

Word Senses and WordNet

Natural Language Processing
CSCI 6350
Dr. Sal Barbosa

Overview

- Terminology and Definitions
- WordNet
- Word Similarity
- Information Content Measures

2

Terminology: Lemma and Wordform

- A **lemma** or **citation form**
 - Same stem, part of speech, rough semantics
- A **wordform**
 - The inflected word as it appears in text

Wordform	Lemma
banks	bank
sung	sing
duermes	dormir

3

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

Lemmas Have senses

- One lemma “bank” can have many meanings:
 - Sense 1: ...a **bank** can hold the investments in a custodial account 1..
 - ...as agriculture burgeons on the east **bank** the
 - Sense 2: river will shrink even more” 2
- **Sense (or word sense)**
 - A discrete representation of an aspect of a word’s meaning.
- The lemma **bank** here has two senses

4

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

Homonymy

Homonyms: words that share a form but have unrelated, distinct meanings:

- **bank**₁: financial institution, **bank**₂: sloping land
 - **bat**₁: club for hitting a ball, **bat**₂: nocturnal flying mammal
1. Homographs (bank/bank, bat/bat)
 2. Homophones:
 1. **Write** and **right**
 2. **Piece** and **peace**

5

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

Homonymy causes problems for NLP applications

- Information retrieval
 - "bat care"
 - Machine Translation
 - **bat**: **murciélago** (animal) or **bate** (for baseball)
- Text-to-Speech
 - **bass** (stringed instrument) vs. **bass** (fish)

6

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

Polysemy

- 1. The **bank** was constructed in 1875 out of local red brick.
- 2. I withdrew the money from the **bank**
- Are those the same sense?
 - Sense 1: "The building belonging to a financial institution"
 - Sense 2: "A financial institution"
- A **polysemous** word has **related** meanings
 - Most non-rare words have multiple meanings

7

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

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)

8

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

How do we know when a word has more than one sense?

- One technique for determining if two senses are distinct is to join two uses of the word in a sentence (this is called **zeugma**)
- The “zeugma” test: Two senses of **serve**?
 - Which flights **serve** breakfast?
 - Does Lufthansa **serve** Philadelphia?
 - ?Does Lufthansa **serve** breakfast and San Jose?
- Since this conjunction sounds weird,
 - we say that these are **two different senses of “serve”**

9

adapted from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

Synonyms

- 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 substituted for each other in all situations
 - If so they have the same **propositional meaning**

10

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

Synonyms

- But there are few (or no) examples of perfect synonymy.
 - Even if many aspects of meaning are identical
 - Still may not preserve the acceptability based on notions of norms, politeness, slang, register, genre, etc.
- Example:
 - Water/H₂O
 - Big/large
 - Brave/courageous

11

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

Synonymy is a relation between senses rather than words

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

12

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

Antonyms

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

13

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

Hyponymy and Hypernymy

- One sense is a **hyponym/subordinate** of another if the first sense is more specific, denoting a subclass of the other
 - car is a hyponym of vehicle
 - mango is a hyponym of fruit
- Conversely **hypernym/superordinate** ("hyper is super")
 - vehicle is a **hypernym** of car
 - fruit is a hypernym of mango

Superordinate/hypernym	vehicle	fruit	furniture
Subordinate/hyponym	car	mango	chair

14

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

Hyponymy more formally

- Extensional:
 - The class denoted by the superordinate extensionally includes the 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
 - A is **subsumed** by B

15

adapted from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

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, ...

16

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

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*.

17

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

WordNet 3.0

- A hierarchically organized lexical database
- On-line thesaurus + aspects of a dictionary
 - Some [other languages](#) available or under development
 - (Arabic, Finnish, German, Portuguese...)

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

18

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

Senses of “bass” in Wordnet

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)
- *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”

19

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

How is “sense” defined in WordNet?

- The **synset (synonym set)**, the set of near-synonyms, instantiates a sense or concept, with a **gloss** (or definition)
- Example: **chump** as a noun with the **gloss**:
“a person who is gullible and easy to take advantage of”
- This sense of “chump” is shared by 9 words:
chump¹, fool², gull¹, mark⁹, patsy¹, fall guy¹, sucker¹, soft touch¹, mug²
- Each of **these** senses have this same gloss
 - (Not every sense; sense 2 of gull is the aquatic bird)

20

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

WordNet Hypernym Hierarchy for “bass”

- *S. (n) bass, basso* (an adult male singer with the lowest voice)
 - *direct hypernym / inherited hypernym / sister term*
 - *S. (n) singer, vocalist, vocaliser* (a person who sings)
 - *S. (n) musician, instrumentalist, player* (someone who plays a musical instrument (as a profession))
 - *S. (n) performer, performing artist* (an entertainer who performs a dramatic or musical work for an audience)
 - *S. (n) entertainer* (a person who tries to please or amuse)
 - *S. (n) person, individual, someone, somebody, mortal, soul* (a human being) “there was too much for one person to do”
 - *S. (n) organism, being* (a living thing that has (or can develop) the ability to act or function independently)
 - *S. (n) living thing, animate thing* (a living (or once living) entity)
 - *S. (n) whole, unit* (an assemblage of parts that is regarded as a single entity) “how big is that part compared to the whole?”, “the team is a unit”
 - *S. (n) object, physical object* (a tangible and visible entity; an entity that can cast a shadow) “it was full of rackets, balls and other objects”
 - *S. (n) physical entity* (an entity that has physical existence)
 - *S. (n) entity* (that which is perceived or known or inferred to have its own distinct existence (living or nonliving))

24

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

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 |
| Instance Hypernym | Instance | From instances to their concepts | <i>Austen</i> ¹ → <i>author</i> ¹ |
| Instance Hyponym | Has-Instance | From concepts to concept instances | <i>composer</i> ¹ → <i>Bach</i> ¹ |
| Member Meronym | Has-Member | From groups to their members | <i>faculty</i> ² → <i>professor</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> ¹ → <i>meal</i> ¹ |
| Substance Meronym | | From substances to their subparts | <i>water</i> ¹ → <i>oxygen</i> ¹ |
| Substance Holonym | | From parts of substances to wholes | <i>gin</i> ¹ → <i>martini</i> ¹ |
| Antonym | | Semantic opposition between lemmas | <i>leader</i> ¹ ↔ <i>follower</i> ¹ |
| Derivationally Related Form | | Lemmas w/same morphological root | <i>destruction</i> ¹ ↔ <i>destroy</i> ¹ |

22

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

WordNet Verb Relations

| Relation | Definition | Example |
|-----------------------------|--|---|
| Hypernym | From events to superordinate events | <i>fly</i> ⁹ → <i>travel</i> ⁵ |
| Troponym | From events to subordinate event (often via specific manner) | <i>walk</i> ¹ → <i>stroll</i> ¹ |
| Entails | From verbs (events) to the verbs (events) they entail | <i>snore</i> ¹ → <i>sleep</i> ¹ |
| Antonym | Semantic opposition between lemmas | <i>increase</i> ¹ ↔ <i>decrease</i> ¹ |
| Derivationally Related Form | Lemmas with same morphological root | <i>destroy</i> ¹ ↔ <i>destruction</i> ¹ |

23

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

WordNet: Viewed as a Graph



24

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

“Supersenses”¹⁷: The top level hypernyms in the hierarchy

(counts from Schneider and Smith 2013’s Streusel corpus)

| Noun | | Verb | |
|------------|-----------------------|-------------|----------------------|
| GROUP | 1469 <i>place</i> | STATIVE | 2922 <i>is</i> |
| PERSON | 1202 <i>people</i> | COGNITION | 1093 <i>know</i> |
| ARTIFACT | 971 <i>car</i> | COMMUNIC. | 974 <i>recommend</i> |
| COGNITION | 771 <i>way</i> | SOCIAL | 944 <i>use</i> |
| FOOD | 766 <i>food</i> | MOTION | 602 <i>go</i> |
| ACT | 700 <i>service</i> | POSSESSION | 309 <i>pay</i> |
| LOCATION | 638 <i>area</i> | CHANGE | 274 <i>fix</i> |
| TIME | 530 <i>day</i> | EMOTION | 249 <i>love</i> |
| EVENT | 431 <i>experience</i> | PERCEPTION | 143 <i>see</i> |
| COMMUNIC. | 417 <i>review</i> | CONSUMPTION | 93 <i>have</i> |
| POSSESSION | 339 <i>price</i> | BODY | 82 <i>get...done</i> |
| ATTRIBUTE | 205 <i>quality</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 — |

25

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

Supersenses

- A word’s supersense can be a useful coarse-grained representation of word meaning for NLP tasks

I googled_{communication} restaurants_{GROUP} in the area_{LOCATION} and Fuji_{GROUP} Sushi_{GROUP} came_{communication} up and reviews_{COMMUNICATION} were_{stative} great so I made_{communication} a carry_{possession} out order_{communication}

26

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

WordNet 3.0

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

27

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

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

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

Synset

28

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

The MeSH Hierarchy

1. • Anatomy [A]
2. • Organisms [B]
3. • Diseases [C]
4. • Chemicals and Drugs [D]
 - Inorganic Chemicals [D011] +
 - Organic Chemicals [D02] +
 - Heterocyclic Compounds [D03] +
 - Polycyclic Compounds [D04] +
 - Macromolecular Substances [D05] +
 - Hormones, Hormone Substitutes, and Hormone Antagonists [D08] +
 - Enzymes and Coenzymes [D09] +
 - Carbohydrates [D09] +
 - Lipids [D10] +
 - Amino Acids, Peptides, and Proteins [D12] +
 - Nucleic Acids, Nucleotides, and Nucleosides [D13] +
 - Complex Mixtures [D20] +
 - Biological Factors [D23] +
 - Biomedical and Dental Materials [D25] +
 - Pharmaceutical Preparations [D26] +
 - Chemical Actions and Uses [D27] +
5. • Analytical, Diagnostic and Therapeutic Techniques and Equipment [E]
6. • Psychiatry and Psychology [F]
7. • Phenomena and Processes [G]

Amino Acids, Peptides, and Proteins [D12]

Proteins [D12.776]

Blood Proteins [D12.776.124]

Acute-Phase Proteins [D12.776.124.050] +

Anion Exchange Protein 1, Erythrocyte [D12.776.124.078]

Ankyrins [D12.776.124.080]

beta 2-Glycoprotein 1 [D12.776.124.117]

Blood Coagulation Factors [D12.776.124.125] +

Cholesterol Ester Transfer Proteins [D12.776.124.197]

Fibrin [D12.776.124.270] +

Glycophorin [D12.776.124.300]

Hemocyanin [D12.776.124.337]

► Hemoglobins [D12.776.124.400]

Carboxyhemoglobin [D12.776.124.400.141]

Erythrocytorins [D12.776.124.400.220]

29

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

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

30

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

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

31

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

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

32

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

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

33

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

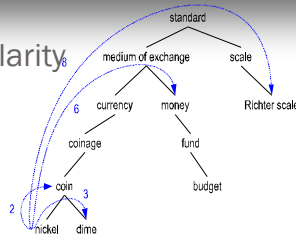
Two classes of similarity algorithms

- 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

34

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

Path based similarity



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

35

adapted from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

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

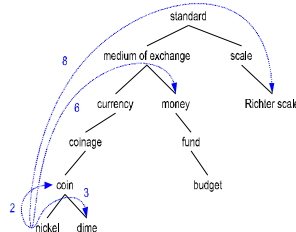
36

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

Example: path-based similarity

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

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



37

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

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

38

adapted from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

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

39

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

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
 - $\text{words}(\text{"geo-formation"}) = \{\text{hill}, \text{ridge}, \text{grotto}, \text{coast}, \text{cave}, \text{shore}, \text{natural elevation}\}$
 - $\text{words}(\text{"natural elevation"}) = \{\text{hill}, \text{ridge}\}$

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

40

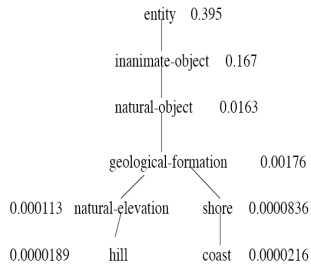
N is the number of words in the corpus that are also in the thesaurus

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

Information content similarity

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

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



44

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

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?
 - $I(\text{heads}) = -\log_2(0.5) = -\log_2(1/2) = \log_2(2) = 1 \text{ bit}$
- I flip a biased coin: $P(\text{heads}) = 0.8$ I don't learn as much
 - $I(\text{heads}) = -\log_2(0.8) = -\log_2(0.8) = .32 \text{ bits}$

42

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

Information content: definitions

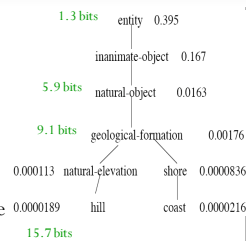
- Information content:

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

- 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



43

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

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
 - $\text{sim}_{\text{resnik}}(c_1, c_2) = -\log P(LCS(c_1, c_2))$

44

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

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

45

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

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

$$\text{sim}_{Lin}(A,B) \propto \frac{IC(\text{common}(A,B))}{IC(\text{description}(A,B))}$$

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

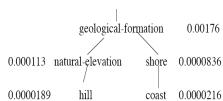
$$\text{sim}_{Lin}(c_1, c_2) = \frac{2 \log P(LCS(c_1, c_2))}{\log P(c_1) + \log P(c_2)}$$

46

from Speech and Language Processing (3rd ed - DRAFT) by Jurafsky and Martin

Lin similarity function

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



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

47

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

48

Summary: Thesaurus-based Similarity

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

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

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

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