

Introduction to Information Retrieval

<http://informationretrieval.org>

IIR 9: Relevance Feedback & Query Expansion

Hinrich Schütze

Institute for Natural Language Processing, Universität Stuttgart

2008.06.03

- 1 Recap
- 2 Relevance feedback: Basics
- 3 Relevance feedback: Details
- 4 Global query expansion

1 / 57

2 / 57

Plan for this lecture

- First: Recap
- Main topic today: How can we improve **recall** in search?
 - "aircraft" in query doesn't match with "plane" in document
 - "heat" in query doesn't match with "thermodynamics" in document
- Options for improving recall
 - Local methods: Do a "local", on-demand analysis for a user query
 - Main local method: **relevance feedback**
 - Global methods: Do a global analysis once (e.g., of collection) to produce **thesaurus**
 - Use thesaurus for **query expansion**

Google example query:
~hospital -hospital -hospitals

3 / 57

4 / 57

- 1 Recap
- 2 Relevance feedback: Basics
- 3 Relevance feedback: Details
- 4 Global query expansion

- 1 Recap
- 2 Relevance feedback: Basics
- 3 Relevance feedback: Details
- 4 Global query expansion

5 / 57

Relevance

- We will evaluate the quality of an information retrieval system and, in particular, its ranking algorithm with respect to **relevance**.
- A document is relevant if it gives the user the information she was looking for.
- To evaluate relevance, we need an **evaluation benchmark** with three elements:
 - A benchmark document collection
 - A benchmark suite of queries
 - An assessment of the relevance of each query-document pair

7 / 57

6 / 57

Relevance: query vs. information need

- The notion of "relevance to the query" is very problematic.
- **Information need i** : You are looking for information on whether drinking red wine is more effective at reducing your risk of heart attacks than white wine.
- **Query q** : WINE AND RED AND WHITE AND HEART AND ATTACK
- Consider document d' : *He then launched into the heart of his speech and attacked the wine industry lobby for downplaying the role of red and white wine in drunk driving.*
- d' is relevant to the query q , but d' is **not** relevant to the information need i .
- User happiness/satisfaction (i.e., how well our ranking algorithm works) can only be measured **by relevance to information needs, not by relevance to queries**.

8 / 57

- Precision (P) is the fraction of retrieved documents that are relevant

$$\text{Precision} = \frac{\#(\text{relevant items retrieved})}{\#(\text{retrieved items})} = P(\text{relevant}|\text{retrieved})$$

- Recall (R) is the fraction of relevant documents that are retrieved

$$\text{Recall} = \frac{\#(\text{relevant items retrieved})}{\#(\text{relevant items})} = P(\text{retrieved}|\text{relevant})$$

- F allows us to trade off precision against recall.

- Balanced F :

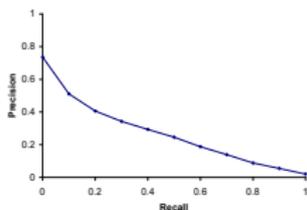
$$F_1 = \frac{2PR}{P+R}$$

- This is a kind of soft minimum of precision and recall.

9 / 57

10 / 57

Averaged 11-point precision/recall graph



- This curve is typical of performance levels for the TREC benchmark.
- 70% chance of getting the first document right (roughly)
- When we want to look at at least 50% of all relevant documents, then for each relevant document we find, we will have to look at about two nonrelevant documents.
- That's not very good.
- High-recall retrieval is an unsolved problem.

11 / 57

Outline

- 1 Recap
- 2 Relevance feedback: Basics
- 3 Relevance feedback: Details
- 4 Global query expansion

12 / 57

Relevance feedback: Basic idea

- User issues a (short, simple) query
- Search engine returns set of docs
- User marks some docs as relevant, some as nonrelevant
- Search engine computes a new representation of information need – better than initial query
- Search engine runs new query and returns new results
- New results have (hopefully) better recall.
- We can iterate this.
- We will use the term *ad hoc retrieval* to refer to regular retrieval without relevance feedback.
- We will now look at three different examples of relevance feedback that highlight different aspects of the process.

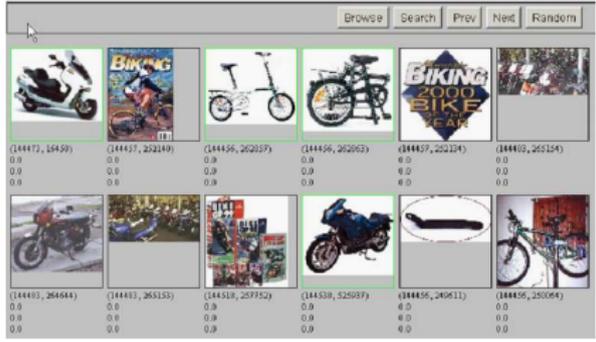
Relevance Feedback: Example



Results for initial query



User feedback: Select what is relevant

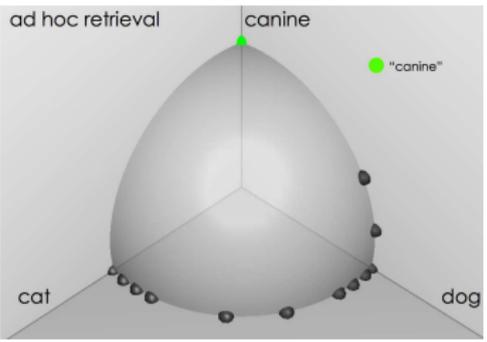


Results after relevance feedback

Browse Search Prev Next Random

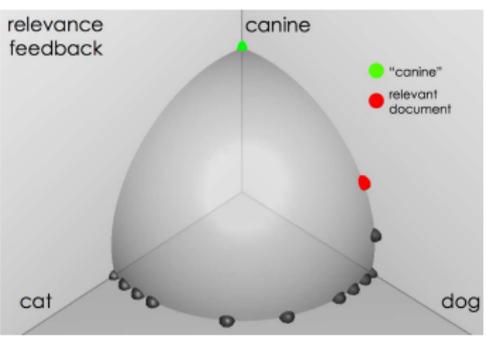
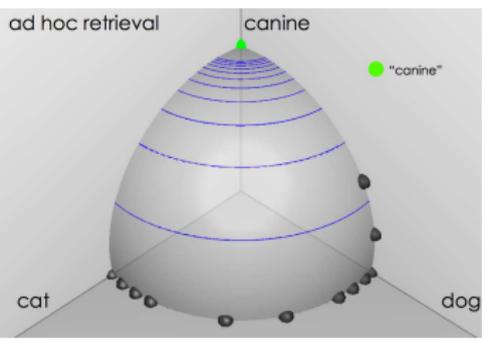
(144338, 328493) 0.54352 0.211944 0.509026	(144338, 328383) 0.50536296 0.357204 0.292069	(144338, 321522) 0.3594179 0.3206051 0.303530	(144436, 233099) 0.743011 0.531945 0.202615	(144436, 233061) 0.659275 0.411745 0.22853	(144338, 328199) 0.6070137 0.53053 0.509059
(144473, 10249) 0.6721 0.59522 0.27826	(144436, 209038) 0.675018 0.4459 0.213116	(144436, 213027) 0.695901 0.47945 0.305451	(144473, 10320) 0.708359 0.509102 0.501337	(144436, 203064) 0.70170296 0.56176 0.303948	(144473, 312419) 0.79207 0.449111 0.235059

Ad hoc retrieval for query "canine" (1)

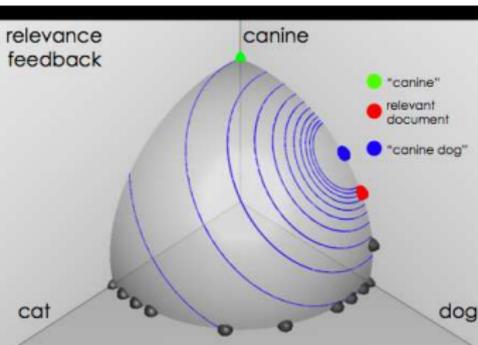


Ad hoc retrieval for query "canine" (2)

User feedback: Select what is relevant



Results after relevance feedback



source:
Fernando Díaz

Results for initial query

Initial query: New space satellite applications Results for initial

query:

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, Arianspace Receives Satellite Launch Pact From Telesat Canada
8. 0.509, 12/02/87, Telecommunications Tale of Two Companies

User then marks relevant documents with "+".

21 / 57

22 / 57

Expanded query after relevance feedback

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

Results for expanded query

- * 1. 0.513, 07/09/91, NASA Scratches Environment Gear From Satellite Plan
- * 2. 0.500, 08/13/91, NASA Hasn't Scrapped Imaging Spectrometer
- 3. 0.493, 08/07/89, When the Pentagon Launches a Secret Satellite, Space Sleuths Do Some Spy Work of Their Own
- 4. 0.493, 07/31/89, NASA Uses 'Warm' Superconductors For Fast Circuit
- * 5. 0.492, 12/02/87, Telecommunications Tale of Two Companies
- 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

23 / 57

24 / 57

- 1 Recap
- 2 Relevance feedback: Basics
- 3 Relevance feedback: Details
- 4 Global query expansion

- The centroid is the center of mass of a set of points.
- Recall that we represent documents as points in a high-dimensional space.
- Thus: we can compute centroids of documents.
- Definition:

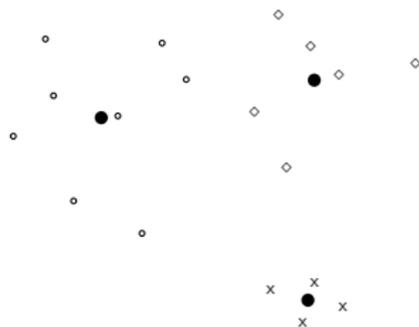
$$\vec{\mu}(D) = \frac{1}{|D|} \sum_{d \in D} \vec{v}(d)$$

where D is a set of documents and $\vec{v}(d) = \vec{d}$ is the vector we use to represent the document d .

25 / 57

26 / 57

Centroid: Examples



27 / 57

Rocchio algorithm

- The Rocchio algorithm implements relevance feedback in the vector space model.
- Rocchio chooses the query \vec{q}_{opt} that maximizes

$$\vec{q}_{opt} = \arg \max_{\vec{q}} [\text{sim}(\vec{q}, D_r) - \text{sim}(\vec{q}, D_{nr})]$$

- Closely related to maximum separation between relevant and nonrelevant docs
- This optimal query vector is:

$$\vec{q}_{opt} = \frac{1}{|D_r|} \sum_{\vec{d}_j \in D_r} \vec{d}_j + \left[\frac{1}{|D_r|} \sum_{\vec{d}_j \in D_r} \vec{d}_j - \frac{1}{|D_{nr}|} \sum_{\vec{d}_j \in D_{nr}} \vec{d}_j \right]$$

D_r : set of relevant docs; D_{nr} : set of nonrelevant docs

28 / 57

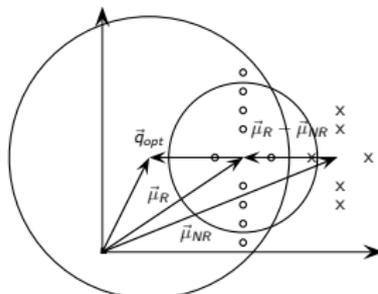
Rocchio algorithm

- The optimal query vector is:

$$\vec{q}_{opt} = \frac{1}{|D_r|} \sum_{\vec{d}_j \in D_r} \vec{d}_j + \left[\frac{1}{|D_r|} \sum_{\vec{d}_j \in D_r} \vec{d}_j - \frac{1}{|D_{nr}|} \sum_{\vec{d}_j \in D_{nr}} \vec{d}_j \right]$$

- q-opt = centroid-rel - (centroid-rel - centroid-nonrel)
- We move the centroid of the relevant documents by the difference between the two centroids.
- We had to assume $|\vec{\mu}_r| = |\vec{\mu}_{nr}| = 1$ for this derivation.

Rocchio illustrated



circles: relevant documents, Xs: nonrelevant documents $\vec{\mu}_R$: centroid of relevant documents $\vec{\mu}_{NR}$ does not separate relevant/nonrelevant. $\vec{\mu}_R - \vec{\mu}_{NR}$: centroid of nonrelevant documents $\vec{\mu}_R - \vec{\mu}_{NR}$: difference vector Add difference vector to $\vec{\mu}_R \dots \dots$ to get \vec{q}_{opt} \vec{q}_{opt} separates relevant/nonrelevant perfectly.

29 / 57

30 / 57

Rocchio 1971 algorithm (SMART)

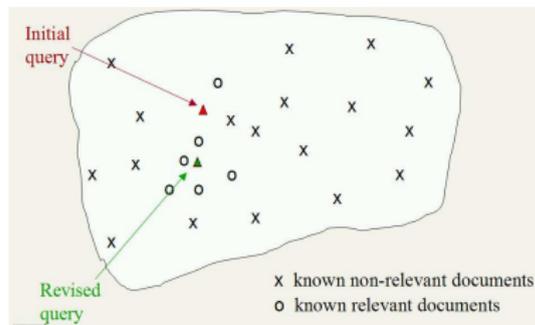
- Used in practice:

$$\vec{q}_m = \alpha \vec{q}_0 + \beta \frac{1}{|D_r|} \sum_{\vec{d}_j \in D_r} \vec{d}_j - \gamma \frac{1}{|D_{nr}|} \sum_{\vec{d}_j \in D_{nr}} \vec{d}_j$$

q_m : modified query vector; q_0 : original query vector; D_r and D_{nr} : sets of known relevant and nonrelevant documents respectively; α , β , and γ : weights attached to each term

- New query moves towards relevant documents and away from nonrelevant documents.
- Tradeoff α vs. β/γ : If we have a lot of judged documents, we want a higher β/γ .
- Set negative term weights to 0.
- "Negative weight" for a term doesn't make sense in the vector space model.

Rocchio relevance feedback illustrated



- Questions?

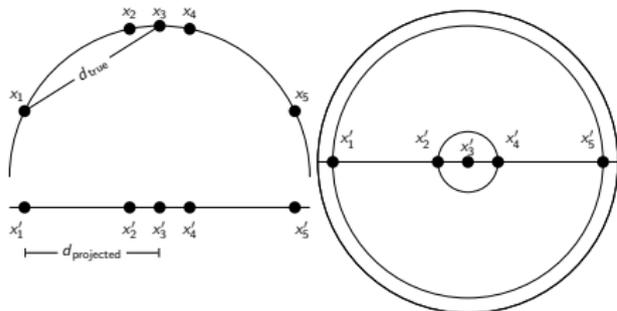
31 / 57

32 / 57

Positive vs. negative relevance feedback

- Positive feedback is more valuable than negative feedback.
- Why?
- For example, set $\beta = 0.75$, $\gamma = 0.25$ to give higher weight to positive feedback.
- Many systems only allow positive feedback.

Aside: 2D/3D graphs can be misleading



Left: A projection of the 2D semicircle to 1D. For the points x_1, x_2, x_3, x_4, x_5 at x coordinates $-0.9, -0.2, 0, 0.2, 0.9$ the distance $|x_2x_3| \approx 0.201$ only differs by 0.5% from $|x'_2x'_3| = 0.2$; but $|x_1x_3|/|x'_1x'_3| = d_{\text{true}}/d_{\text{projected}} \approx 1.06/0.9 \approx 1.18$ is an example of a large distortion (18%) when projecting a large area. *Right:* The corresponding projection of the 3D hemisphere to 2D.

33 / 57

34 / 57

Relevance feedback: Assumptions

- When can relevance feedback enhance recall?
- Assumption A1: The user knows the terms in the collection well enough for an initial query.
- Assumption A2: Relevant documents contain similar terms (so I can “hop” from one relevant document to a different one when giving relevance feedback).

Violation of A1

- Violation of assumption A1: The user knows the terms in the collection well enough for an initial query.
- Mismatch of searcher’s vocabulary and collection vocabulary
- Example: cosmonaut / astronaut

35 / 57

36 / 57

- Violation of A2: Relevant documents are not similar.
- Example query: contradictory government policies
- Why is relevance feedback unlikely to increase recall substantially for this query?
- Several unrelated “prototypes”
 - Subsidies for tobacco farmers vs. anti-smoking campaigns
 - Aid for developing countries vs. high tariffs on imports from developing countries
- Relevance feedback on tobacco docs will not help with finding docs on developing countries.

- Pick one of the evaluation measures from last lecture, e.g., precision in top 10: $P@10$
- Compute $P@10$ for original query q_0
- Compute $P@10$ for modified relevance feedback query q_1
- In most cases: q_1 is spectacularly better than q_0 !
- Is this a fair evaluation?

37 / 57

38 / 57

- Fair evaluation must be on “residual” collection: docs not yet judged by user.
- Studies have shown that relevance feedback is successful when evaluated this way.
- Empirically, one round of relevance feedback is often very useful. Two rounds are marginally useful.

- True evaluation of usefulness must compare to other methods taking the same amount of time.
- Alternative to relevance feedback: User revises and resubmits query.
- Users may prefer revision/resubmission to having to judge relevance of documents.
- There is no clear evidence that relevance feedback is the “best use” of the user’s time.

39 / 57

40 / 57

Do search engines use relevance feedback?

“similar pages” at Google

41 / 57

42 / 57

Relevance feedback: Problems

- Relevance feedback is expensive.
 - Relevance feedback creates long modified queries.
 - Long queries are expensive to process.
- Users are reluctant to provide explicit feedback.
- It's often hard to understand why a particular document was retrieved after applying relevance feedback.
- Excite had full relevance feedback at one point, but abandoned it later.

43 / 57

Other use of relevance feedback

- Maintaining a **standing query**
- Example: “multicore computer chips”
- I want to receive each morning a list of news articles published in the previous 24 hours on “multicore computer chips”.
- Relevance feedback can be used to refine this standing query over time.
- Many spam filters offer a similar functionality.
- For standing queries, relevance feedback is more practical than in web search.
- We'll revisit this issue in IIR 13.

44 / 57

Pseudo-relevance feedback

- Pseudo-relevance feedback automates the “manual” part of true relevance feedback.
- Pseudo-relevance algorithm:
 - Retrieve a ranked list of hits for the user’s query
 - Assume that the top k documents are relevant.
 - Do relevance feedback (e.g., Rocchio)
- Works very well on average
- But can go horribly wrong for some queries.
- Several iterations can cause *query drift*.
- Why?

45 / 57

Outline

- 1 Recap
- 2 Relevance feedback: Basics
- 3 Relevance feedback: Details
- 4 Global query expansion

47 / 57

Pseudo-relevance feedback at TREC4

- Cornell SMART system
- Results show number of relevant documents out of top 100 for 50 queries (so total number of documents is 5000):

method	number of relevant documents
Inc.Itc	3210
Inc.Itc-PsRF	3634
Lnu.Itu	3709
Lnu.Itu-PsRF	4350

- Results contrast two length normalization schemes (L vs. I) and pseudo-relevance feedback (PsRF).
- The pseudo-relevance feedback method used added only 20 terms to the query. (Rocchio will add many more.)
- This demonstrates that pseudo-relevance feedback is effective on average.

46 / 57

Global query expansion

- Query expansion is another method for [increasing recall](#).
- We use “global query expansion” to refer to “global methods for query reformulation”.
- In global query expansion, the query is modified based on some global resource, i.e. a resource that is not query-dependent.
- Main information we use: (near-)synonymy
- A publication or database that collects (near-)synonyms is called a [thesaurus](#).
- We will look at two types of thesauri: manually created and automatically created.

48 / 57

“Global” query expansion: Example

The screenshot shows a Yahoo! search results page for the query "palm". At the top, there are navigation links for Web, Images, Video, Audio, Directory, Local, News, Shopping, and More. Below the search bar, the search results are displayed. The main results section includes:

- Official Palm Store**: Free shipping on all handbells and more at the official Palm store.
- Palms Hotel - Best Rate Guarantee**: Book the Palms Hotel Casino with our best rate guarantee at VEGAS.com, the official Vegas travel site.
- Palm Pilots - Palm Downloads**: Yahoo! Shortcut - About
- Palm, Inc.**: Maker of handheld PDA devices that allow mobile users to manage schedules, contacts, and other personal and business information. Category: [ECS > Personal Digital Assistants \(PDAs\)](#). [www.palm.com](#) - 20k - [Cached](#) - [More from this site](#) - [Save](#)

Sponsor results include:

- Palm Memory**: Memory Giant is fast and easy. Guaranteed compatible memory. Great... [www.memorygiant.com](#)
- The Palms, Turks and Caicos Islands**: Resort/Condo photos, rates, availability and reservations... [www.worldwidereservationsystems.com](#)
- The Palms Casino Resort, Las Vegas**: Low price guarantee at the Palms Casino resort in Las Vegas. Book... [lasvegas.hotelscorp.com](#)

At the bottom of the page, it says "49 / 57".

Types of user feedback

- User gives feedback on **documents**.
 - More common in relevance feedback
- User gives feedback on **words** or **phrases**.
 - More common in query expansion
- Relevance feedback can also be thought of as a type of query expansion.
- We add terms to the query.
- The terms added in relevance feedback are based on “local” information in the result list.
- The terms added in query expansion are often based on “global” information that is not query-specific.

49 / 57

50 / 57

Types of query expansion

- Manual thesaurus (maintained by editors, e.g., PubMed)
- Automatically derived thesaurus (e.g., based on co-occurrence statistics)
- Query-equivalence based on query log mining (common on the web as in the “palm” example)

51 / 57

Thesaurus-based query expansion

- For each term t in the query, expand the query with words the thesaurus lists as semantically related with t .
- Example from earlier: HOSPITAL \rightarrow MEDICAL
- Generally increases recall
- May significantly decrease precision, particularly with ambiguous terms
 - INTEREST RATE \rightarrow INTEREST RATE FASCINATE EVALUATE
- Widely used in specialized search engines for science and engineering
- It's very expensive to create a manual thesaurus and to maintain it over time.
- A manual thesaurus is roughly equivalent to annotation with a *controlled vocabulary*.

52 / 57

- Attempt to generate a thesaurus automatically by analyzing the distribution of words in documents
- Fundamental notion: similarity between two words
- Definition 1: Two words are similar if they co-occur with similar words.
- Definition 2: Two words are similar if they occur in a given grammatical relation with the same words.
 - You can harvest, peel, eat, prepare, etc. apples and pears, so apples and pears must be similar.
- Co-occurrence is more robust, grammatical relations are more accurate.

53 / 57

54 / 57

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, parasitic
senses	grasp, psyche, truly, clumsy, naive, innate

- Relevance feedback and query expansion increase recall.
- In many cases, precision is decreased, often significantly.
- Log-based query modification (which is more complex than simple query expansion) is more common on the web than relevance feedback.

55 / 57

56 / 57

Resources

- Chapter 9 of IIR
- Resources at <http://ifnlp.org/ir>
- Salton and Buckley 1990 (original relevance feedback paper)
- Spink, Jansen, Ozmultu 2000: Relevance feedback at Excite
- Schütze 1998: Automatic word sense discrimination (describes a simple method for automatic thesaurus generation)