$
- Difference: $IC(\text{description}(A,B)) - IC(\text{common}(A,B))$



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

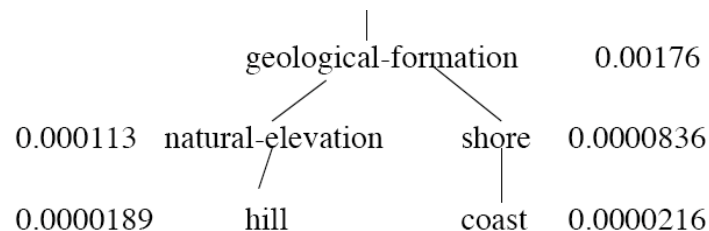
$$sim_{Lin}(A, B) \propto \frac{IC(common(A, B))}{IC(description(A, B))}$$

- 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 function



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

$$sim_{Lin}(\text{hill}, \text{coast}) = \frac{2 \log P(\text{geological-formation})}{\log P(\text{hill}) + \log P(\text{coast})}$$

$$= \frac{2 \ln 0.00176}{\ln 0.0000189 + \ln 0.0000216}$$

$$= .59$$



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 **specially prepared** for use in drafting
 - **Decal**: the art of transferring designs from **specially 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 **specially prepared** for $1 + 2^2 = 5$
 - Compute overlap also for other relations
 - glosses of hypernyms and hyponyms



Summary: thesaurus-based similarity

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

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

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

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



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/>
 - Web-based interface:
 - <http://marimba.d.umn.edu/cgi-bin/similarity/similarity.cgi>



Evaluating similarity

- Intrinsic Evaluation:
 - Correlation between algorithm and human word similarity ratings
- Extrinsic (task-based, end-to-end) Evaluation:
 - Malapropism (spelling error) detection
 - WSD
 - Essay grading
 - Taking TOEFL multiple-choice vocabulary tests

Levied is closest in meaning to:

imposed, believed, requested, correlated

[illegible]