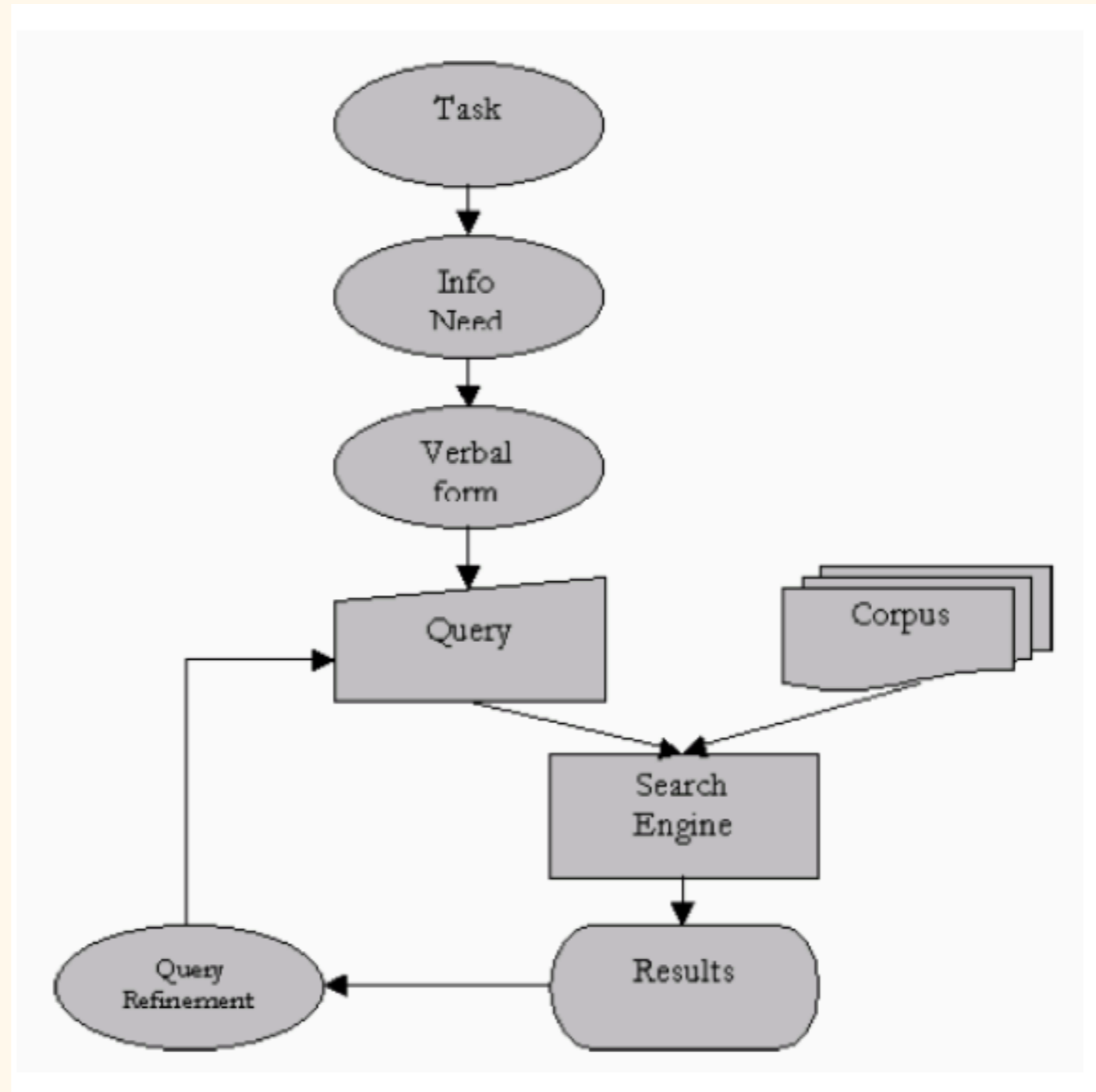# Relevance Feedback

- Concepts

- Global methods

  - Query expansion

    - Thesauri

- Local methods

  - Rocchio's algorithm

  - Probabilistic formulations

  - Implicit/Blind Relevance Feedback

- Papers

# Relevance Feedback



Recall that the information-seeking process is iterative...

Recall that the information-seeking process is iterative...

Why?

Users might:

Not know how to optimally word a query...

Not be sure what they're looking for, but will "know it when I see it"...

Not be sure what the collection contains...

All of these forces drive the need to iterate!

Remember: "recognition over recall"!

Classical relevance feedback flow:

1. User issues query;
2. System returns initial set of results;
3. User marks some number of documents as +/- relevant
4. System uses this data to compute a better representation of the user's information need
5. Repeat from step 2.

There are three major ways of doing relevance feedback:

## 1. Explicit Feedback

Users explicitly mark results as +/- relevant.

## 2. Implicit Feedback

The system attempts to infer +/- relevance from observable user behavior

## 3. Blind Feedback

The system attempts to infer +/- relevance blindly (no evidence)

Another way to conceptualize relevance feedback:

"Global" methods adjust a query independent of its results;

*Thesaurus-based query expansion, etc.*

"Local" methods adjust a query relative to its results.

*Explicit/pseudo relevance feedback, etc.*

(144473, 16458)
0.0
0.0
0.0

(144457, 252140)
0.0
0.0
0.0

(144456, 262857)
0.0
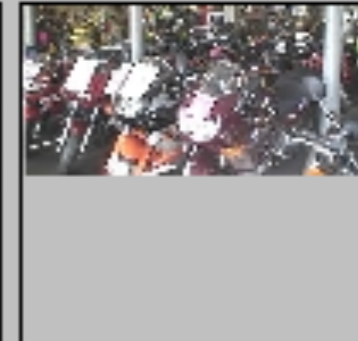0.0
0.0

(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

(144518, 257752)
0.0
0.0
0.0

(144538, 525937)
0.0
0.0
0.0

(144456, 249611)
0.0
0.0
0.0

(144456, 250064)
0.0
0.0
0.0

(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

(144518, 257752)
0.0
0.0
0.0

(144538, 525937)
0.0
0.0
0.0

(144456, 249611)
0.0
0.0
0.0

(144456, 250064)
0.0
0.0
0.0

(144538, 523493)
0.54182
0.231944
0.309876

(144538, 523835)
0.56319296
0.267304
0.295889

(144538, 523529)
0.584279
0.280881
0.303398

(144456, 253569)
0.64501
0.351395
0.293615

(144456, 253568)
0.650275
0.411745
0.23853

(144538, 523799)
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

Another way to conceptualize relevance feedback:

Let's hypothesize an optimal query to represent a given information need…

… the goal of relevance feedback is to adjust the user's query to more closely match the optimal query.

We can do this by adding terms or otherwise altering the user's query based on +/- relevant documents.

# "Global" methods adjust a query independent of its results.

# Query expansion: increase query effectiveness by adding new, hopefully relevant search terms

"Global" methods adjust a query independent of its results.

Query expansion: increase query effectiveness by adding new, hopefully relevant search terms

Supervised approaches:
  Thesauri, controlled vocabularies, etc.

Unsupervised approaches:
  Corpus-based association mining, vector embeddings, etc.

"Local" methods adjust a query relative to its results.

# The classical approach: Rocchio's algorithm

## Rocchio's algorithm operates on a vector space model.

The theory: the optimal query vector maximizes similarity with the relevant documents while minimizing similarity with non-relevant documents.

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

$$\vec{q}_{opt} = \frac{1}{|C_r|} \sum_{\vec{d}_j \in C_r} \vec{d}_j - \frac{1}{|C_{nr}|} \sum_{\vec{d}_j \in C_{nr}} \vec{d}_j$$

relevant document vectors                    non-relevant document vectors

Of course, we don't know the full set of relevant and non-relevant documents!

# The classical approach: Rocchio's algorithm

The theory: the optimal query vector maximizes similarity with the relevant documents while minimizing similarity with non-relevant documents.

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

$$\vec{q}_{opt} = \frac{1}{|C_r|} \sum_{\vec{d}_j \in C_r} \vec{d}_j - \frac{1}{|C_{nr}|} \sum_{\vec{d}_j \in C_{nr}} \vec{d}_j$$

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

Original query vector

small set of known-relevant docs

small set of known-non-relevant docs

$\alpha$, $\beta$, $\gamma$ are weights for how much to rely on the various components.

How should they be set?

# The classical approach: Rocchio's algorithm

The theory: the optimal query vector maximizes similarity with the relevant documents while minimizing similarity with non-relevant documents.

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

$$\vec{q}_{opt} = \frac{1}{|C_r|} \sum_{\vec{d}_j \in C_r} \vec{d}_j - \frac{1}{|C_{nr}|} \sum_{\vec{d}_j \in C_{nr}} \vec{d}_j$$

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

small set of known-relevant docs      small set of known-non-relevant docs

α, β, γ are weights for how much to rely on the various components.

When we have a lot of judged documents, β and γ should be large...

# The classical approach: Rocchio's algorithm

The theory: the optimal query vector maximizes similarity with the relevant documents while minimizing similarity with non-relevant documents.

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

$$\vec{q}_{opt} = \frac{1}{|C_r|} \sum_{\vec{d}_j \in C_r} \vec{d}_j - \frac{1}{|C_{nr}|} \sum_{\vec{d}_j \in C_{nr}} \vec{d}_j$$

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

small set of known-relevant docs          small set of known-non-relevant docs

$\alpha$, $\beta$, $\gamma$ are weights for how much to rely on the various components.

When we have fewer judged documents, $\alpha$ should dominate...

# The classical approach: Rocchio's algorithm

The theory: the optimal query vector maximizes similarity with the relevant documents while minimizing similarity with non-relevant documents.

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

$$\vec{q}_{opt} = \frac{1}{|C_r|} \sum_{\vec{d}_j \in C_r} \vec{d}_j - \frac{1}{|C_{nr}|} \sum_{\vec{d}_j \in C_{nr}} \vec{d}_j$$

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

small set of known-relevant docs          small set of known-non-relevant docs

$\alpha$, $\beta$, $\gamma$ are weights for how much to rely on the various components.

In any event, positive feedback is almost always more useful, so $\beta > \gamma$.

# The classical approach: Rocchio's algorithm

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

| 0 | 4 | 0 | 8 | 0 | 0 |
|---|---|---|---|---|---|

Original query vector

| 2 | 4 | 8 | 0 | 0 | 2 |
|---|---|---|---|---|---|

Known-relevant centroid

| 8 | 0 | 4 | 4 | 0 | 16 |
|---|---|---|---|---|---|

Known-non-relevant centroid

# The classical approach: Rocchio's algorithm

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

$\alpha = 1$

$\beta = 0.5$

$\gamma = 0.25$

| 0 | 4 | 0 | 8 | 0 | 0 |
|---|---|---|---|---|---|

$+$

| 1 | 2 | 4 | 0 | 0 | 1 |
|---|---|---|---|---|---|

$-$

| 2 | 0 | 1 | 1 | 0 | 4 |
|---|---|---|---|---|---|

| -1 | 6 | 3 | 7 | 0 | -3 |
|----|---|---|---|---|----|

| 0 | 6 | 3 | 7 | 0 | 0 |
|---|---|---|---|---|---|

Under Rocchio's formulation, negative weights get zeroed out.

# Another formulation: Probabilistic approaches

A simple naïve Bayesian model:

$$\hat{P}(x_t = 1 | R = 1) \quad = \quad |VR_t| / |VR|$$
$$\hat{P}(x_t = 1 | R = 0) \quad = \quad (df_t - |VR_t|) / (N - |VR|)$$

# known-relevant documents with $t$     # known-relevant documents

These probabilities can be used to re-weight any other probabilistic retrieval model (e.g., the Lee & Croft paper).

Explicit relevance feedback imposes work on users...

Automated local analysis methods can help!

Pseudo-relevance feedback

Implicit feedback

Local clustering

Pseudo-relevance feedback is the simplest form of automated local analysis.

The idea: Do retrieval as normal;

Blindly take the top $k$ returned documents, assume relevance;

Proceed with relevance feedback as normal.

Amazingly, this (kind of) works!

Problems: susceptible to topic drift, etc.

Implicit feedback: use click-through data, etc. to determine relevance.

(discussed previously, and again in a future lecture)

Local clustering:

Identify possible query expansion terms by doing associational clustering within the retrieved result set.

Can use local context to generate more useful expansions.

Can operate at the document level, or even at the passage level within a document.

# Manning and Schütze make a very good point:

"Relevance feedback can improve both recall and precision. But, in practice, it has been shown to be most useful for increasing recall in situations where recall is important. This is partly because the technique expands the query, but it is also partly an effect of the use case: when they want high recall, users can be expected to take time to review results and to iterate on the search."

How should we evaluate relevance feedback?

One idea: look at MAP before and after feedback, on the same set of documents.

What's wrong with this picture?

It's cheating! For the second query, we know *a priori* some documents to include!

How should we evaluate relevance feedback?

Second idea: Evaluate on the *residual collection* (the set of documents excluding those judged by the user)

What's wrong with this picture?

We get artificially degraded performance, since we exclude some of the most relevant documents.

If our goal is to compare two different RF approaches, this may not be an issue.

How should we evaluate relevance feedback?

Third idea: Use multiple sets of documents ("train"/"test"), examine post-RF query on both.

Fourth idea: Extrinsic evaluation!

Does RF help the user complete a task more quickly? Do they identify more relevant documents using RF? Etc. etc.

Problems with Relevance Feedback:

RF does not help with:

    Mis-spellings

    Mismatch of searcher's vocabulary to collection's ("laptop" vs. "notebook")

    Documents not clustering "naturally"

    Relevance not being related to term distribution of documents

Also: users often don't like giving explicit feedback!

# A deterministic resampling method using overlapping document clusters for pseudo-relevance feedback

Kyung Soon Lee [a,*], W. Bruce Croft [b]

[a] Division of Computer Science and Engineering, CAIIT, Chonbuk National University, 567 Baekje-daero, Deokjin-gu, Jeonju, Jeollabuk-do 561-756, Republic of Korea
[b] Center for Intelligent Information Retrieval, Department of Computer Science, University of Massachusetts Amherst, 140 Governors Drive, Amherst, MA 01003-9264, USA

## ARTICLE INFO

## ABSTRACT

Typical pseudo-relevance feedback methods assume the top-retrieved documents are relevant and use these pseudo-relevant documents to expand terms. The initial retrieval set can, however, contain a great deal of noise. In this paper, we present a cluster-based resampling method to select novel pseudo-relevant documents based on Lavrenko's relevance model approach. The main idea is to use overlapping clusters to find dominant documents for the initial retrieval set, and to repeatedly use these documents to emphasize the core topics of a query.

The proposed resampling method can skip some documents in the initial high-ranked documents and deterministically construct overlapping clusters as sampling units. The hypothesis behind using overlapping clusters is that a good representative document for a query may have several nearest neighbors with high similarities, participating in several different clusters. Experimental results on large-scale web TREC collections show significant improvements over the baseline relevance model.

To justify the proposed approach, we examine the relevance density and redundancy ratio of feedback documents. A higher relevance density will result in greater retrieval accuracy, ultimately approaching true relevance feedback. The resampling approach shows higher relevance density than the baseline relevance model on all collections, resulting in better retrieval accuracy in pseudo-relevance feedback.

© 2013 Elsevier Ltd. All rights reserved.

## 1. Introduction

Most pseudo-relevance feedback methods (e.g., Attar & Fraenkel, 1977; Buckley, Salton, Allan, & Singhal, 1995; Croft & Harper, 1979; Lavrenko & Croft, 2001; Robertson, Walker, Beaulieu, Gatford, & Payne, 1996) assume that a set of top-retrieved documents is relevant and then learn from the pseudo-relevant documents to expand terms or to assign better weights to the original query. This is similar to the process used in relevance feedback, when actual relevant documents are used (Salton & Buckley, 1990). In general, however, the top retrieved documents contain noise: when the precision of the top 10 documents (P@10) is 0.5, this means that five of them are non-relevant. This is common and even expected in all retrieval models. When combined with pseudo-relevance feedback, this noise, however, can cause the query representation to "drift" away from the original query.

This paper describes *a deterministic sampling method based on overlapping clusters* to select better documents for pseudo-relevance feedback. The sampling unit is a document cluster from the initial retrieval set which can represent an aspect of a

The basic idea:

Pseudo-relevance feedback (PRF) is easy to do...

... but obviously depends heavily on the relevance of the first several documents.

Lee & Croft's approach: be more clever about selecting which documents to use for PRF.

The main idea:

Cluster result documents according to similarity;

Identify "dominant" documents (those that appear in multiple clusters)...

... use those to select query expansion terms.

The initial set of results is produced using standard language model-based retrieval:

$$P(Q|D) = \prod_{i=1}^{m} P(q_i|D)$$

The initial set of results is produced using standard language model-based retrieval.

The results are *k*-means clustered, and then each cluster is sampled...

$$P(Q|Clu) = \prod_{i=1}^{m} P(q_i|Clu)$$

| Collection | LM | Rerank | RM | Resampling | TrueRF |
|---|---|---|---|---|---|
| GOV2 | 0.3258 | $0.3406^{\alpha}$ | $0.3581^{\alpha\beta}$ | $\mathbf{0.3806}^{\alpha\beta\gamma}$ | $0.4315^{\alpha\beta\gamma\delta}$ |
| WT10g | 0.1861 | $0.2044^{\alpha}$ | 0.1966 | $\mathbf{0.2352}^{\alpha\beta\gamma}$ | $0.4030^{\alpha\beta\gamma\delta}$ |
| ROBUST | 0.2920 | $0.3206^{\alpha}$ | $\mathbf{0.3591}^{\alpha\beta}$ | $0.3515^{\alpha\beta}$ | $0.5351^{\alpha\beta\gamma\delta}$ |
| AP | 0.2077 | $0.2361^{\alpha}$ | $0.2803^{\alpha\beta}$ | $\mathbf{0.2906}^{\alpha\beta}$ | $0.4253^{\alpha\beta\gamma\delta}$ |
| WSJ | 0.3258 | $0.3611^{\alpha}$ | $0.3967^{\alpha\beta}$ | $\mathbf{0.4033}^{\alpha\beta}$ | $0.5306^{\alpha\beta\gamma\delta}$ |