Semantics I IN4325 - Information Retrieval



Today

- We finally finish the evaluation part
- Knowledge-based semantics
 - applications
- Statistics-based "semantics"
 - applications



Evaluation: state-of-the-art



Turning away from binary qrels

1960s









Järvelin & Kekäläinen, 2002 [5]

- Graded relevance scales (e.g. 0-3)
- Used in Web search evaluations
- NDCG measures the "gain" of documents
- Assumptions
 - highly relevant documents are more valuable than marginally relevant documents
 - The greater the ranked position of a relevant document, the less valuable it is for the user (less likely that the user will examine it)
 - User has limited time, may have seen the information already, cognitive load, etc.



Normalized Discounted Cumulative Gain Järvelin & Kekäläinen, 2002 [5]

Direct cumulative gain can be defined iteratively

$$G' = <3,2,3,0,0,1,2,2,3,0,\ldots>$$

$$CG_{i} = \begin{cases} G_{1}, if \ i = 1 \\ CG_{i-1} + G_{i}, otherwise \end{cases}$$

$$CG' = <3,5,8,8,8,9,11,13,16,16,\ldots>$$

CG at each rank can be read directly



Järvelin & Kekäläinen, 2002 [5]

- **Discounted cumulative gain**: reduce the document score as ist rank increases (but not too steeply)
 - Divide the document score by the log of ist rank
 - Base of the logarithm determines discount factor

$$DCG_{i} = \begin{cases} CG_{i}, \text{ if } i < b \\ CG_{i-1} + G_{i} / \log_{b} i, \text{ if } i \ge b \end{cases}$$

 $\log_2 1 = 0$ $\log_2 2 = 1$ $\log_2 1024 = 10$



Järvelin & Kekäläinen, 2002 [5]

- Normalized DCG: compare DCG to the theoretically best possible
- Ideal vectors constructed as follows:

$$BV_{i} = \begin{cases} 3, if \ i \leq m, \\ 2, if \ m < i \leq m+l, \\ 1, if \ m+l < i \leq m+l+k, \\ 0, otherwise \\ I' = <3,3,3,2,2,2,1,1,1,1,0,0,0,...> \\ CG_{I'} = <3,6,9,11,13,15,16,17,18,19,19,19,19,19,...> \\ DCG_{I'} = <3,6,7.89,8.89,9.75,10.52,10.88,11.21,11.53,11.83,...> \end{cases}$$



Järvelin & Kekäläinen, 2002 [5]

- Normalized DCG: compare DCG to the theoretically best possible
 - The DCG vectors are divided component-wise by the corresponding ideal DCG vectors
 - 1: ideal performance
 - [0,1): share of ideal performance
- In practice NDCG for query *Q* at rank *k*:

rel. score assessors gave to D at query j

$$NDCG(Q,k) = \frac{1}{|Q|} \sum_{j=1}^{|Q|} Z_{kj} \sum_{m=1}^{k} \frac{2^{R(j,m)} - 1}{\log_2(1+m)}$$

normalization factor so that a perfect ranking's NDCG at *k* for query *j* is 1





TUDelft (trec_eval scores)

Latent spaces & concepts





What is wrong with keyword-based IR?

Relies on syntax alone (string matching)

 Shortcomings: multi-word expressions, synonymy, polysemy, generalizations (vocabulary problem)

<DOC>Space station makes progress with computers</DOC>
Query: "ISS"

<DOC>Apple PCs gaining gaming credibility</DOC>
<DOC>Canada imposes apple moth restrictons</DOC>
Query: "apple"

<DOC>Germany's Merkel meets treaty opponent Poland</DOC>
Query: "German chancellor"



What is wrong with keyword-based IR?



<DOC>Germany's Merkel meets treaty opponent Poland</DOC>
Query: "German chancellor"



The advantages of concepts/semantics

- If the machine "knows" the meaning of terms and phrases, the retrieval effectiveness improves
 - The **vocabulary gap** between users and text authors vanishes

Our world is being revolutionized by data-driven methods: access to large amounts of data has generated new insights and opened exciting new opportunities in commerce, science, and computing applications. Processing the enormous quantities of data necessary for these advances requires large clusters, making distributed computing paradigms more crucial than ever. **MapReduce** is a programming model for expressing distributed computations on massive datasets and an execution framework for large-scale data processing on clusters of commodity servers. The programming model provides an easy-to-understand abstraction for designing scalable algorithms, while the execution framework transparently handles many system-level details, ranging from scheduling to synchronization to fault tolerance. This book focuses on **MapReduce** algorithm design, with an emphasis on text processing algorithms common in natural language processing, information retrieval, and machine learning. We introduce the notion of **MapReduce** design patterns, which represent general reusable solutions to commonly occurring problems across a variety of problem domains. This book not only intends to help the reader "think in **MapReduce**", but also discusses limitations of the programming model as well.

query: hadoop

Goal: retrieval based on the semantics (the topics) of documents and queries



Concept extraction from documents

 Source of extra-document knowledge: knowledge bases vs. statistics

- Knowledge bases
 - Examples: WordNet, UMLS
 - Advantages: precise, ~faultless
 - Disadvantages: maintainance, up-to-datedness, corpusdependency
- Statistics (LSA, PLSA, LDA, HAL)
 - Derived from the corpus itself or external corpora (the Web)
 - Advantages: automatic generation, corpus-dependency
 - Disadvantages: noisy, error-prone











WordNet

http://wordnet.princeton.edu/

- Lexical database of the English language
 - Models common-sense world knowledge
- Inspired by psycholinguistic theories of human lexical memory
- Manual effort led by George A. Miller (Princeton University)
- Nouns, adjectives, adverbs and verbs are organized by lexical concepts (sets of synonyms -> synsets)
- The syntactic categories form largely separate networks
- Lexical concepts linked through relations



WordNet 3.0 statistics

Synset distribution



3644 adverb synsets

- 18156 adjective synsets
 - 13767 verb synsets
- 82115 noun synsets

Most forms have few senses, a few have many senses (Zipf distribution)



Example: instrument

WordNet's sense entries consist of a set of synonyms, dictionary style definitions (gloss) and example uses

Noun

- S: (n) instrument (a device that requires skill for proper use)
- <u>S:</u> (n) instrument, tool (the means whereby some act is accomplished) "my" greed was the instrument of my destruction"; "science has given us new tools to fight disease"
- <u>S:</u> (n) instrument, <u>pawn</u>, <u>cat's-paw</u> (a person used by another to gain an end)
- <u>S:</u> (n) <u>legal document</u>, <u>legal instrument</u>, <u>official document</u>, **instrument** ((law) a document that states some contractual relationship or grants some right)
- <u>S:</u> (n) <u>instrumental role</u>, **instrument** (the semantic role of the entity (usually inanimate) that the agent uses to perform an action or start a process)
- <u>S:</u> (n) <u>musical instrument</u>, **instrument** (any of various devices or contrivances that can be used to produce musical tones or sour

Verb

- <u>S:</u> (v) instrument (equip with instruments for mea controlling)
- S: (v) instrument, instrumentate (write an instrum
- <u>S:</u> (v) instrument (address a legal document to)

- A word sense is made up of synonyms (words with similar/identical meaning)
- Each word sense is mapped to a synset -> a synset is a set of words that are interchangeable in some context



Important WordNet relations

- Noun relations
 - Hyponymy / hypernymy ("is-a")
 - red is a hyponym of color; color is a hypernym of red
 - Holonymy / meronymy ("part-of")
 - a toe is a meronym of leg, leg is a holonym of toe
 - Antonymy
 - king is an antonym of queen, queen is an antonym of king
- Verb relations
 - Hypernymy, antonymy
 - Entailment
 - to snore lexically entails to sleep (you have to sleep to be able to snore)
 - Troponymy
 - to limp is a troponym of to walk (a more specific description of walking)
- Adverb/adjective relations
 - Antonymy



WordNet as a graph



- In-degree d_{in} of a node: number of incoming edges $d_{in}(v_1)=0$, $d_{in}(v_2)=1$
- Out-degree d_{out} of a node: number of outgoing edges $d_{out}(v_1)=2$, $d_{out}(v_2)=0$
- Path *P*: sequence of nodes with an edge between neighbouring nodes $P = \{v_1, v_3, v_4\}$
- Distance $d(v_i, v_j)$ between nodes v_i and v_j : length of the shortest path between them $d(v_1, v_4) = 2$, $d(v_1, v_2) = 1$, $d(v_1, v_1) = 0$, $d(v_1, v_5) = inf$.

Relatedness measures on WordNet

- Calculate the distance/relatedness between words or synsets within the same syntactic category
- Proposed measures rely on different resources and apply to different categories
- Semantic *relatedness*: relationship between concepts is determiend by any kind of relation
- Semantic *similarity*: considers the IS-A relation only
- Two types
 - Based on senses
 - Based on words (usually involves all senses and min/max/av. scores)



Leacock-Chodorow similarity [5]

- Measurement *based on the IS-A shortest path* $len(c_1, c_2)$ *between two synsets* c_1, c_2
- Path length scaled by the overall depth D (~18 for nouns) of the taxonomy

$$sim(c_1, c_2) = -\log \frac{len(c_1, c_2)}{2D}$$

- Examples (maximum score across senses)
 - sim(car#n,vehicle#n)=2.59
 - sim(car#n,boat#n)=2.08
 - sim(car#n,birthday_cake#n)=1.05



Hirst-St-Onge relatedness [5]

- Two concepts c_1 and c_2 are semantically close, if their synsets are connected by a path that is not too long and does not change direction too often
- Approach utilizes all relations within WordNet



Resnik similarity [5]

sim(car#n,vehicle#n)=5.53
sim(car#n,boath#n)=5.53
sim(car#n,birthday_cake#n)=0.61

- Information content *IC* is a measure of specificity of a concept
 - High IC indicates specificity
 - Low IC indicates generality
- *IC* estimated by counting concept frequencies in a corpus
 - Sense-tagged corpus: use counts directly
 - Untagged corpus: assign counts to all synsets that contain the word
 - Counts trickle up the IS-A hierarchy (giant counts towards grownup)
- Similarity between two concepts is the *IC* of their lowest common superordinate

$$IC(c_i) = -\log(P_{ml}(c_i))$$
 $sim(c_1, c_2) = IC(lso(c_1, c_2))$



Resnik similarity [5]

sim(car#n,vehicle#n)=5.53
sim(car#n,boath#n)=5.53



• Similarity between two concepts is the *IC* of their lowest common superordinate

$$IC(c_i) = -\log(P_{ml}(c_i))$$
 $sim(c_1, c_2) = IC(lso(c_1, c_2))$



Which measure performs best? [5]

• Set of ground truth judgments

- Rubenstein & Goodenough (1965) created 65 pairs of words; 51 human assessors rated their similarity of meaning (scale 0.0-4.0)
- Evaluation in terms of the correlation between measure based similarity and user based similarity judgments

Leacock-Chodorow *	0.84
Hirst-St-Onge	0.79
Resnik	0.78



Word-based distances

 Graph with words as nodes; edges between those words that appear in one synset





Word-based distances

Kamps et al, 2004



Source: http://staff.science.uva.nl/~kamps/wordnet/



WordNet applications

- Word sense disambiguation
- Refinement of search queries
- Semantic orientation of texts
- Detection and correction of spelling errors
- Summarization

Scientific papers @ Google Scholar mentioning 'WordNet'



Search query refinement

Nemrava, 2006

 P_0 such as P_1 , P_2 , ..., P_{n-1} (and or) P_n

- Motivation: keyword searches are often ambiguous
 - "Java", "Python", "Pluto"
- Idea: give the user a choice between possible hypernyms
 *if you do not have a query log
 - 1. Look up WordNet synsets and glosses for the query
 - 2. Process the glosses (POS tagger, stopword removal) and keep the nouns as potential hypernyms
 - 3. Apply Hearst patterns on the Web & retrieve #result pages for each candidate
 - 4. Consider the candidate with the highest score as hypernym



Search query refinement

1. WordNet glosses for query term "Pluto"

- **SYN1** a small planet and the farthes known planet from the sun; has the most elliptical orbit of all the planets
- SYN2 (Greek mythology) the god of the underworld in ancient mythology; brother of Zeus and husband of Persephone
- SYN3 a cartoon character created by Walt Disney

2. Candidate nouns

- **SYN1** planet, sun, orbit, planets
- SYN2 Greek, god, underworld, mythology, brother, Zeus, husband, Persephone
- SYN3 cartoon, character, Walt, Disney
- 3. Hearst patterns and page counts (shown for SYN1 only)
 - "Pluto is a planet" (1550), "Pluto is planet" (145) "Pluto is a sun" (2), "Pluto is sun" (0)
 - "Pluto is an orbit" (1), "Pluto is orbit" (1) "Pluto is a planets" (0), "Pluto is planets" (0)
- 4. Refinement offers: "Pluto planet", "Pluto god", "Pluto cartoon"



Detection & correction of spelling errors

 Intuition: in coherent texts, many instances of related pairs of words can be found

 Words are disambiguated where possible by accepting senses that are semantically related to possible senses of other nearby words

I walk along the river \rightarrow I talk along the river His dog barks \rightarrow His dog parks The sun shines \rightarrow The son shines



Detection & correction of spelling errors Error detection approach: Hirst & Budanitsky, 2005 [7]

- 1. Consider word *w* that appears only once in the text and is semantically unrelated to nearby words
- 2. Create spelling variations w_1 , w_2 , w_3 , ... of w by a single insertion, deletion or transposition



- 3. Check w_1 , w_2 , w_3 , ... for their semantic relatedness to nearby words
- 4. Suggest the most closely related found variation to the user



Drawbacks of WordNet

Coverage and named entities (extra lecture about NEs)

- George W. Bush, Barack Obama, Queen Beatrix
- York ("House of York")
- The structure is somewhat arbitrary
 - {water, body of water}
 - {thing}
 - {physical entity}
- Synsets can be difficult to distinguish, e.g. fear#v
 - **SYN1** be afraid or feel anxious or apprehensive about a possible or probable situation or event
 - **SYN2** be afraid or scared of; be frightened of
 - **SYN3** be uneasy or apprehensive about



Nguyen et al., 2009

- First of all: Wikipedia != WordNet (EuroWordNet)
- But: when relying on Wikipedia concepts and their hyperlinks, the graph structures appear similar
 - Concept: Wikipedia article
 - Edge between concepts: hyperlinks between articles
- Given a query in language A, it is translated to language B using **only** Wikipedia
- Idea: use cross-lingual links for translation

Languages العربية Aragonés Azərbaycanca Bân-lâm-gú Беларуская (тарашкевіца) Български Bosanski Català Česky Dansk Deutsch Eesti Español فارسى Français Galego 한국어 ldo different languages Bahasa Indonesia Íslenska



cross-lingual links:

same concept in

Nguyen et al., 2009

- Why Wikipedia
 - Better coverage of named entities, many domain-specific concepts
 - Community keeps it up-to-date
 - Wikipedia articles provide more context (than glosses)
 - Synonyms are provided "for free" (redirect pages)
 - Disadvantage: coverage of common terms (e.g. "house" has many unusual senses)

• Approach:

$$query_{English} \longrightarrow WP_{English} \longrightarrow WP_{German} \longrightarrow query_{German}$$



Nguyen et al., 2009

- 1. Determine the most relevant concepts to the original query after a keyword-based search with the whole query
- 2. Search on every single query term using
 - The **internal links** from the concepts retrieved in 1.
 - The **text** and title of the articles retrieveed in 1.
- 3. Add articles that redirect to the found concepts (synonyms, spelling variants)
- 4. Create new translated query

$$query_{English} \longrightarrow WP_{English} \longrightarrow WP_{German} \longrightarrow query_{German}$$



Nguyen et al., 2009

CLEF 2004 <title>Atlantis-Mir Koppeling</title> <desc>Vind documenten over de eerste space shuttle aankoppeling tussen de Amerikaanse shuttle Atlantis en het Mir ruimte station</desc> search WP with full query WP:space shuttle_atlantis, WP:mir_(ruimtestation) search WP with query terms on links and content amerika, atlantis (disambiguation), koppeling, mir mir (disambiguation), roskosmos, atlantis (ruimteveer), ...



Nguyen et al., 2009





Nguyen et al., 2009

- Evaluation: retrieval of English documents using Dutch, French and Spanish queries
 - Optimum: same effectiveness as the monolingual run

Language	МАР
English (monolingual)	0.3407
French	0.2278 (66.9%)
Spanish	0.2181 (64.0%)
Dutch	0.2038 (59.8%)



Nguyen et al., 2009

 Evaluation: retrieval of English documents using Dutch, French and Spanish queries

When does it fail?

"fictives" (fr.) means "fictional" (engl.)

- Planets_in_science_fiction
- Fictional_brands





Alternatives to WordNet or Wikipedia

• Yago (Suchanek et al., 2007)

- Combines WordNet and Wikipedia
- Contains 2 million entities, "knows" 20 million facts
- Manually confirmed accuracy of 95%
- Queries have the form <entity><relation><entity>
- WordNet Affect
 - Manual labeling of 3000 synsets with affective concepts
 - 11 affective labels, e.g.
 - Emotion (anger#n#1, fear#v#1)
 - Cognitive state (confusion#n#2, dazed#adj#2)
 - Emotional response (cold_sweat#n#1, tremble#v#1)



IM A GENET

ImageNet

Deng et al., 2009

- Image retrieval (visual codewords)
- Text-to-image synthesis



- cross-modal semantic relationships
- decrease semantic gap (MM lecture)

http://www.image-net.org/



Latent Semantic Analysis

Deerwester & Dumais, 1990

- Find the documents' hidden meaning (obscured by noise)
- Distributional hypothesis: terms with similar/related meanings tend to co-occur with the same other words



The **Netherlands** is a constituent country of the Kingdom of the Netherlands, located mainly in North-West Europe and with several islands in the Caribbean. Mainland Netherlands borders the North Sea to the north and west, Belgium to the south, and Germany to the east, and shares maritime borders with Belgium, Germany and the United Kingdom. It is a parliamentary democracy organised as a unitary state. The country capital is Amsterdam and the seat of government is The Hague.

The main cities in **Holland** are Amsterdam, Rotterdam and The Hague. Amsterdam is formally the capital of the Netherlands and its largest city. The Port of Rotterdam is Europe's largest and most important harbour and port. The Hague is the seat of government of the Netherlands. These cities, combined with Utrecht and other smaller municipalities, effectively form a single city—a conurbation called Randstad.



Latent Semantic Analysis

Deerwester & Dumais, 1990

- Idea: transformation of the term space into a concept space that captures most of the variance in the collection
 - Documents and terms now indirectly related through concepts
 - The originally sparse space is now dense (dimensionality reduction)



 Documents sharing often co-occurring terms are close to each other in concept space

 $D_1 = \{Holland\}, D_2 = \{Netherlands\}$

• Distant in term space, close in concept space



LSA & Singular value decomposition

Deerwester & Dumais, 1990

• Factorization of term document matrix leads to concept space

- SVD to find correlations among the rows and columns
 - Every real matrix has a SVD decomposition



LSA & Singular value decomposition

Deerwester & Dumais, 1990

Perfect reconstruction of the term-document matrix (sparse) is possible

$$U\Sigma V^T = T$$

 Denser and smaller matrix if only the first k singular values are used (dimensionality reduction)

$$U\Sigma_k V^T \approx T$$
, with $k \ll m$

- Benefits of SVD
 - Transformation of correlated variables into uncorrelated ones
 - Ordering of dimensions according to the data's greatest variability
 - Best approximation of original data with fewer dimensions



LSA & Singular value decomposition

Deerwester & Dumais, 1990

- Keep the *k* most important concepts (often [100,300])
 - Best approximation of the term-document matrix in the least square sense
 - Each document is explained by concept vectors
- For retrieval purposes, the (keyword) query needs to be translated into the concept space: $q'=q^T U_k$
 - Retrieval as in the vector space model
- Good retrieval performance for small corpora
- Concept vectors not always easily explainable



Latent Semantic Analysis

Example concepts



highest scoring terms of three LSI concept vectors; corpus: GH95 (GeoCLEF)



Hofman, 1999 [4]

 Co-occurrence statistics modelled through latent class variables ("topics") in a statistical model

 Each word is generated from a single topic; different words in a document may be generated from different topics



Hofman, 1999 [4]

 Co-occurrence statistics modelled through latent class variables ("topics") in a statistical model Brazil has become the sixth-biggest economy in the world, overtaking the UK, the country's Each word is get finance minister says. a document may Carmaker Nissan is to build a new model at its Sunderland factory, with investment of £125m, which it says will create 2,000 jobs. Arsene Wenger's Arsenal will give it "a real go" as they try to overturn a 4-0 Champions League deficit against AC Milan. **Sports** Economy Cars



Hofman, 1999 [4]

- Co-occurrence statistics modelled through latent class variables ("topics") in a statistical model
- Each word is generated from a single topic; different words in a document may be generated from different topics
- Generative process of observing word *w* in document *D*
 - **1.** Select document *D* with probability P(D)
 - **2.** Pick latent class *z* from $\{z_1, ..., z_k\}$ with probability P(z|D)
 - **3.** Generate a word *w* with "factor" probability P(w|z)

$$P(w,D) = P(D) \sum_{z \in \mathbb{Z}} P(w \mid z) P(z \mid D)$$



Factor examples [4]

4 factors from a 128 factor decomposition (TDT-1 corpus) 8 most likely terms per factor (i.e. P(w|z))





Summary

- Semantics
 - Important
 - And cool 🙂
- Two camps
 - Knowledge-based
 - Statistics-based
- More on Monday!



Sources

- 1 Introduction to Information Retrieval. Manning et al. 2008
- 2 Information retrieval. Keith van Rijsbergen. 1979
- 3 Managing gigabytes, Witten et al. 1999
- 4 Unsupervised learning by probabilistic latent semantic analysis. Hofmann. 2001
- Cumulated gain-based evaluation of IR techniques. Järvelin
 & Kekäläinen. 2002
- 6 Semantic distance in WordNet: an experimental, application-oriented evaluation of five measures. Budanitsky & Hirst. 2001
- Correcting real-world spelling errors by restoring lexical cohesion. Hirst & Budanitsky. 2005

