

Introduction to Information Retrieval

CS276
Information Retrieval and Web Search
Pandu Nayak and Prabhakar Raghavan
Lecture 9: Query expansion

Reminder

- Midterm in class on Thursday 28th
- Material from first 8 lectures
- Open book, open notes
- You can use (and should bring!) a basic calculator
- You cannot use any wired or wireless communication. Use of such communication will be regarded as an Honor Code violation.
- You *can* preload the pdf of the book on to your laptop which you can use disconnected in the room.

Recap of the last lecture

- Evaluating a search engine
 - Benchmarks
 - Precision and recall
- Results summaries

Recap: Unranked retrieval evaluation: Precision and Recall

- **Precision**: fraction of retrieved docs that are relevant = $P(\text{relevant} | \text{retrieved})$
- **Recall**: fraction of relevant docs that are retrieved = $P(\text{retrieved} | \text{relevant})$

	Relevant	Nonrelevant
Retrieved	tp	fp
Not Retrieved	fn	tn

- Precision $P = \text{tp}/(\text{tp} + \text{fp})$
- Recall $R = \text{tp}/(\text{tp} + \text{fn})$

4

Recap: A combined measure: F

- Combined measure that assesses precision/recall tradeoff is **F measure** (weighted harmonic mean):

$$F = \frac{1}{\alpha \frac{1}{P} + (1-\alpha) \frac{1}{R}} = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$$

- People usually use balanced F_1 measure
 - i.e., with $\beta = 1$ or $\alpha = \frac{1}{2}$
- Harmonic mean is a conservative average
 - See CJ van Rijsbergen, *Information Retrieval*

5

This lecture

- Improving results
 - For high recall. E.g., searching for *aircraft* doesn't match with *plane*; nor *thermodynamic* with *heat*
- Options for improving results...
 - Global methods
 - Query expansion
 - Thesauri
 - Automatic thesaurus generation
 - Local methods
 - Relevance feedback
 - Pseudo relevance feedback

Introduction to Information Retrieval Sec. 9.1

Relevance Feedback

- Relevance feedback: user feedback on relevance of docs in initial set of results
 - User issues a (short, simple) query
 - The user marks some results as relevant or non-relevant.
 - The system computes a better representation of the information need based on feedback.
 - Relevance feedback can go through one or more iterations.
- Idea: it may be difficult to formulate a good query when you don't know the collection well, so iterate

Introduction to Information Retrieval Sec. 9.1

Relevance feedback

- We will use *ad hoc retrieval* to refer to regular retrieval without relevance feedback.
- We now look at four examples of relevance feedback that highlight different aspects.

Introduction to Information Retrieval Sec. 9.1.1

Similar pages

Google sarah brightman Search [Advanced Search](#) [Preferences](#)

Web Video Music

Sarah Brightman Official Website - Home Page
 Official site of world's best-selling soprano. Join FAN AREA free to access exclusive perks, photo diaries, a global forum community and more.
www.sarah-brightman.com/ - 4k - [Cached](#) [Similar pages](#)

Introduction to Information Retrieval Sec. 9.1.1

Relevance Feedback: Example

- Image search engine <http://nayana.ece.ucsb.edu/imsearch/imsearch.html>

Shopping related 607,000 images are indexed and classified in the database
 Only One keyword is allowed!!!

bike Search

Designed by [Baris Sumengen](#) and [Shawn Newsam](#)

Powered by [JLAMP2000](#) (Java, Linux, Apache, Mysql, Perl, Windows2000)

Introduction to Information Retrieval Sec. 9.1.1

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, 262857) 0.0 0.0 0.0	(144456, 262863) 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	(144510, 237752) 0.0 0.0 0.0	(144530, 52937) 0.0 0.0 0.0	(144456, 249611) 0.0 0.0 0.0	(144456, 250064) 0.0 0.0 0.0

Introduction to Information Retrieval Sec. 9.1.1

Relevance Feedback













Browse Search Prev Next Random

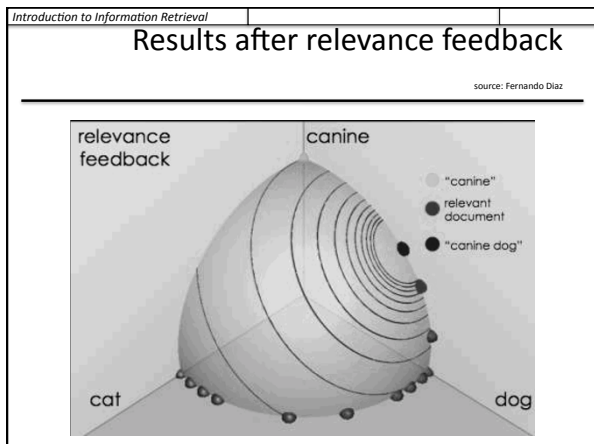
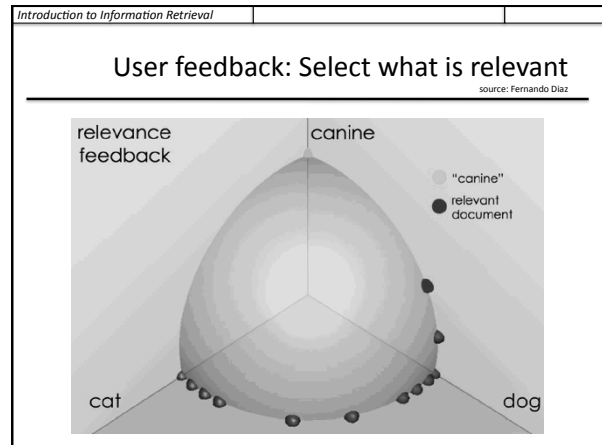
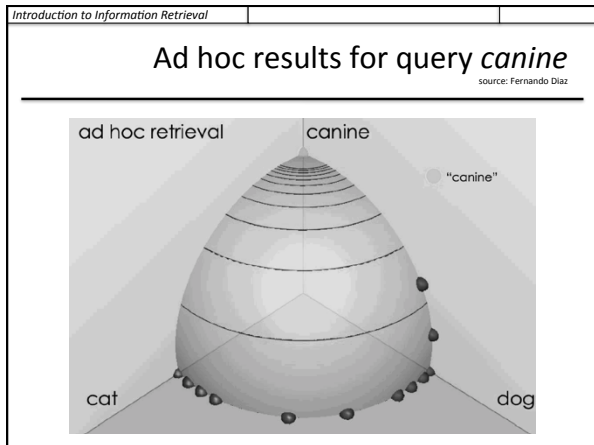
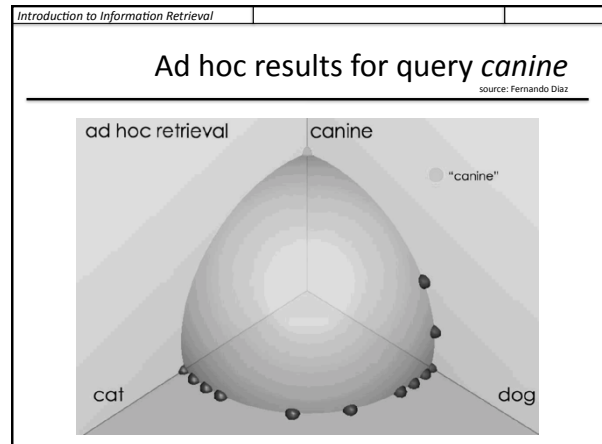
(144473, 16458) 0.0 0.0 0.0	(144457, 252140) 0.0 0.0 0.0	(144456, 262857) 0.0 0.0 0.0	(144456, 262863) 0.0 0.0 0.0	(144457, 252134) 0.0 0.0 0.0	(144483, 265154) 0.0 0.0 0.0
(144483, 264644) 0.0 0.0 0.0	(144483, 265153) 0.0 0.0 0.0	(144510, 237752) 0.0 0.0 0.0	(144530, 52937) 0.0 0.0 0.0	(144456, 249611) 0.0 0.0 0.0	(144456, 250064) 0.0 0.0 0.0

Introduction to Information Retrieval | Sec. 9.1.1

Results after Relevance Feedback

Browse Search Prev Next Random

					
(144538, 523493) 0.54182 0.231944 0.309876	(144538, 523835) 0.56319286 0.267304 0.293889	(144538, 523529) 0.584279 0.280881 0.303398	(14456, 253569) 0.64501	(14456, 253568) 0.630275 0.411745 0.23835	(144538, 523799) 0.66709197 0.338833 0.309059
					
(144473, 16249) 0.6721 0.393922 0.278178	(144456, 246534) 0.675018 0.4639 0.211118	(144456, 253693) 0.676001 0.47645 0.200451	(144473, 16320) 0.700339 0.309002 0.391337	(144483, 265264) 0.7017096 0.36176 0.339948	(144478, 512410) 0.70297 0.469111 0.233859



Introduction to Information Retrieval | Sec. 9.1.1

Initial query/results

- 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, 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 "+".

Introduction to Information Retrieval	Sec. 9.1.1
Expanded query after relevance feedback	
<ul style="list-style-type: none"> ▪ 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 	

Introduction to Information Retrieval	Sec. 9.1.1
Results for expanded query	
<ol style="list-style-type: none"> 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 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 8 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 	

Introduction to Information Retrieval	Sec. 9.1.1
Key concept: Centroid	
<ul style="list-style-type: none"> ▪ The <u>centroid</u> is the center of mass of a set of points ▪ Recall that we represent documents as points in a high-dimensional space ▪ Definition: Centroid 	
$\bar{\mu}(C) = \frac{1}{ C } \sum_{\vec{d} \in C} \vec{d}$	
where C is a set of documents.	

Introduction to Information Retrieval	Sec. 9.1.1
Rocchio Algorithm	
<ul style="list-style-type: none"> ▪ The Rocchio algorithm uses the vector space model to pick a relevance feedback query ▪ Rocchio seeks the query \vec{q}_{opt} that maximizes 	
$\vec{q}_{opt} = \arg \max_{\vec{q}} [\cos(\vec{q}, \bar{\mu}(C_r)) - \cos(\vec{q}, \bar{\mu}(C_{nr}))]$	
<ul style="list-style-type: none"> ▪ Tries to separate docs marked relevant and non-relevant 	
$\vec{q}_{opt} = \frac{1}{ C_r } \sum_{\vec{d}_j \in C_r} \vec{d}_j - \frac{1}{ C_{nr} } \sum_{\vec{d}_j \in C_{nr}} \vec{d}_j$	
<ul style="list-style-type: none"> ▪ Problem: we don't know the truly relevant docs 	

Introduction to Information Retrieval	Sec. 9.1.1
The Theoretically Best Query	
Optimal query	<ul style="list-style-type: none"> x non-relevant documents o relevant documents

Introduction to Information Retrieval	Sec. 9.1.1
Rocchio 1971 Algorithm (SMART)	
<ul style="list-style-type: none"> ▪ Used in practice: 	
$\vec{q}_m = \alpha \vec{q}_0 + \beta \frac{1}{ D_r } \sum_{\vec{d}_j \in D_r} \vec{d}_j - \gamma \frac{1}{ D_{nr} } \sum_{\vec{d}_j \in D_{nr}} \vec{d}_j$	
<ul style="list-style-type: none"> ▪ D_r = set of <u>known</u> relevant doc vectors ▪ D_{nr} = set of <u>known</u> irrelevant doc vectors <ul style="list-style-type: none"> ▪ Different from C_r and C_{nr} Δ ▪ q_m = modified query vector; q_0 = original query vector; α, β, γ: weights (hand-chosen or set empirically) ▪ New query moves toward relevant documents and away from irrelevant documents 	

Introduction to Information Retrieval | Sec. 9.1.1

Subtleties to note

- Tradeoff α vs. β/γ : If we have a lot of judged documents, we want a higher β/γ .
- Some weights in query vector can go negative
 - Negative term weights are ignored (set to 0)

Introduction to Information Retrieval | Sec. 9.1.1

Relevance feedback on initial query

Initial query

Revised query

x known non-relevant documents
o known relevant documents

Introduction to Information Retrieval | Sec. 9.1.1

Relevance Feedback in vector spaces

- We can modify the query based on relevance feedback and apply standard vector space model.
- Use only the docs that were marked.
- Relevance feedback can improve recall and precision
- Relevance feedback is most useful for increasing *recall* in situations where recall is important
 - Users can be expected to review results and to take time to iterate

Introduction to Information Retrieval | Sec. 9.1.1

Positive vs Negative Feedback

- Positive feedback is more valuable than negative feedback (so, set $\gamma < \beta$; e.g. $\gamma = 0.25$, $\beta = 0.75$).
- Many systems only allow positive feedback ($\gamma=0$).

Introduction to Information Retrieval

Aside: Vector Space can be Counterintuitive.

Doc

"J. Snow & Cholera"

Query

"cholera"

q1 query "cholera"

o www.ph.ucla.edu/epi/snow.html

x other documents

Introduction to Information Retrieval

High-dimensional Vector Spaces

- The queries "cholera" and "john snow" are far from each other in vector space.
- How can the document "John Snow and Cholera" be close to both of them?
- Our intuitions for 2- and 3-dimensional space don't work in $>10,000$ dimensions.
- 3 dimensions: If a document is close to many queries, then some of these queries must be close to each other.
- Doesn't hold for a high-dimensional space.

Introduction to Information Retrieval	Sec. 9.1.3
---------------------------------------	------------

Relevance Feedback: Assumptions

- A1: User has sufficient knowledge for initial query.
- A2: Relevance prototypes are "well-behaved".
 - Term distribution in relevant documents will be similar
 - Term distribution in non-relevant documents will be different from those in relevant documents
 - Either: All relevant documents are tightly clustered around a single prototype.
 - Or: There are different prototypes, but they have significant vocabulary overlap.
 - Similarities between relevant and irrelevant documents are small

Introduction to Information Retrieval	Sec. 9.1.3
---------------------------------------	------------

Violation of A1

- User does not have sufficient initial knowledge.
- Examples:
 - Misspellings (Brittany Speers).
 - Cross-language information retrieval (hígado).
 - Mismatch of searcher's vocabulary vs. collection vocabulary
 - Cosmonaut/astronaut

Introduction to Information Retrieval	Sec. 9.1.3
---------------------------------------	------------


Violation of A2

- There are several relevance prototypes.
- Examples:
 - Burma/Myanmar
 - Contradictory government policies
 - Pop stars that worked at Burger King
- Often: instances of a general concept
- Good editorial content can address problem
 - Report on contradictory government policies

Introduction to Information Retrieval	
---------------------------------------	--

Relevance Feedback: Problems

- Long queries are inefficient for typical IR engine.
 - Long response times for user.
 - High cost for retrieval system.
 - Partial solution:
 - Only reweight certain prominent terms
 - Perhaps top 20 by term frequency
- Users are often reluctant to provide explicit feedback
- It's often harder to understand why a particular document was retrieved after applying relevance feedback



Introduction to Information Retrieval	Sec. 9.1.5
---------------------------------------	------------

Evaluation of relevance feedback strategies

- Use q_0 and compute precision and recall graph
- Use q_m and compute precision recall graph
 - Assess on all documents in the collection
 - Spectacular improvements, but ... it's cheating!
 - Partly due to known relevant documents ranked higher
 - Must evaluate with respect to documents not seen by user
 - Use documents in residual collection (set of documents minus those assessed relevant)
 - Measures usually then lower than for original query
 - But a more realistic evaluation
 - Relative performance can be validly compared
- Empirically, one round of relevance feedback is often very useful. Two rounds is sometimes marginally useful.

Introduction to Information Retrieval	Sec. 9.1.5
---------------------------------------	------------

Evaluation of relevance feedback

- Second method – assess only the docs *not* rated by the user in the first round
 - Could make relevance feedback look worse than it really is
 - Can still assess relative performance of algorithms
- Most satisfactory – use two collections each with their own relevance assessments
 - q_0 and user feedback from first collection
 - q_m run on second collection and measured

Introduction to Information Retrieval Sec. 9.1.3

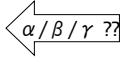
Evaluation: Caveat

- 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.

Introduction to Information Retrieval Sec. 9.1.4

Relevance Feedback on the Web

- Some search engines offer a similar/related pages feature (this is a trivial form of relevance feedback)
 - Google (link-based)
 - Altavista
 - Stanford WebBase
- But some don’t because it’s hard to explain to average user:
 - Alltheweb
 - Bing
 - Yahoo
- Excite initially had true relevance feedback, but abandoned it due to lack of use.



Introduction to Information Retrieval Sec. 9.1.4

Excite Relevance Feedback

Spink et al. 2000

- Only about 4% of query sessions from a user used relevance feedback option
 - Expressed as “More like this” link next to each result
- But about 70% of users only looked at first page of results and didn’t pursue things further
 - So 4% is about 1/8 of people extending search
- Relevance feedback improved results about 2/3 of the time

Introduction to Information Retrieval Sec. 9.1.6

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?

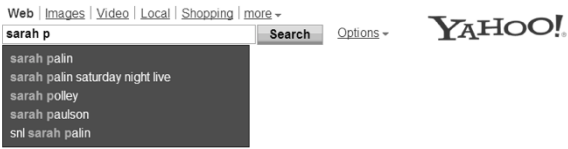
Introduction to Information Retrieval Sec. 9.2.2

Query Expansion

- In relevance feedback, users give additional input (relevant/non-relevant) on documents, which is used to reweight terms in the documents
- In query expansion, users give additional input (good/bad search term) on words or phrases

Introduction to Information Retrieval

Query assist



Web | Images | Video | Local | Shopping | more -

sarah p **YAHOO!**

sarah palin
sarah palin saturday night live
sarah polley
sarah paulson
snl sarah palin

Would you expect such a feature to increase the query volume at a search engine?

Introduction to Information Retrieval | Sec. 9.2.2

How do we augment the user query?

- Manual thesaurus
 - E.g. MedLine: physician, syn: doc, doctor, MD, medico
 - Can be query rather than just synonyms
- Global Analysis: (static; of all documents in collection)
 - Automatically derived thesaurus
 - (co-occurrence statistics)
 - Refinements based on query log mining
 - Common on the web
- Local Analysis: (dynamic)
 - Analysis of documents in result set

Introduction to Information Retrieval | Sec. 9.2.2

Example of manual thesaurus

Introduction to Information Retrieval | Sec. 9.2.2

Thesaurus-based query expansion

- For each term, t , in a query, expand the query with synonyms and related words of t from the thesaurus
 - feline → feline cat
- May weight added terms less than original query terms.
- Generally increases recall
- Widely used in many science/engineering fields
- May significantly decrease precision, particularly with ambiguous terms.
 - "interest rate" → "interest rate fascinate evaluate"
- There is a high cost of manually producing a thesaurus
 - And for updating it for scientific changes

Introduction to Information Retrieval | Sec. 9.2.3

Automatic Thesaurus Generation

- Attempt to generate a thesaurus automatically by analyzing the collection of 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 based is more robust, grammatical relations are more accurate. ← Why?

Introduction to Information Retrieval | Sec. 9.2.3

Co-occurrence Thesaurus

- Simplest way to compute one 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)

What does C contain if A is a term-doc incidence (0/1) matrix?

- For each t_i , pick terms with high values in C

Introduction to Information Retrieval | Sec. 9.2.3

Automatic Thesaurus Generation Example

word	ten nearest neighbors
absolutely	absurd whatsoever totally exactly nothing
bottomed	dip copper drops topped slide trimmed slig
captivating	shimmer stunningly superbly plucky witty
doghouse	dog porch crawling beside downstairs gaze
Makeup	repellent lotion glossy sunscreen Skin gel p
mediating	reconciliation negotiate cease conciliation p
keeping	hoping bring wiping could some would othe
lithographs	drawings Picasso Dali sculptures Gauguin
pathogens	toxins bacteria organisms bacterial parasite
senses	grasp psyche truly clumsy naive innate awl

Introduction to Information Retrieval	Sec. 9.2.3
<h2>Automatic Thesaurus Generation</h2> <h3>Discussion</h3> <hr/> <ul style="list-style-type: none"> ■ Quality of associations is usually a problem. ■ Term ambiguity may introduce irrelevant statistically correlated terms. <ul style="list-style-type: none"> ▪ "Apple computer" → "Apple red fruit computer" ■ Problems: <ul style="list-style-type: none"> ▪ False positives: Words deemed similar that are not ▪ False negatives: Words deemed dissimilar that are similar ■ Since terms are highly correlated anyway, expansion may not retrieve many additional documents. 	

Introduction to Information Retrieval	
<h2>Indirect relevance feedback</h2> <hr/> <ul style="list-style-type: none"> ■ On the web, DirectHit introduced a form of indirect relevance feedback. ■ DirectHit ranked documents higher that users look at more often. <ul style="list-style-type: none"> ▪ Clicked on links are assumed likely to be relevant <ul style="list-style-type: none"> ▪ Assuming the displayed summaries are good, etc. ■ Globally: Not necessarily user or query specific. <ul style="list-style-type: none"> ▪ This is the general area of <i>clickstream mining</i> ■ Today – handled as part of machine-learned ranking 	

Introduction to Information Retrieval	
<h2>Resources</h2> <hr/> <p>IIR Ch 9 MG Ch. 4.7 MIR Ch. 5.2 – 5.4</p>	