# Introduction to
# **Information Retrieval**

## Lecture 11: Relevance Feedback & Query Expansion - II

# Take-away today

- **Interactive relevance feedback:** improve initial retrieval results by telling the IR system which docs are relevant / nonrelevant

- Best known relevance feedback method: Rocchio feedback

- **Query expansion:** improve retrieval results by adding synonyms / related terms to the query

  - **Sources for related terms:** Manual thesauri, automatic thesauri, query logs

# Rocchio 1971 algorithm (SMART)

Used in practice:

$$\begin{aligned} \vec{q}_m &= \alpha\vec{q}_0 + \beta\mu(D_r) - \gamma\mu(D_{nr}) \\ &= \alpha\vec{q}_0 + \beta\frac{1}{|D_r|}\sum_{\vec{d}_j \in D_r}\vec{d}_j - \gamma\frac{1}{|D_{nr}|}\sum_{\vec{d}_j \in D_{nr}}\vec{d}_j \end{aligned}$$

$q_m$: modified query vector; $q_0$: original query vector; $D_r$ and $D_{nr}$ : sets of known relevant and nonrelevant documents respectively; $\alpha$, $\beta$, and $\gamma$: weights

- New query moves towards relevant documents and away from nonrelevant documents.
- Tradeoff $\alpha$ vs. $\beta/\gamma$: If we have a lot of judged documents, we want a higher $\beta/\gamma$.
- Set negative term weights to 0.
- "Negative weight" for a term doesn't make

3

# Positive vs. negative relevance feedback

- Positive feedback is more valuable than negative feedback.

- For example, set $\beta = 0.75$, $\gamma = 0.25$ to give higher weight to positive feedback.

- Many systems only allow positive feedback.

# 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

- 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

# Violation of A2

- Assumption A2: Relevant documents are similar.
- Example for violation: [contradictory government policies]
- 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.

# 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 q1
- In most cases: $q_1$ is spectacularly better than $q_0$!
- Is this a fair evaluation?

# Evaluation: Caveat

- True evaluation of usefulness <span style="color:red">must compare to other methods taking the same amount of time</span>.

- 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.

# Relevance feedback: Problems

- Relevance feedback is expensive.
  - Relevance feedback creates long modified queries.
  - Long queries are expensive to process.
- <span style="color:red">Users are reluctant to provide explicit feedback</span>.
- 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.

# 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*.

# 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 |
|---|---|
| lnc.ltc | 3210 |
| lnc.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)
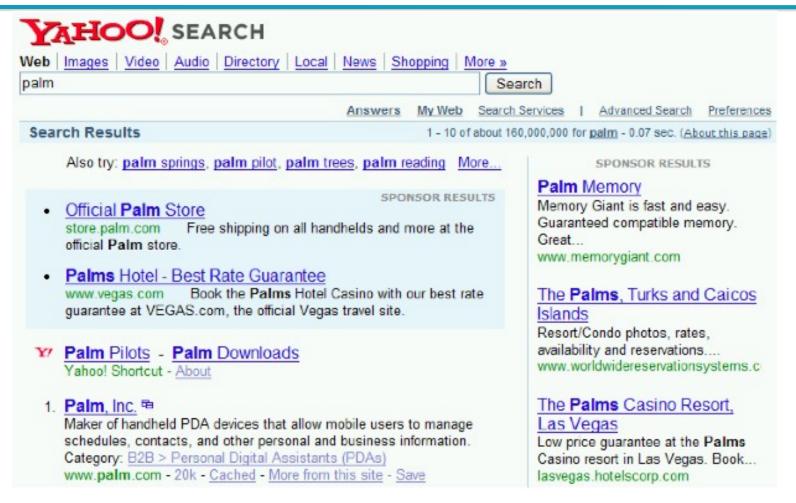- Demonstrates that pseudo-relevance feedback is effective on average

# Outline

❶ Motivation

❷ Relevance feedback: Basics

❸ Relevance feedback: Details

❹ Query expansion

# 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.

# Query expansion: Example

# 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

# Types of query expansion

- Manually constructed thesaurus (maintained by editors, e.g., Unified Medical Language System)

- Automatically derived thesaurus (e.g., based on co-occurrence statistics of terms)

- Query-equivalence based on query log mining (common on the web as in the "palm" example)

# 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 → MEDICAL
- Generally increases recall
- May significantly decrease precision, particularly with ambiguous terms: INTEREST RATE → INTEREST RATE FASCINATE

- Widely used in specialized search for science & 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

# 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" ≈ "motorcycle" because both occur 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"

# Co-occurence-based thesaurus construction

$$PMI(w_1, w_2) = log_2 \frac{P_{corpus}(w_1, w_2)}{P_{corpus}(w_1)P_{corpus}(w_2)}$$

$$P_{corpus}(w_1, w_2) = \frac{freq(w_1, w_2)}{N} \qquad P_{corpus}(w) = \frac{freq(w)}{N}$$

Statistically measure whether two words co-occur frequently (relative to their global frequencies)

20

# Co-occurence-based thesaurus: Examples

| petroleum | oil:0.032 gas:0.029 crude:0.029 barrels:0.028 exploration:0.027 barrel:0.026 opec:0.026 refining:0.026 gasoline:0.026 fuel:0.025 natural:0.025 exporting:0.025 |
|---|---|
| drug | trafficking:0.029 cocaine:0.028 narcotics:0.027 fda:0.026 police:0.026 abuse:0.026 marijuana:0.025 crime:0.025 colombian:0.025 arrested:0.025 addicts:0.024 |
| insurance | insurers:0.028 premiums:0.028 lloyds:0.026 reinsurance:0.026 underwriting:0.025 pension:0.025 mortgage:0.025 credit:0.025 investors:0.024 claims:0.024 benefits:0.024 |
| forest | timber:0.028 trees:0.027 land:0.027 forestry:0.026 environmental:0.026 species:0.026 wildlife:0.026 habitat:0.025 tree:0.025 mountain:0.025 river:0.025 lake:0.025 |
| robotics | robots:0.032 automation:0.029 technology:0.028 engineering:0.026 systems:0.026 sensors:0.025 welding:0.025 computer:0.025 manufacturing:0.025 automated:0.025 |

$$PMI(w_1, w_2) = log_2 \frac{P_{corpus}(w_1, w_2)}{P_{corpus}(w_1)P_{corpus}(w_2)}$$

$$P_{corpus}(w_1, w_2) = \frac{freq(w_1, w_2)}{N} \qquad P_{corpus}(w) = \frac{freq(w)}{N}$$

# Query Expansion: Examples

**TREC Topic 104:** *catastrophic health insurance*

**Query Representation:** surtax:1.0 hcfa:0.97 medicare:0.93 hmos:0.83 medicaid:0.8 hmo:0.78 beneficiaries:0.75 ambulatory:0.72 premiums:0.72 hospitalization:0.71 hhs:0.7 reimbursable:0.7 deductible:0.69

- Broad expansion terms: **medicare, beneficiaries, premiums** …

- Specific domain terms: **HCFA** (Health Care Financing Administration), **HMO** (Health Maintenance Organization), **HHS** (Health and Human Services)

**TREC Topic 355:** *ocean remote sensing*

**Query Representation:** radiometer:1.0 landsat:0.97 ionosphere:0.94 cnes:0.84 altimeter:0.83 nasda:0.81 meterology:0.81 cartography:0.78 geostationary:0.78 doppler:0.78 oceanographic:0.76

- Broad expansion terms: **radiometer, landsat, ionosphere** …

- Specific domain terms: **CNES** (Centre National dÉtudes Spatiales) and **NASDA** (National Space Development Agency of Japan)

# 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. Users searching for [flower clipart] frequently click on the same URL.
  - → "flower clipart" and "flower pix" are potential