

Digital Libraries: Relevance Feedback and Query Expansion

RAFFAELLA BERNARDI

UNIVERSITÀ DEGLI STUDI DI TRENTO

P.ZZA VENEZIA, ROOM: 2.05, E-MAIL: BERNARDI@DISI.UNITN.IT

Contents

1	Improving Recall	4
1.1	Options for improving recall	5
2	Relevance feedback: Basics	6
2.1	Rocchio' algorithm	7
2.2	Relevance Feedback: Example 1	9
2.3	Results for initial query	11
2.4	User feedback: Select what is relevant	13
2.5	Results after relevance feedback	15
2.6	Example 2: A real (non-image) example	16
2.7	Expanded query after relevance feedback	17
2.8	Results for expanded query	18
2.9	Relevance feedback: Assumptions	19
2.10	Violation of A1	20
2.11	Violation of A2	21
2.12	Relevance feedback: Evaluation	22
2.13	Relevance feedback: Evaluation	23
2.14	Relevance feedback: Problems	24

2.15	Pseudo-relevance feedback	25
2.16	Pseudo-relevance feedback at TREC4	26
3	Global method: Query Expansion	27
3.1	Types of user feedback	28
3.2	Query expansion: Example	29
3.3	Types of query expansion	30
3.4	Thesaurus-based query expansion	31
3.5	Example for manual thesaurus: PubMed	32
3.6	Automatic thesaurus generation	33
3.7	Co-occurrence-based thesaurus: Examples	34
3.8	Query expansion at search engines	35
4	What do user wants?	36

1. Improving Recall

- Consider query q : [aircraft] ...
- ... and document d containing “plane”, but not containing “aircraft”
- A simple IR system will not return d for q .
- Even if d is the most relevant document for q !
- We want to return relevant documents even if there is no term match with the (original) query

1.1. Options for improving recall

Loose definition of recall in this lecture: “increasing the number of relevant documents returned to user”

- Local: Do a “local”, on-demand analysis for a user query
 - Main local method: *relevance feedback*
 - Part 1
- Global: Do a global analysis once (e.g., of collection) to produce *thesaurus*
 - Use thesaurus for query expansion
 - Part 2

2. Relevance feedback: Basics

Idea: You may not know what you are looking for, but you'll know when you see it.

- The user issues a (short, simple) query.
- The search engine returns a set of documents.
- User marks some docs as relevant, some as nonrelevant.
- Search engine computes a new representation of the information need. Hope: better than the initial query.
- Search engine runs new query and returns new results.
- New results have (hopefully) better recall.
- We can iterate this: several rounds of relevance feedback.

2.1. Rocchio' algorithm

We want to find a query vector that maximizes similarity with relevant documents (C_r) while minimizing similarity with nonrelevant documents (C_{nr}).

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

$$\vec{q}_{opt} = \operatorname{argmax}_{\vec{q}} [\operatorname{sim}(\vec{q}, C_r) - \operatorname{sim}(\vec{q}, C_{nr})]$$

- Intent: \vec{q}_{opt} is the vector that separates relevant and nonrelevant docs maximally.

the optimal query is the vector difference between the centroids of the relevant and non-relevant documents. (note we only have a partial knowledge of these two sets.)

2.2. Rocchio in a picture

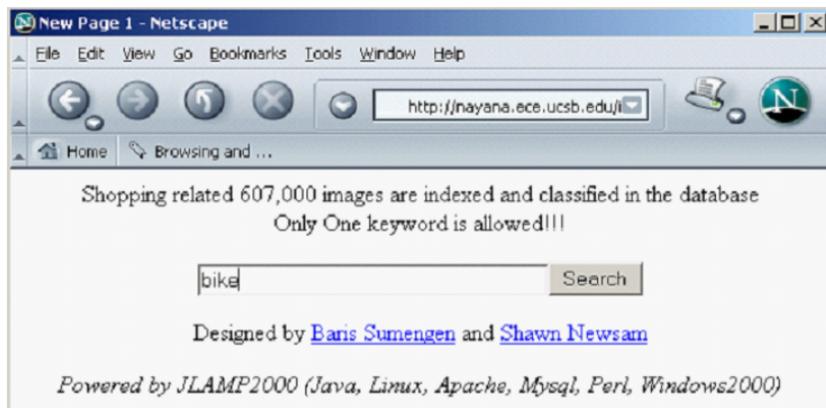
Rocchio in Pictures

query vector = $\alpha \cdot$ original query vector
+ $\beta \cdot$ positive feedback vector
- $\gamma \cdot$ negative feedback vector

Typically, $\gamma < \beta$

Original query	<table border="1"><tr><td>0</td><td>4</td><td>0</td><td>8</td><td>0</td><td>0</td></tr></table>	0	4	0	8	0	0	$\alpha = 1.0$	<table border="1"><tr><td>0</td><td>4</td><td>0</td><td>8</td><td>0</td><td>0</td></tr></table>	0	4	0	8	0	0
0	4	0	8	0	0										
0	4	0	8	0	0										
Positive Feedback	<table border="1"><tr><td>2</td><td>4</td><td>8</td><td>0</td><td>0</td><td>2</td></tr></table>	2	4	8	0	0	2	$\beta = 0.5$	<table border="1"><tr><td>1</td><td>2</td><td>4</td><td>0</td><td>0</td><td>1</td></tr></table> (+)	1	2	4	0	0	1
2	4	8	0	0	2										
1	2	4	0	0	1										
Negative feedback	<table border="1"><tr><td>8</td><td>0</td><td>4</td><td>4</td><td>0</td><td>16</td></tr></table>	8	0	4	4	0	16	$\gamma = 0.25$	<table border="1"><tr><td>2</td><td>0</td><td>1</td><td>1</td><td>0</td><td>4</td></tr></table> (-)	2	0	1	1	0	4
8	0	4	4	0	16										
2	0	1	1	0	4										
			<hr/>												
	New query		<table border="1"><tr><td>-1</td><td>6</td><td>3</td><td>7</td><td>0</td><td>-3</td></tr></table>	-1	6	3	7	0	-3						
-1	6	3	7	0	-3										

2.3. Relevance Feedback: Example 1



2.4. Results for initial query

Browse Search Prev Next Random

					
(144473, 16458) 0.0 0.0 0.0	(144457, 252140) 0.0 0.0 0.0	(144456, 262057) 0.0 0.0 0.0	(144456, 262063) 0.0 0.0 0.0	(144457, 252134) 0.0 0.0 0.0	(144483, 265154) 0.0 0.0 0.0
					
(144483, 264644) 0.0 0.0 0.0	(144483, 265153) 0.0 0.0 0.0	(144518, 257752) 0.0 0.0 0.0	(144538, 525937) 0.0 0.0 0.0	(144456, 249611) 0.0 0.0 0.0	(144436, 250064) 0.0 0.0 0.0

2.5. User feedback: Select what is relevant

Navigation buttons: [Browse](#) [Search](#) [Prev](#) [Next](#) [Random](#)

					
(144473, 16450)	(144457, 252140)	(144456, 262057)	(144456, 262063)	(144457, 252134)	(144403, 265154)
0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.0
					
(144483, 264644)	(144483, 265153)	(144518, 257752)	(144538, 525937)	(144456, 249611)	(144456, 250064)
0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.0

2.6. Results after relevance feedback

Browse Search Prev Next Random					
					
(144538, 528493) 0.54182 0.231944 0.309876	(144538, 528335) 0.56319296 0.267304 0.295869	(144538, 523529) 0.584279 0.280881 0.303598	(144456, 253569) 0.64301 0.351395 0.293615	(144456, 253568) 0.650275 0.411745 0.23853	(144538, 528799) 0.66709197 0.358033 0.309059
					
(144473, 16249) 0.6721 0.393922 0.278178	(144456, 249634) 0.675018 0.4639 0.211118	(144456, 253693) 0.676901 0.47645 0.200451	(144473, 16328) 0.700339 0.309002 0.391337	(144483, 265264) 0.70170796 0.36176 0.339948	(144478, 512410) 0.70297 0.469111 0.233859

2.7. Example 2: A real (non-image) example

Initial query: [new space satellite applications]

Results for initial query: (r = rank)

	r		
+	1	0.539	NASA Hasn't Scrapped Imaging Spectrometer
+	2	0.533	NASA Scratches Environment Gear From Satellite Plan
	3	0.528	Science Panel Backs NASA Satellite Plan, But Urges Launches of Smaller Probes
	4	0.526	A NASA Satellite Project Accomplishes Incredible Feat: Staying Within Budget
	5	0.525	Scientist Who Exposed Global Warming Proposes Satellites for Climate Research
	6	0.524	Report Provides Support for the Critics Of Using Big Satellites to Study Climate
	7	0.516	Arianespace Receives Satellite Launch Pact From Telesat Canada
+	8	0.509	Telecommunications Tale of Two Companies

User then marks relevant documents with “+”.

2.8. 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	arianespace
3.004	bundespost	2.806	ss
2.790	rocket	2.053	scientist
2.003	broadcast	1.172	earth
0.836	oil	0.646	measure

Compare to original query: [new space satellite applications]

2.9. Results for expanded query

	<i>r</i>	
*	1 0.513	NASA Scratches Environment Gear From Satellite Plan
*	2 0.500	NASA Hasn't Scrapped Imaging Spectrometer
	3 0.493	When the Pentagon Launches a Secret Satellite, Space Sleuths Do Some Spy Work of Their Own
	4 0.493	NASA Uses 'Warm' Superconductors For Fast Circuit
*	5 0.492	Telecommunications Tale of Two Companies
	6 0.491	Soviets May Adapt Parts of SS-20 Missile For Commercial Use
	7 0.490	Gaping Gap: Pentagon Lags in Race To Match the Soviets In Rocket Launchers
	8 0.490	Rescue of Satellite By Space Agency To Cost \$90 Million

* marks the documents which were judged as relevant.

2.10. 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).

2.11. Violation of A1

Assumption A1: The user knows the terms in the collection well enough for an initial query.

2.11. Violation of A1

Assumption A1: The user knows the terms in the collection well enough for an initial query.

- Violation: Mismatch of searcher's vocabulary and collection vocabulary
- Example: cosmonaut / astronaut

2.12. Violation of A2

Assumption A2: Relevant documents are not similar.

2.12. Violation of A2

Assumption A2: Relevant documents are not similar.

- Example for violation: [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.

2.13. Relevance feedback: Evaluation

- 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?

2.14. Relevance feedback: Evaluation

- 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.

2.14. Relevance feedback: Evaluation

- 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.

Evaluation: Caveat

- 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.

2.15. 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.
- The search engine Excite had full relevance feedback at one point, but abandoned it later.

2.16. 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*.

2.17. 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.ltc	3210
Inc.ltc-PsRF	3634
Lnu.ltu	3709
Lnu.ltu-PsRF	4350

- Results contrast two length normalization schemes (L vs. l) 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.

3. Global method: Query Expansion

- Query expansion is another method for increasing recall.
- We use “global query expansion” to refer to “global methods for query reformulation”.

3. Global method: 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.

3. Global method: 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.

3.1. 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.

3.2. Query expansion: Example

YAHOO! SEARCH

Web | [Images](#) | [Video](#) | [Audio](#) | [Directory](#) | [Local](#) | [News](#) | [Shopping](#) | [More »](#)

[Answers](#) | [My Web](#) | [Search Services](#) | [Advanced Search](#) | [Preferences](#)

Search Results 1 - 10 of about 160,000,000 for [palm](#) - 0.07 sec. ([About this page](#))

Also try: [palm springs](#), [palm pilot](#), [palm trees](#), [palm reading](#) [More...](#)

SPONSOR RESULTS

- [Official Palm Store](#)
[store.palm.com](#) Free shipping on all handhelds and more at the official **Palm** store.
- [Palms Hotel - Best Rate Guarantee](#)
[www.vegas.com](#) Book the **Palms** Hotel Casino with our best rate guarantee at VEGAS.com, the official Vegas travel site.

Y [Palm Pilots](#) - [Palm Downloads](#)
Yahoo! Shortcut - [About](#)

1. [Palm, Inc.](#) 
Maker of handheld PDA devices that allow mobile users to manage schedules, contacts, and other personal and business information.
Category: [B2B > Personal Digital Assistants \(PDAs\)](#)
[www.palm.com](#) - 20k - [Cached](#) - [More from this site](#) - [Save](#)

SPONSOR RESULTS

[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.c](#)

[The Palms Casino Resort, Las Vegas](#)
Low price guarantee at the **Palms** Casino resort in Las Vegas. Book...
[lasvegas.hotelscorp.com](#)

3.3. 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)

3.4. 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: hospital \rightarrow medical
- Generally increases recall
- May significantly decrease precision, particularly with ambiguous terms
 - interest rate \rightarrow interest rate fascinate
- 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 has an effect roughly equivalent to annotation with a controlled vocabulary.

3.5. Example for manual thesaurus: PubMed

The screenshot displays the PubMed search interface. At the top, the NCBI logo is on the left, and the PubMed logo and National Library of Medicine (NLM) logo are on the right. Below the logos is a navigation bar with tabs for PubMed, Nucleotide, Protein, Genome, Structure, PopSet, and Taxonomy. The search bar contains the text "Search PubMed" followed by a dropdown menu set to "PubMed", the word "for", and the search term "cancer". To the right of the search bar are "Go" and "Clear" buttons. Below the search bar are links for "Limits", "Preview/Index", "History", "Clipboard", and "Details".

On the left side, there is a vertical menu with the following items: "About Entrez", "Text Version", "Entrez PubMed", "Overview", "Help | FAQ", "Tutorial", "New/Noteworthy", "E-Utilities", "PubMed Services", "Journals Database", "MeSH Browser", "Single Citation", and "Multiple Citations".

The main content area shows the "PubMed Query:" section with the following query text:

```
("neoplasms"[MeSH Terms] OR cancer[Text Word])
```

At the bottom of the query area, there are "Search" and "URL" buttons.

3.6. Automatic thesaurus generation

- 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.

3.6. Automatic thesaurus generation

- 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.
 - “car” \approx “motorcycle” because both with “road”, “gas” and “license”, so they must be similar.
- 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.

3.7. Co-occurrence-based thesaurus: Examples

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

3.8. Query expansion at search engines

- Main source of query expansion at search engines: query logs
- Example 1: After issuing the query [herbs], users frequently search for [herbal remedies].
 - → “herbal remedies” is potential expansion of “herb”.
- Example 2: Users searching for [flower pix] frequently click on the URL [photobucket.com/flower](https://www.photobucket.com/flower). Users searching for [flower clipart] frequently click on the same URL.
 - → “flower clipart” and “flower pix” are potential expansions of each other.

4. What do user wants?

