# Introduction to
# **Information Retrieval**

## Relevance Feedback

# 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

# 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.

# Similar pages

# Relevance Feedback: Example

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

# Results for Initial Query

# Relevance Feedback

# Results after Relevance Feedback

# Key concept: Centroid

- 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

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

where C is a set of documents.

# The Theoretically Best Query



Optimal query

x  non-relevant documents
o  relevant documents

# Rocchio 1971 Algorithm

- Used in practice:

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

- *$D_r$ = set of <u>known</u> relevant doc vectors*
- *$D_{nr}$ = set of <u>known</u> irrelevant doc vectors*
- *$q_m$ = modified query vector; $q_0$ = original query vector; $\alpha,\beta,\gamma$: weights (hand-chosen or set empirically)(typical values **a = 1**, **b = 0.8**, and **ϒ = 0.1**.*
- New query moves toward relevant documents and away from irrelevant documents

# Relevance feedback on initial query



Initial query

Revised query

x  known non-relevant documents
o  known relevant documents

# Relevance Feedback in vector spaces

- We can modify the query based on relevance feedback and apply standard vector space model.

- <span style="color:red">Use only the docs that were marked.</span>

- Relevance feedback can improve recall and precision

- <span style="color:red">Relevance feedback is most useful for increasing *recall* in situations where recall is important</span>

  - Users can be expected to review results and to take time to iterate

# 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 ($\gamma=0$).

# 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

# 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

# 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

# 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

# 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

α/β/γ ??

- 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.

# 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

# 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?