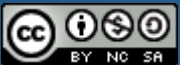


7. Query Expansion and Relevance Feedback

Prof. Dr. Goran Glavaš

Data and Web Science Group
Fakultät für Wirtschaftsinformatik und Wirtschaftsmathematik
Universität Mannheim



Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International

After this lecture, you'll...

2

- Know about methods for improving the recall of IR engines
- Know what (pseudo) relevance feedback is
- Understand different methods for expanding initial user queries
- Learn how to exploit „click data” for query expansion

Outline

3

- [Recap of Lecture #6](#)
- Improving recall in IR
- Relevance feedback
- Query expansion
- User interaction data („click” data)

Recap of the previous lecture

4

- Probabilistic retrieval
 - **Q:** Why probability theory in IR, and why probabilistic ranking?
 - **Q:** What are the uncertainties of the IR process that we model probabilistically?
- Language modeling
 - **Q:** What is a language model?
 - **Q:** How do we estimate probabilities of sequences of words?
 - **Q:** What language models are used in IR? Why? Explain the sparseness issue of LMs.
- Language modeling for IR
 - **Q:** What probability does the query likelihood model estimate?
 - **Q:** How do we estimate the probability of the query given the document?
 - **Q:** What happens if some query term is not present in the document?
 - **Q:** Explain the differences between smoothing techniques.

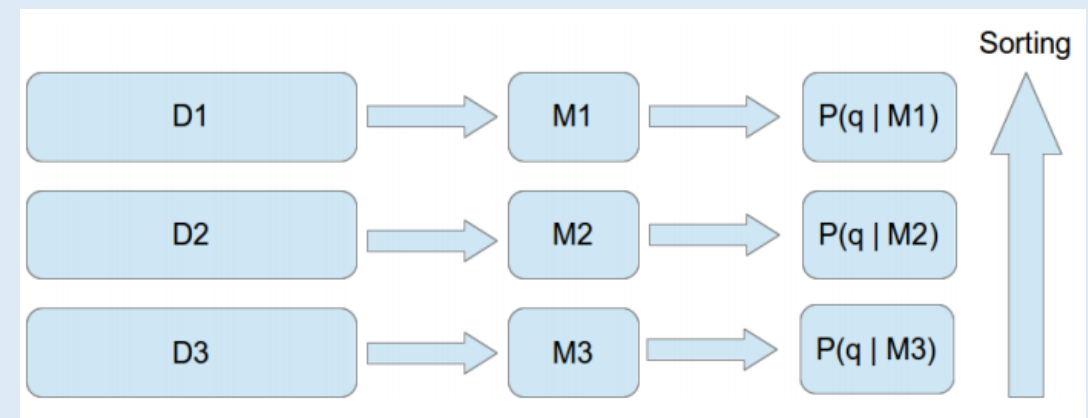
Query likelihood model for ranking

5

- Given a document collection D and a query q we need to estimate the probability $P(q | d)$ for every document d in D
- In the **query likelihood model**, we estimate the probability $P(q | d)$ as the probability that the language model built from d generates the query q

- Algorithm

1. Compute the language model M_i for every document d_i in D
2. Compute the probability $P(q | M_i)$ for every language model M_i



- **Intuition:** Language models of relevant documents should assign higher probability for the query

Query likelihood model for ranking – example

6

- We are given a toy collection consisting of three documents
 - d_1 : „Sam chased the orc with the sword”
 - d_2 : „Frodo and Sam stabbed orcs”
 - d_3 : „Sam took the sword”
- We are given the query „Sam and orc and sword”
- Let’s rank the documents according to **unigram LM for IR** (ignore stopwords)
- **Step 1:** Compute language models of individual documents
 - M_1 : $P(\text{„sam"}) = 0.25, P(\text{„chase"}) = 0.25, P(\text{„orc"}) = 0.25, P(\text{„sword"}) = 0.25$
 - M_2 : $P(\text{„frodo"}) = 0.25, P(\text{„sam"}) = 0.25, P(\text{„stab"}) = 0.25, P(\text{„orc"}) = 0.25$
 - M_3 : $P(\text{„sam"}) = 0.33, P(\text{„took"}) = 0.33, P(\text{„sword"}) = 0.33$

Jelinek-Mercer smoothing

7

- Laplace smoothing assumes that **all unseen words** are **equally likely**
- **Jelinek-Mercer smoothing** (also known as **interpolated smoothing**)
 1. Additionally builds a language model M_D from the whole document collection D
 2. Interpolates between probabilities of the query according to the
 - **Local LM** – language model M_d built from the particular document d
 - **Global LM** – language model M_D built from the whole collection

$$P(t_i|M_d) = \lambda \cdot P(t_i|M_d) + (1 - \lambda) \cdot P(t_i|M_D)$$

- The probability of a word unseen in the document d still gets **some probability** from the global language model
 - Probability of an unseen word depends on its frequency in whole collection
- **Q:** What is some query term doesn't appear in the whole collection D ?

Outline

8

- Recap of Lecture #6
- **Improving recall in IR**
- Relevance feedback
- Query expansion
- User interaction data („click” data)

Improving recall of IR systems

9

- Most ranked retrieval systems optimize **precision**
 - Implicit assumptions:
 1. Plenty of relevant documents
 2. Users only look at top-ranked results
 - **Top-ranked results should be relevant**
 - Not so bad if we miss some of the relevant documents
- However, sometimes **recall matters more** than precision
 - Security & intelligence
 - Patent & publication search

Improving recall of IR systems

10

- Sometimes **recall matters more** than precision
 - Retrieving **all** relevant documents is essential
 - Even at cost of generating many high-ranked irrelevant documents
- **Q:** How to increase recall?
- **A:** By modifying the initial query
 1. By adding new terms
 - E.g., terms that are semantically related to the terms of the original query
 2. By making the query vector more similar to vectors of relevant documents
 - We need **some indication of relevance**
 - Either provided by the user or assumed by the model

Improving recall of IR systems

11

- Methods for improving recall of IR systems

1. **Global method:** query expansion

- Adding new terms to the initial query
- New terms somehow related to original terms
- New terms identified using thesauri, distributional semantics, or lexical association

2. **Local methods**

- Relevance feedback – requires additional input from the user
- Pseudo relevance feedback (automated assumption of relevance)
- Relevance model (automated assumption of relevance)

Outline

12

- Recap of Lecture #6
- Improving recall in IR
- **Relevance feedback**
- Query expansion
- User interaction data („click” data)

Relevance feedback

13

- **Relevance feedback** means improving the query based on **user feedback on relevance** of the documents in the initial set of results
- Workflow:
 1. User issues a (usually short and simple) initial query
 2. Search engine retrieves and ranks the results
 3. User **explicitly marks** some of the results as relevant and/or non-relevant
 4. The search engine computes the **new query** (i.e., a better representation of the information need) **based on the feedback**
- There may be more than one iteration of the steps listed above
- **Assumption** motivating relevance feedback: it is often **difficult to formulate a good (discriminative) query** when you don't know the collection well

Relevance feedback

14

- Some definitions
 - **Ad hoc retrieval** – regular retrieval without relevance feedback
 - We'll call the first (uninformed) query being executed an **initial query** (q_0)
 - **Feedback set** – the first set of documents that are retrieved

Relevance feedback

15

■ **Algorithm** (direct addition of query terms)

1. Retrieve an initial ranked list of hits for the user's initial query q_0
2. **Ask user**: which documents in the feedback set are relevant?
3. Identify the best terms with which to extend the query
 - Good terms occur frequently in relevant documents
 - Good terms occur infrequently in non-relevant documents
 - Terms can have different importance weights
 - Form the vector q_{RF} containing the weighted new terms
4. Expand the query with new terms
 - New terms could be additionally weighted with respect to original terms
 - $q' = f(q_0, q_{RF})$
5. Retrieve the ranking for the expanded query q'

Initial query and results

16

- Initial query: „*New space satellite applications*”
 - + 1. 0.539, 08/13/91, [NASA Hasn't Scrapped Imaging Spectrometer](#)
 - + 2. 0.533, 07/09/91, [NASA Scratches Environment Gear From Satellite Plan](#)
 - 3. 0.528, 04/04/90, [Science Panel Backs NASA Satellite Plan, But Urges Launches of Smaller Probes](#)
 - 4. 0.526, 09/09/91, [A NASA Satellite Project Accomplishes Incredible Feat: Staying Within Budget](#)
 - 5. 0.525, 07/24/90, [Scientist Who Exposed Global Warming Proposes Satellites for Climate Research](#)
 - 6. 0.524, 08/22/90, [Report Provides Support for the Critics Of Using Big Satellites to Study Climate](#)
 - 7. 0.516, 04/13/87, [Arianespace Receives Satellite Launch Pact From Telesat Canada](#)
 - + 8. 0.509, 12/02/87, [Telecommunications Tale of Two Companies](#)
- User then marks relevant documents (+).

Expanded query after relevance feedback

17

- 2.074 new
- 30.816 satellite
- 5.991 nasa
- 4.196 launch
- 3.516 instrument
- 3.004 bundespost
- 2.790 rocket
- 2.003 broadcast
- 0.836 oil
- 15.106 space
- 5.660 application
- 5.196 eos
- 3.972 aster
- 3.446 arianespace
- 2.806 ss
- 2.053 scientist
- 1.172 earth
- 0.646 measure

Importance weight

Results for the expanded query

18

- 2 1. 0.513, 07/09/91, [NASA Scratches Environment Gear From Satellite Plan](#)
- 1 2. 0.500, 08/13/91, [NASA Hasn't Scrapped Imaging Spectrometer](#)
- 8 3. 0.492, 12/02/87, [Telecommunications Tale of Two Companies](#)
4. 0.493, 08/07/89, [When the Pentagon Launches a Secret Satellite, Space Sleuths Do Some Spy Work](#)
5. 0.493, 07/31/89, [NASA Uses 'Warm' Superconductors For Fast Circuit](#)
6. 0.491, 07/09/91, [Soviets May Adapt Parts of SS-20 Missile For Commercial Use](#)
7. 0.490, 07/12/88, [Gaping Gap: Pentagon Lags in Race To Match the Soviets In Rocket Launchers](#)
8. 0.490, 06/14/90, [Rescue of Satellite By Space Agency To Cost \\$90 Million](#)

- Documents **explicitly marked** as relevant by users will be ranked higher
 - So will the documents *similar* to them

Relevance feedback in probabilistic retrieval

19

- **Recap**: In probabilistic retrieval, relevance feedback can be incorporated into estimations of probabilities $P(D_t \mid Q, r)$ and $P(D_t \mid Q, \neg r)$
- Estimates **without** relevance feedback
 - $P(D_t \mid Q, r) = 0.5$
 - $P(D_t \mid Q, \neg r) = N_t / N$
- Estimates **with** relevance feedback
 - $P(D_t \mid Q, r) = (r_t + 0.5) / (R + 1)$
 - $P(D_t \mid Q, \neg r) = (N_t - r_t + 0.5) / (N - R + 1)$

Relevance feedback in VSM

20

- **Rocchio algorithm** –used to incorporate relevance feedback into the vector space model
- This algorithm is looking to rewrite the query so that it's vector is:
 - As **similar** as possible to the vectors of relevant documents
 - As **dissimilar** as possible to the vectors of non-relevant documents
- **Open questions**
 - How do we aggregate vectors of all (non-)relevant docs in one vector?
 - Do we take the original query into account?

Rocchio algorithm

21

- How do we aggregate vectors of all (non-)relevant docs in one vector?
 - Key concept: **centroid**
- Centroid is the center of mass of a set of points
 - In VSM, we represent documents as points in a high-dimensional space
 - So, given a set of documents (their VSM vectors), we can compute their centroid
- Centroid μ of the set of document vectors C is computed simply as the **mean of vectors** of individual documents

$$\frac{1}{|C|} \sum_{d \in C} d$$

Rocchio algorithm

22

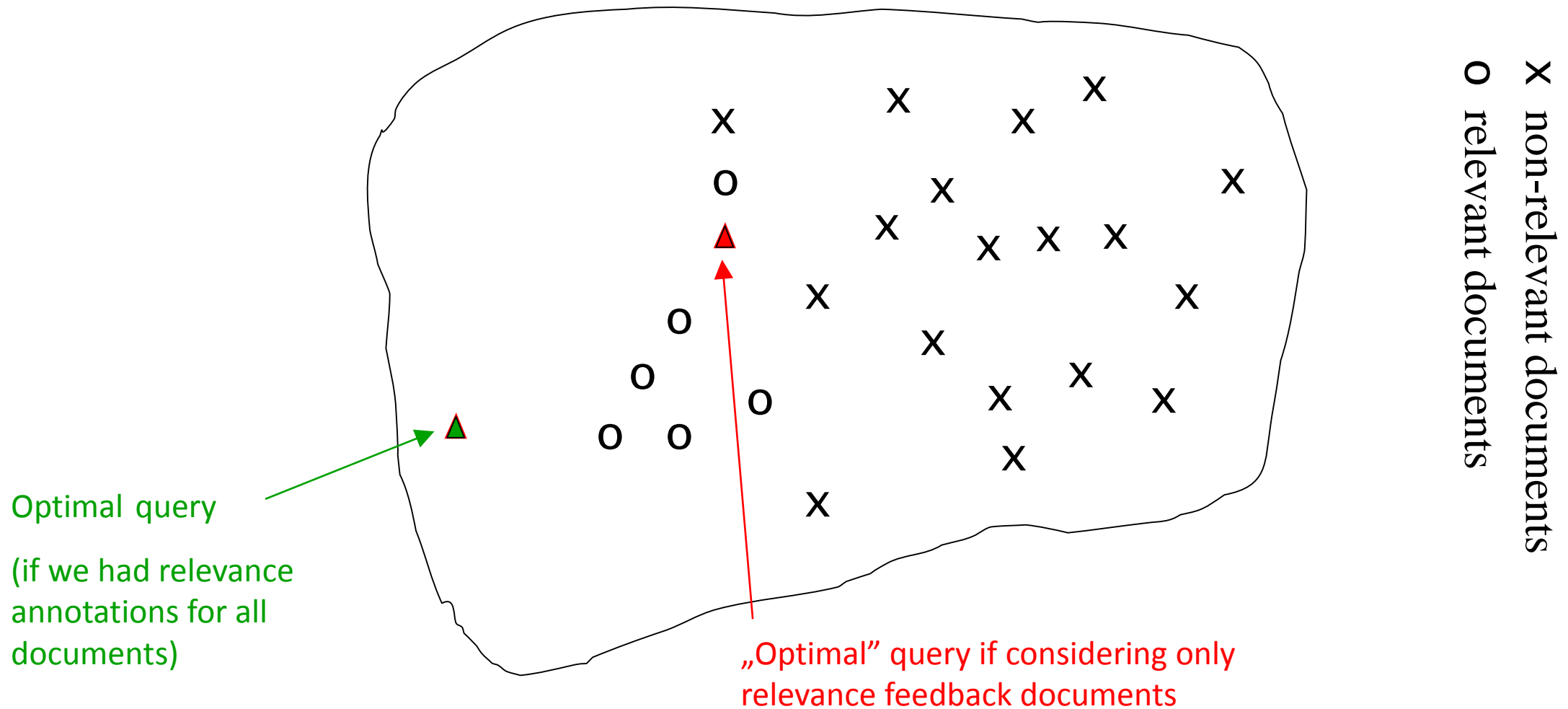
- **Q:** Do we take the original query q_o into account?
- **A:** Theoretically, we don't have to
 - Because Rocchio just needs relevant and non-relevant documents to compute the „optimal” query:

$$q_{opt} = \operatorname{argmax}_q \left(\cos \left(q, \frac{1}{|C_r|} \sum_{d_r \in C_r} d_r \right) - \cos \left(q, \frac{1}{|C_{nr}|} \sum_{d_{nr} \in C_{nr}} d_{nr} \right) \right)$$

- **Q:** What might be the problem with computing the optimal query only from relevance feedback?
- **A:** Users provide relevance judgements only for **a small number of documents**, (not all documents in the collection!)

Rocchio algorithm

23



Rocchio algorithm

24

- We are given **only a handful** of relevance feedback annotations
- Thus, we re-estimate the query by combining
 1. Centroid of relevant documents
 2. Centroid of non-relevant documents
 3. **Initial query vector** q_0

$$q_m = \alpha \cdot q_0 + \left(\beta \cdot \frac{1}{|D_r|} \sum_{d_j \in D_r} d_j \right) - \left(\gamma \cdot \frac{1}{|D_{nr}|} \sum_{d_j \in D_{nr}} d_j \right)$$

- D_r is the set of vectors of known relevant documents (different from C_r)
- D_{nr} is the set of vectors of known non-relevant documents (different from C_{nr})
- α , β , and γ are weights, determining the contribution of each component (set beforehand or empirically)
- New query moves **towards the relevant** and **away from non-relevant** documents

Rocchio algorithm

25

- Tradeoff between α and β/γ depends on the number of judged documents
 - A few documents judged – $\alpha > \beta$ and γ
 - More emphasis on the **original query**
 - A lot of documents judged – $\alpha < \beta$ and γ
 - More emphasis on **relevance feedback**
- In the modified query, weights of some terms may become negative
 - Terms that appear more frequently in **non-relevant** queries
 - We just **ignore** negative weights (i.e., we reset them to **0**)

Rocchio algorithm

26

- Positive vs. Negative feedback
 - **Q:** Should we emphasize more relevant or non-relevant documents?
 - I.e., should β be larger than γ or vice-versa?
 - **A:** positive feedback is more valuable than negative feedback
 - So, we should set $\gamma < \beta$, e.g., $\beta = 0.75$ and $\gamma = 0.25$
- Many systems allow only positive feedback (i.e., $\gamma = 0$)
 - **Q:** Why?

Relevance feedback – issues

27

- **Q:** When does relevance feedback work?
- **A:** Relevance feedback works when
 1. User has **sufficient knowledge** for initial query
 2. Relevance for initial query „**behaves nicely**”
 - **Case #1:** All relevant documents are tightly clustered around a single relevance vector
 - **Case #2:** There are several different clusters of relevant documents, but they have significant vocabulary overlap
 - **Similarities** (term overlap) between relevant and non-relevant documents are **small**

Relevance feedback – issues

28

- Cases when relevance feedback **won't work**:
 1. User's knowledge **insufficient** for initial query
 - Misspellings (e.g., „**Brittany Speers**”), vocabulary mismatches („**cosmonaut**” vs. „**astronaut**”)
 2. Multiple relevance vectors (groups) with little lexical overlap
 - When query contains **general concepts with multiple senses**
 - E.g., „**contradictory government policies**”
„**pop stars that worked at Burger King**”

Relevance feedback – issues

29

■ Efficiency

- Extending the queries with additional terms (directly or indirectly, e.g., via Rocchio) may lead to long queries
- Long queries are **inefficient** to process
- Solution: limit the number of terms to be added
 - Criteria – e.g., collection-wide frequency

■ User issues

- Users often **reluctant** to provide explicit feedback
 - **Time consuming**, sometimes also quite **cognitively demanding**
- **Harder** for users to understand why a particular document was retrieved after applying relevance feedback

Pseudo-relevance feedback

30

- **Pseudo-relevance feedback** is relevance feedback **without the human in the loop**
 - I.e., we automate the „manual” part of relevance feedback
 - Eliminate the need for human relevance judgements
- Algorithm
 1. Retrieve an initial ranked list of hits for the user's initial query q_0
 2. Assume that the top K documents are relevant
 3. Compute expansions terms, q
 4. Expand the query (the vector q_{RF} containing the weighted new terms)
 5. Retrieve the ranking for the expanded query q'

$$q' = \lambda \cdot q_0 + (1 - \lambda) \cdot q_{RF}$$

- **Q:** How do we compute the terms with which to expand the query?

Relevance model (Lavrenko, 2001)

31

- Input
 - Initial query q_0
 - Top K documents in the ranking for initial query – d_1, d_2, \dots, d_K
 - Relevance probabilities of top ranked documents for the initial query – $P(d_i|q_0)$
- Output
 - A distribution of terms denoting how well they describe the initial query q_0
 - An importance/probability of term w for q_0 query is computed as follows:

$$P(w|q_0) = \sum_{i=1}^K P(w|d_i) \cdot P(d_i|q_0)$$

- Rank the terms in decreasing order of $P(w|q_0)$, take top N terms and combine them into a weighted expansion query q_{PRF}

Relevance model – expansion lists

32

- From: [Croft et al.](#) „Search Engines – Information Retrieval in Practice“ (Chapter 7)

<i>president lincoln</i>	<i>abraham lincoln</i>	<i>fishing</i>	<i>tropical fish</i>
lincoln	lincoln	fish	fish
president	america	farm	tropic
room	president	salmon	japan
bedroom	faith	new	aquarium
house	guest	wild	water
white	abraham	water	species
america	new	caught	aquatic
guest	room	catch	fair
serve	christian	tag	china
bed	history	time	coral
washington	public	eat	source
old	bedroom	raise	tank
office	war	city	reef
war	politics	people	animal
long	old	fishermen	tarpon
abraham	national	boat	fishery

Relevance model vs. Rocchio algorithm

33

- Let's compare Lavrenko's relevance model with Rocchio algorithm
 - Assume Rocchio considers top K initially ranked documents as relevant (D_r) and does not consider non-relevant documents ($\gamma = 0$)

- Lavrenko's relevance model

$$q' = \lambda \cdot q_0 + (1 - \lambda) \cdot q_{RF} \quad P(w|q_0) = \sum_{i=1}^K P(w|d_i) \cdot P(d_i|q_0)$$

- Rocchio algorithm

$$q_m = \alpha \cdot q_0 + \beta \cdot \frac{1}{|D_r|} \sum_{d_j \in D_r} d_j$$

Rocchio uses all terms, **RM uses only top N terms**

Rocchio computes simple average, **RM weighted average with document relevances for query $P(d_i|q_0)$ as weight**

Rocchio uses TF-IDF weights, **RM uses $P(w|d_i)$**

Pseudo-relevance feedback

34

- Often works very well, but **may go wrong**
 - **Crucially depends** on how good the initially top-ranked K documents are
 - On average **works better** than query expansion
 - As evaluated on standard benchmark datasets
- **Pitfall** of pseudo-relevance feedback: **query drift**
 - E.g., imagine that for initial query „copper mines“, top K documents retrieved mostly describe **copper mines in Chile**
 - The final result may be more about **Chile** than about **copper mines**
- Query drift is the reason why typically **only one iteration** of pseudo-relevance feedback is performed

Outline

36

- Recap of Lecture #6
- Improving recall in IR
- Relevance feedback
- Query expansion
- User interaction data („click” data)

Query expansion

37

- **Query expansion** techniques try to improve retrieval results by adding terms to the user query
- Relevance model is a form of automated query expansion
- Next, we will consider other ways of expanding queries
 - Independent of the initial ranking of documents
 - I.e., when the content of top-ranked documents **is not used** for query expansion
 - We use **external resources** to expand the query
 1. Dictionaries / thesauri
 2. Large corpora

Query expansion

38

- There are two types of query expansions:
 1. **Explicit** query expansion
 - Initial query is executed and results retrieved and ranked
 - Expanded version(s) of the initial query proposed to the user for subsequent searches
 2. **Implicit** query expansion
 - The original query is implicitly expanded and results are obtained by executing the expanded query
- **Q:** How do we identify good expansion words?

- The following are common methods of query expansion:

1. Controlled vocabulary

- There is a canonical term for each concept, all other terms for the same concept are replaced by the canonical term
- E.g., {„hotel”, „apartment”, „room”} → „accommodation”; „Burma” → „Myanmar”

2. Manual thesaurus

- Human annotators build sets of synonymous terms for concepts, without designating a canonical terms
- E.g., {„hotel”, „apartment”, „room”, „accommodation”}; {„Burma”, „Myanmar”}
- WordNet (<https://wordnet.princeton.edu/>)
 - General large lexical database for English language, contains 117K synonym sets (*synsets*)
- Domain specific Thesauri – e.g., NIH for medicine (<https://ncit.nci.nih.gov/ncitbrowser>)
- Knowledge bases, e.g., DBPedia or Yago can also be used for query expansion

- The following are common methods of query expansion:

3. Automatically generated thesaurus

- Co-occurrence statistics (lexical association metrics) on a large external corpus (e.g., Wikipedia) are used to automatically induce a large thesaurus
- E.g., „hotel” often co-appears in text with „accommodation”

4. Query log mining

- Query reformulations done by users are logged
- These reformulations can be used for query expansion for similar queries of new users
- Requires huge amounts of logged queries – feasible only in web search

Thesaurus-based query expansion

41

- **Algorithm** is simple
 - For each term t from the initial query q_0 , look up synonyms and related terms in the thesaurus
 - E.g., „frodo” -> „frodo baggins”, „hobbit frodo”
 - Optionally, assign new terms lower weights than the original terms
- Thesaurus-based query expansion
 - **Generally increases recall**
 - Especially in specific domains, with rich domain-specific thesauri
 - May **significantly decrease precision**
 - Particularly when expanding an ambiguous query term
 - E.g., „interest rate” -> „interest, fascinate, rate, evaluate”

Thesaurus-based query expansion

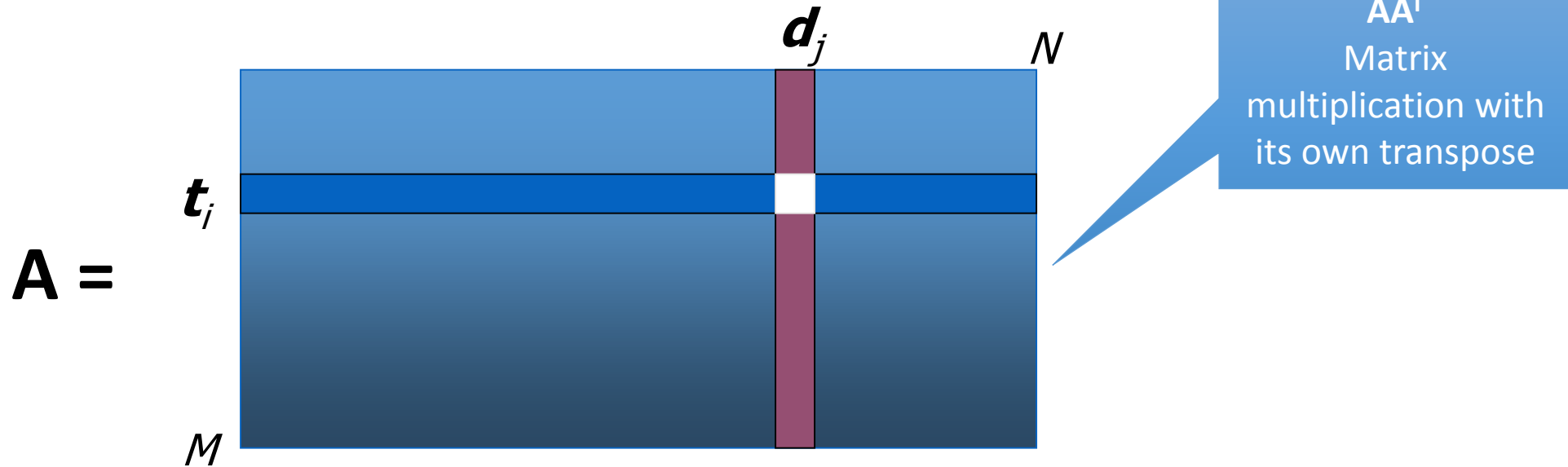
42

- Manually producing a thesaurus is **time-consuming** and **expensive**
 - Additionally, it needs to be **constantly updated** to reflect changes in the domain
- **Automated thesaurus generation**
 - Generating thesaurus by detecting similarity/relatedness of terms in a large corpora
 - **Distributional hypothesis** – words are similar if they occur in similar contexts
 - E.g., „apple” is similar to „pear” as you can both *harvest*, *peel*, *prepare* and *eat* both
 - **Related words** – words that often co-appear are semantically related
 - E.g., „pilot” and „airplane”

Co-occurrence thesaurus

43

- Simplest way to compute a thesaurus is based on **term-term similarities** in $C = AA^T$ where A is **term-document matrix**
- $w_{i,j}$ = (normalized) weight for (t_i, d_j)



- For each t_i , pick terms with high values in C (term-term matrix)

Automated thesaurus generation

44

- Expansion based on distributional vectors has **shortcomings**
 - Distributional vectors cannot distinguish **similarity** from **relatedness**
 - E.g., two **synonyms** might be as similar as two **antonyms**
 - Making distributional vectors encode **only similarity**, **not relatedness** – active research area
 - Semantic similarity and relatedness are **vague**, **not crisp** relations
 - Word with **multiple senses** especially problematic

- Quality of produced term associations often **not good enough** for IR tasks

Word	Nearest neighbors
absolutely	absurd, whatsoever, totally, exactly, nothing
bottomed	dip, copper, drops, topped, slide, trimmed
captivating	shimmer, stunningly, superbly, plucky, witty
doghouse	dog, porch, crawling, beside, downstairs
makeup	repellent, lotion, glossy, sunscreen, skin, gel
mediating	reconciliation, negotiate, case, conciliation
keeping	hoping, bring, wiping, could, some, would
lithographs	drawings, Picasso, Dali, sculptures, Gauguin
pathogens	toxins, bacteria, organisms, bacterial, parasite
senses	grasp, psyche, truly, clumsy, naive, innate

Outline

45

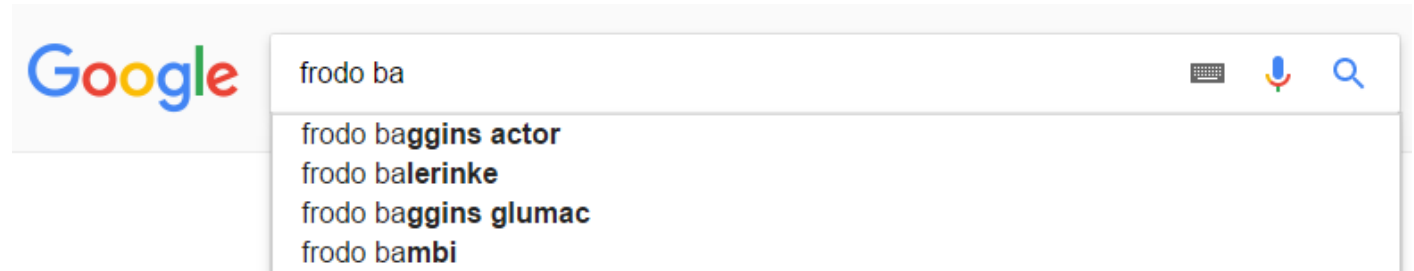
- Recap of Lecture #6
- Improving recall in IR
- Relevance feedback
- Query expansion
- User interaction data („click” data)

Interaction data

46

- Using **data collected** from **users' interaction** with the search engine to improve the search process
 - Can be used to assist the user in formulating the query

- Common example
 - Query auto-completion



- **Indirect relevance feedback**
 - Assumption: ranked documents **clicked** by the user are **relevant** for the query
 - **Clickstream mining**:
 - Documents are considered more relevant, the more clicked they are
 - Global metric of document relevance on the web (e.g., like PageRank)
 - One feature in machine learning-based ranking functions

Click log data for re-ranking

47

ALL RESULTS 1-10 of 131,000 results - [Advanced](#)

RELATED SEARCHES
[CIKM 2008](#)

SEARCH HISTORY
Turn on search history to start remembering your searches.
[Turn history on](#)

ALL RESULTS

[CIKM 2008 | Home](#)
Napa Valley Marriott Hotel & Spa: Napa Valley, California October 26-30, 2008
[cikm2008.org](#) - [Cached page](#)

[Papers](#) [Program Committee](#)
[Themes](#) [News](#)
[Important Dates](#) [Napa Valley](#)
[Banquet](#) [Posters](#)
[Show more results from cikm2008.org](#)

[Conference on Information and Knowledge Management \(CIKM\)](#)
Provides an international forum for presentation and discussion of research on information and knowledge management, as well as recent advances on data and knowledge bases ...
[www.cikm.org](#) - [Cached page](#)

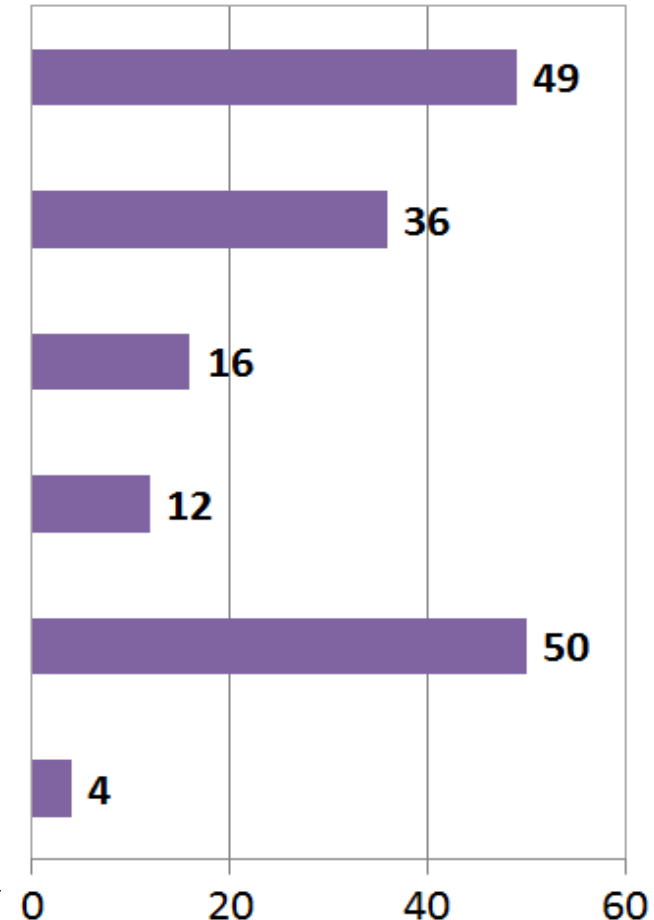
[Conference on Information and Knowledge Management \(CIKM'02\)](#)
SAIC Headquarters, McLean, Virginia, USA, 4-9 November 2002.
[www.cikm.org/2002](#) - [Cached page](#)

[ACM CIKM 2007 - Lisbon, Portugal](#)
News and announcements: 12/02 - Best interdisciplinary paper award at CIKM 2007 went to Fei Wu and Daniel Weld for Autonomously Semantifying Wikipedia.
[www.fc.ul.pt/cikm2007](#) - [Cached page](#)

[CIKM 2009 | Home](#)
CIKM 2009 (The 18th ACM Conference on Information and Knowledge Management) will be held on November 2-6, 2009, Hong Kong. Since 1992, CIKM has successfully brought together ...
[www.comp.polyu.edu.hk/conference/cikm2009](#) - [Cached page](#)

[Conference on Information and Knowledge Management \(CIKM\)](#)
CIKM Conference on Information and Knowledge Management The Conference on Information and Knowledge Management (CIKM) provides an international forum for presentation and ...
[cikmconference.org](#) - [Cached page](#)

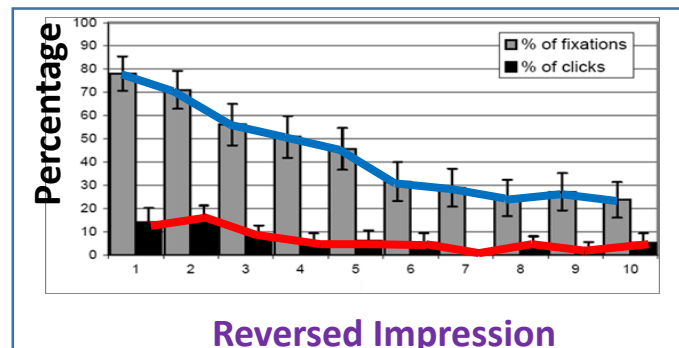
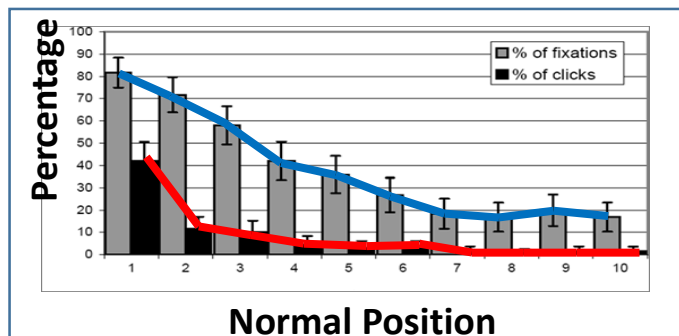
of clicks received



Interaction data

48

- How large is the click-log?
 - **bing** search logs grow **10+ TB every day**
 - **Efficient algorithms** for mining click-logs are needed
- Click-position bias
 - Higher positions receive more **user attention (eye fixation)** and **clicks** than lower positions
 - True, even when the **order** or results is **reversed**!
 - Clicks are **informative**, but **biased** – top-ranked documents clicked even when non-relevant



Interaction data – conclusion

49

- User behavior is an **intriguing source of relevance data**
 - Users make (**somewhat**) informed choices when they interact with search engines
 - Potentially **a lot of data available** in search logs
- But there are **significant caveats**
 - User behavior data can be very **noisy**
 - **Interpreting** user behavior can be **tricky**
 - **Spam** can be a significant problem
 - Not all queries will have user behavior

Now you...

50

- Know about methods for improving the recall of IR engines
- Know what (pseudo) relevance feedback is
- Understand different methods for expanding initial user queries
- Have an idea on how to exploit „click data” for query expansion