

Relevance Feedback and Query Expansion

Berlin Chen

Department of Computer Science & Information Engineering
National Taiwan Normal University

References:

1. *Modern Information Retrieval*, Chapter 5 & Teaching material
2. *Introduction to Information Retrieval*, Chapter 9

Introduction

- Users have no detailed knowledge of
 - The collection makeup
 - The retrieval environment

} Difficult to formulate queries
- Yet, most users often need to reformulate their queries to obtain the results of their interest
 - Thus, the first query formulation should be treated as an initial attempt to retrieve relevant information
 - Documents initially retrieved could be analyzed for relevance and used to improve the initial query

Introduction (cont.)

- The process of query modification is commonly referred as
 - **Relevance feedback**, when the user provides information on relevant documents to a query, or
 - **Query expansion**, when information related to the query is used to expand it
- We refer to both of them as feedback methods

Note also that, in most collections, the same concept may be referred to using different words

- This issue, known as synonymy, has an impact on the recall of most IR systems

Introduction (cont.)

- Two basic approaches of feedback methods:
 - **Explicit feedback**, in which the information for query reformulation is provided directly by the users, and
 - **Implicit feedback**, in which the information for query reformulation is implicitly derived by the system

Typical Framework for Feedback Methods

- Consider a set of documents D_r that are known to be relevant to the current query q
- In relevance feedback, the documents in D_r are used to transform q into a modified query q_m
- However, obtaining information on documents relevant to a query requires the direct interference of the user
 - Most users are unwilling to provide this information, particularly in the Web

Typical Framework for Feedback Methods (cont.)

- Because of this high cost, the idea of relevance feedback has been relaxed over the years
- Instead of asking the users for the relevant documents, we could:
 - Look at documents they **have clicked on**; or
 - Look at terms **belonging to the top documents in the result set**
- In both cases, it is expect that the feedback cycle will produce results of higher quality

Typical Framework for Feedback Methods (cont.)

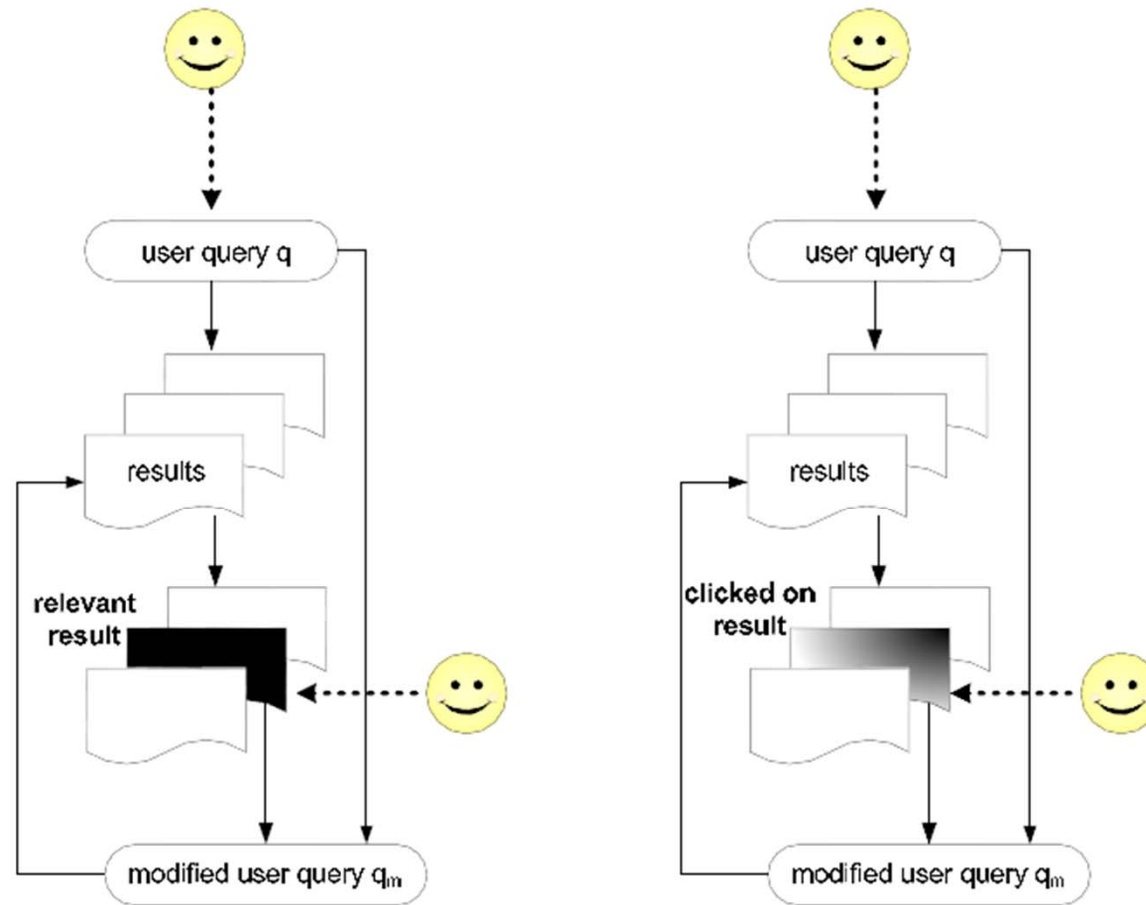
- A **feedback cycle** is composed of two basic steps:
 - Determine feedback information that is either related or expected to be related to the original query q and
 - Determine how to transform query q to take this information effectively into account
- The first step can be accomplished in two distinct ways:
 - Obtain the feedback information explicitly from the users
 - Obtain the feedback information implicitly from the query results or from external sources such as a **thesaurus**

Typical Framework for Feedback Methods (cont.)

- In an **explicit relevance feedback** cycle, the feedback information is provided directly by the users
- However, collecting feedback information is expensive and time consuming
- In the Web, **user clicks** on search results constitute a new source of feedback information
- A click indicate a document that is of interest to the user in the context of the current query
 - Notice that a click does not necessarily indicate a document that is relevant to the query

Typical Framework for Feedback Methods (cont.)

- Explicit Feedback Information



(a) relevance feedback

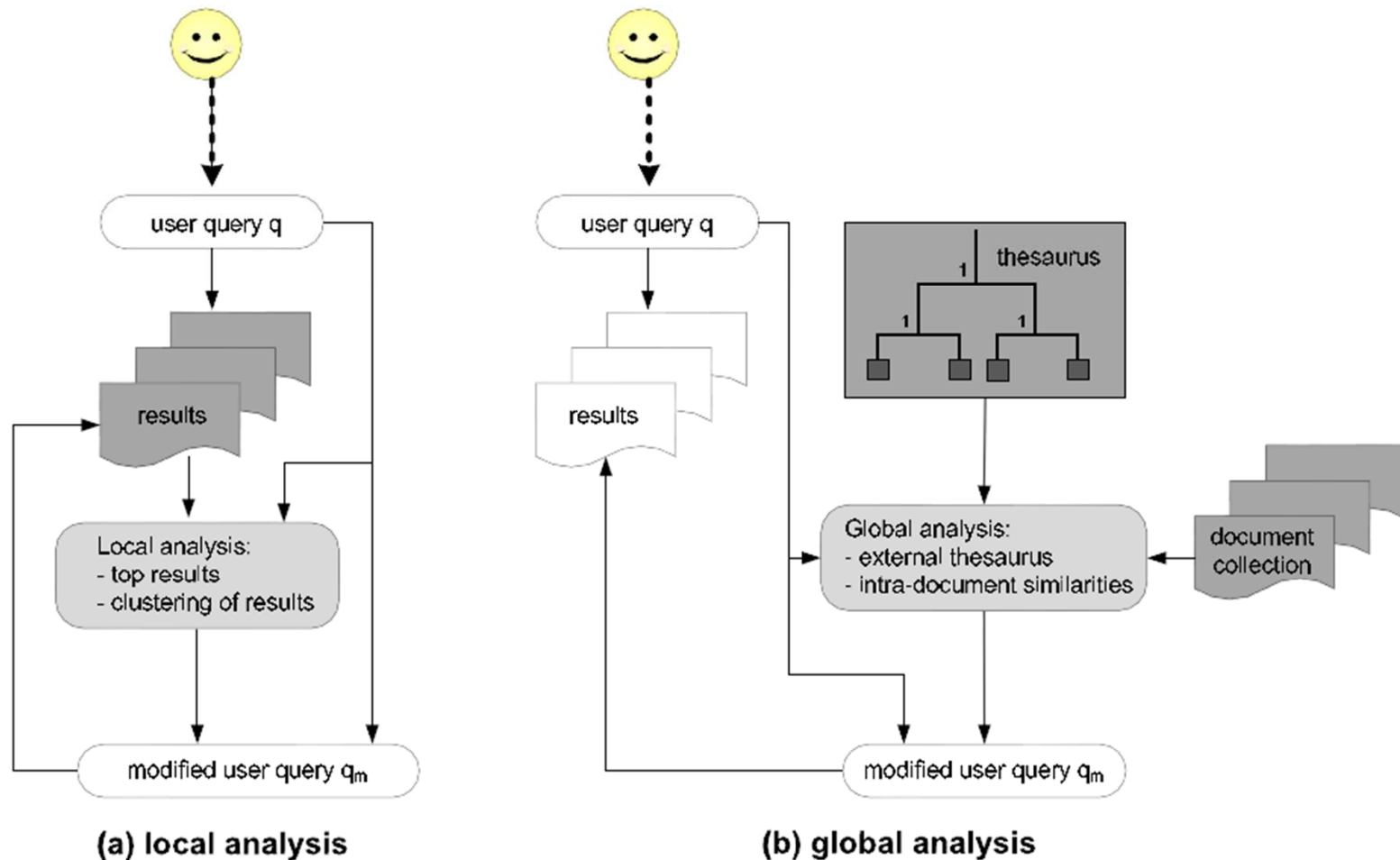
(b) click feedback

Typical Framework for Feedback Methods (cont.)

- In an **implicit relevance feedback** cycle, the feedback information is derived implicitly by the system
- There are two basic approaches for compiling implicit feedback information:
 - **Local analysis**, which derives the feedback information from the top ranked documents in the result set
 - **Global analysis**, which derives the feedback information from external sources such as a thesaurus

Typical Framework for Feedback Methods (cont.)

- Implicit Feedback Information



Typical Framework for Feedback Methods (cont.)

- Implicit Feedback Information
 - Obviously, the feedback information is not necessarily related to the current query, which makes its utilization more challenging than information provided explicitly by the users
 - Despite that, since implicit information is abundant and can be gathered at very low cost, there has been a persistent interest in using implicit information to improve query results.

Summary of Feedback Methods

- Feedback information from the user
 - **Relevance feedback**
 - With vector, probabilistic models et al.
- Information derived from the set of documents initially retrieved (called local set of documents)
 - **Local analysis**
 - Local clustering, local context analysis
- Global information derived from document collection
 - **Global analysis**
 - Similar thesaurus or statistical thesaurus

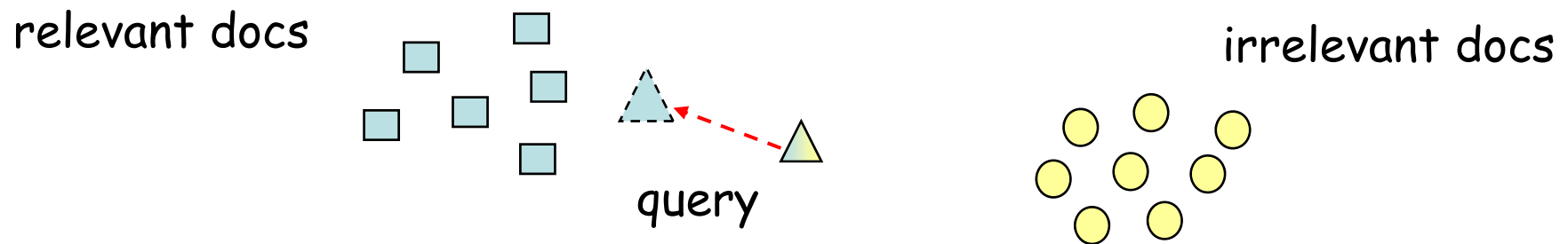
Explicit Relevance Feedback

Explicit Relevance Feedback

- In a classic relevance feedback cycle, the user is presented with a list of the retrieved documents
- Then, the user examines them and marks those that are relevant
- In practice, only the **top 10 (or 20) ranked documents** need to be examined
- The main idea consists of
 - Selecting important terms from the documents that have been identified as relevant, and
 - Enhancing the importance of these terms in a new query formulation

Explicit Relevance Feedback (cont.)

- **Expected effect:** the new query will be moved towards the relevant docs and away from the non-relevant ones
- Early experiments have shown good improvements in precision for small test collections



Explicit Relevance Feedback (cont.)

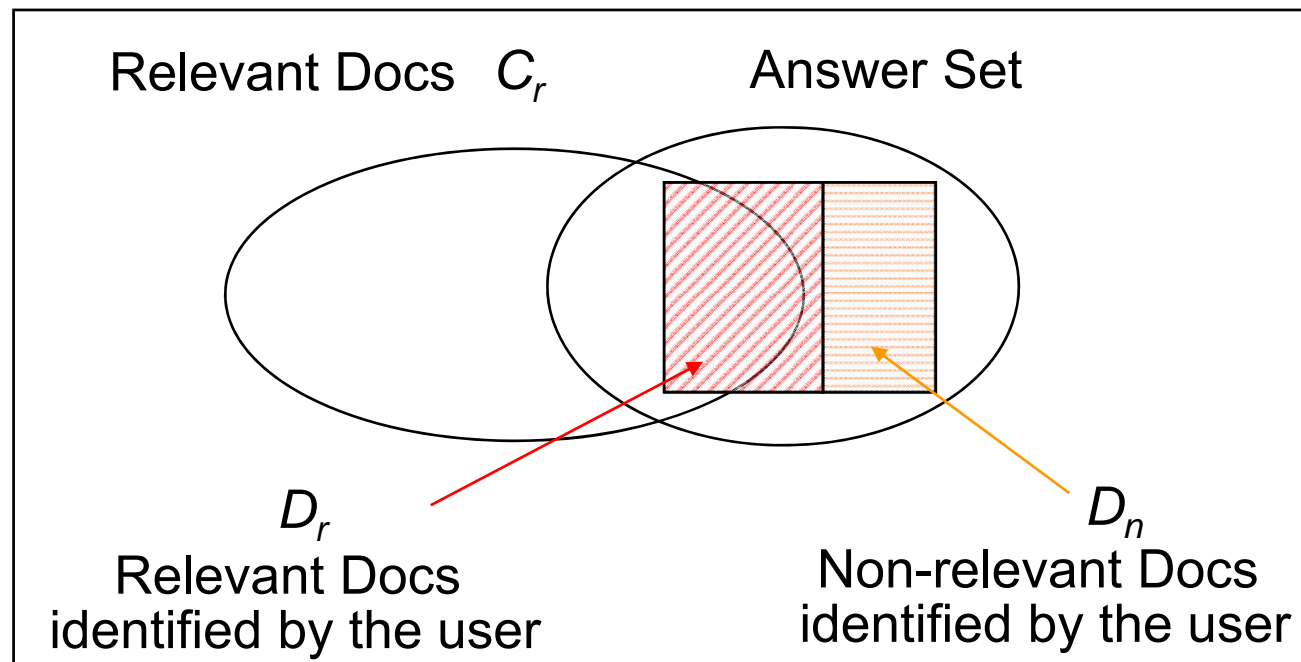
- Relevance feedback presents the following characteristics:
 - It shields the user from the details of the query reformulation process (**all the user has to provide is a relevance judgment**)
 - It breaks down the whole searching task into a sequence of small steps which are easier to grasp
 - Provide a controlled process designed to emphasize some terms (relevant ones) and de-emphasize others (non-relevant ones)

Vector Model: The Rocchio Method

- Premises
 - Documents identified as relevant (to a given query) have similarities among themselves
 - Further, non-relevant docs have term-weight vectors which are dissimilar from the relevant documents
 - The basic idea of the Rocchio Method is to reformulate the query such that it gets:
 - Closer to the neighborhood of the relevant documents in the vector space, and
 - Away from the neighborhood of the non-relevant documents

Vector Model: The Rocchio Method (cont.)

- Terminology



Doc Collection with size N

Vector Model: The Rocchio Method (cont.)

- **Optimal Condition**

- The complete set of relevant docs C_r to a given query q is known in advance

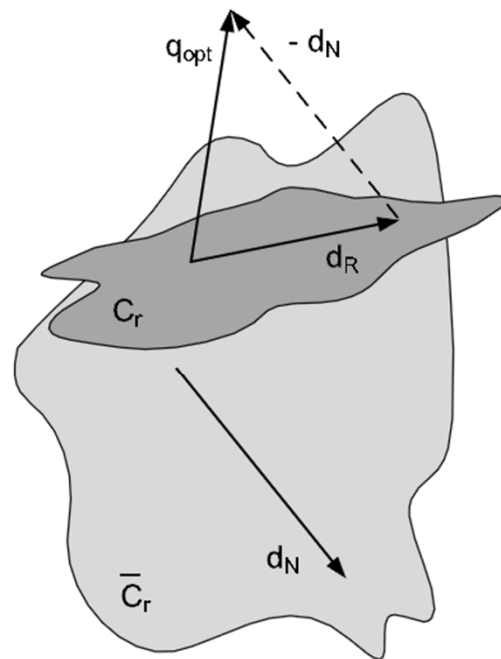
$$\vec{q}_{opt} = \frac{1}{|C_r|} \sum_{\forall \vec{d}_i \in C_r} \vec{d}_i - \frac{1}{N - |C_r|} \sum_{\forall \vec{d}_j \notin C_r} \vec{d}_j$$

Elements in the final vector representation should be kept nonnegative (to be in the positive quadrant of the vector space)

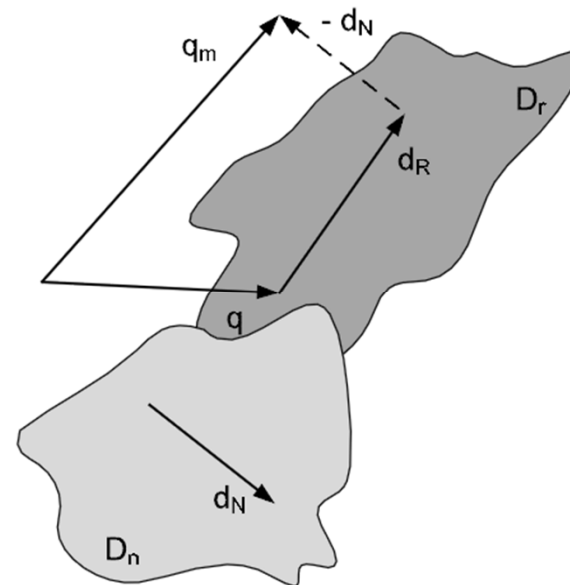
- **Problem:** the complete set of relevant docs C_r are not known a priori

Vector Model: The Rocchio Method (cont.)

- **Solution:** we can formulate an initial query and to incrementally change the initial query vector



(a)



(b)

Vector Model: The Rocchio Method (cont.)

- **In Practice:** There are three classic and similar ways to calculate the modified query

Rocchio 1965

1. Standard_Rocchio

$$\vec{q}_m = \alpha \cdot \vec{q} + \frac{\beta}{|D_r|} \cdot \sum_{\forall \vec{d}_i \in D_r} \vec{d}_i - \frac{\gamma}{|D_n|} \cdot \sum_{\forall \vec{d}_j \in D_n} \vec{d}_j$$

modified query \vec{q}_m initial/original query \vec{q}

2. Ide_Regular

$$\vec{q}_m = \alpha \cdot \vec{q} + \beta \cdot \sum_{\forall \vec{d}_i \in D_r} \vec{d}_i - \gamma \cdot \sum_{\forall \vec{d}_j \in D_n} \vec{d}_j$$

Elements in the final vector representation should be kept nonnegative (to be in the positive quadrant of the vector space)

3. Ide_Dec_Hi

$$\vec{q}_m = \alpha \cdot \vec{q} + \beta \cdot \sum_{\forall \vec{d}_i \in D_r} \vec{d}_i - \gamma \cdot \max_{non-relevant} (\vec{d}_j)$$

The highest ranked non-relevant doc

Positive feedback turns out to be much more valuable than negative feedback.

Vector Model: The Rocchio Method (cont.)

- **Some Observations**

- Similar results were achieved for the above three approach (Dec-Hi slightly better in the past)
- Usually, constant β is bigger than γ (why?)

- **In Practice (cont.)**

- More about the constants
 - Rocchio, 1971: $\alpha=1$
 - Ide, 1971: $\alpha=\beta=\gamma=1$
 - **Positive feedback strategy: $\gamma=0$**

In implementation, terms occurring in the relevant or non-relevant documents can be used in toto or selectively to reweight/argument or be moved from the initial query.

More on Explicit Relevance Feedback

- Advantages

- Simple, good results

- Modified term weights are computed directly from the retrieved docs

- Disadvantages

- No optimality criterion

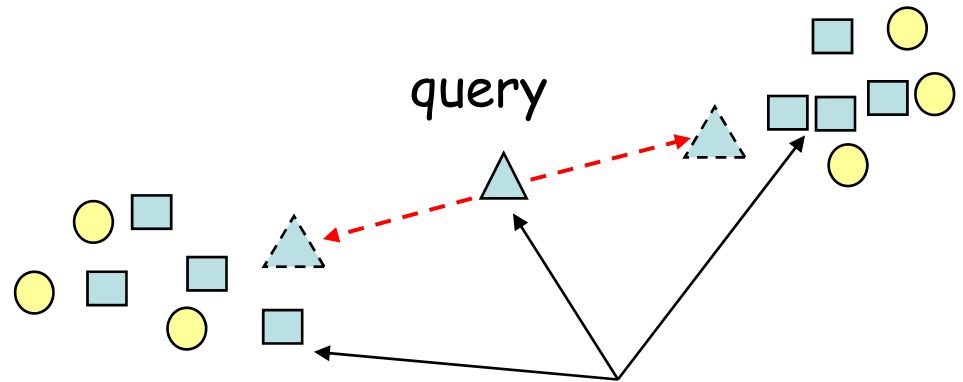
- Empirical and heuristic

(what if relevant documents belong to multiple clusters?)

- High computing cost (potentially long response time)

- Only reweight certain prominent terms in relevant docs?

- There are still cases where relevance feedback alone is not sufficient: e.g., misspellings, cross-language IR, mismatch of searcher's versus collection vocabularies

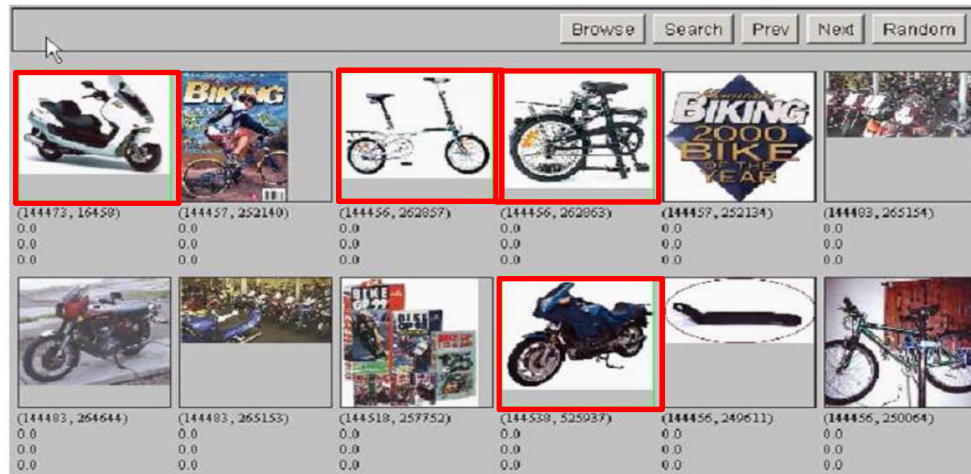


More on Explicit Relevance Feedback (cont.)

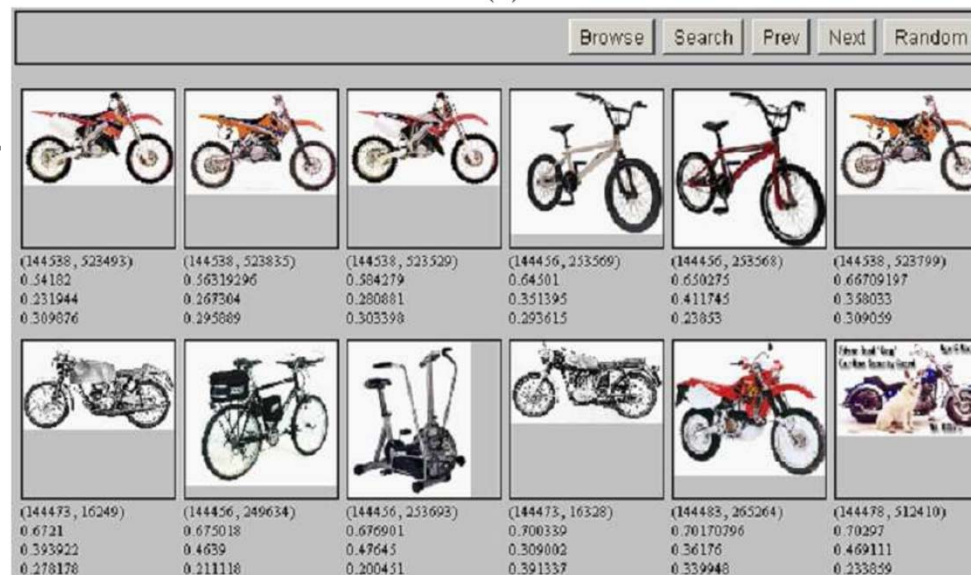
- Have a side effect:
 - Tack a user's evolving information need
 - Seeing some documents may lead users to refine their understanding of the information they are seeking
- However, most Web search users would like to complete their search in a single interaction
 - Relevance feedback is mainly a recall enhancing strategy and Web search users are only rarely concerned with getting sufficient recall
 - An important more recent thread of work is the use of clickthrough data (through query log mining or clickstream mining) to provide indirect/implicit relevance feedback (to be discussed later on!)

Explicit Relevance Feedback for Image Search

The retrieved results with the initial text query "bike"



The new top-ranked results calculated after a round of relevance feedback



Relevance Feedback for Probabilistic Model

Roberston & Sparck Jones 1976

- **Similarity Measure**

$$sim(d_j, q) \approx \sum_{i=1}^t w_{i,q} \times w_{i,j} \times \left[\log \frac{P(k_i | R)}{1 - P(k_i | R)} + \log \frac{1 - P(k_i | \bar{R})}{P(k_i | \bar{R})} \right]$$

Binary weights (0 or 1) are used

prob. of observing term k_i in the set of relevant docs

- **Initial Search (with some assumptions)**

- $P(k_i | R) = 0.5$: is constant for all indexing terms

- $P(k_i | \bar{R}) = \frac{n_i}{N}$: approx. by doc freq. of index terms

$$\Rightarrow sim(d_j, q) \approx \sum_{i=1}^t w_{i,q} \times w_{i,j} \times \left[\log \frac{0.5}{1 - 0.5} + \log \frac{1 - \frac{n_i}{N}}{\frac{n_i}{N}} \right]$$

$$= \sum_{i=1}^t w_{i,q} \times w_{i,j} \times \log \frac{N - n_i}{n_i}$$

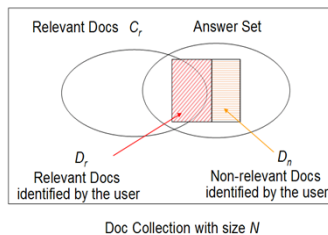
Relevance Feedback for Probabilistic Model (cont.)

- **Relevance feedback** (term reweighting alone)

$$P(k_i | R) = \frac{|D_{r,i}|}{|D_r|}$$

← Relevant docs containing term k_i
← Relevant docs

$$P(k_i | \bar{R}) = \frac{n_i - |D_{r,i}|}{N - |D_r|}$$



Approach 1

$$P(k_i | R) = \frac{|D_{r,i}| + 0.5}{|D_r| + 1}$$

$$P(k_i | \bar{R}) = \frac{n_i - |D_{r,i}| + 0.5}{N - |D_r| + 1}$$

$$P(\bar{k}_i | R) = \frac{|D_r| - |D_{r,i}| + 0.5}{|D_r| + 1}$$

Approach 2

$$P(k_i | R) = \frac{|D_{r,i}| + \frac{n_i}{N}}{|D_r| + 1}$$

$$P(k_i | \bar{R}) = \frac{n_i - |D_{r,i}| + \frac{n_i}{N}}{N - |D_r| + 1}$$

Or $\frac{n_i - |D_{r,i}|}{N - |D_r|}$

$$\begin{aligned} \text{sim}(d_j, q) &\approx \sum_{i=1}^t w_{i,q} \times w_{i,j} \times \left[\log \frac{\frac{|D_{r,i}|}{|D_r|}}{1 - \frac{|D_{r,i}|}{|D_r|}} + \log \frac{1 - \frac{n_i - |D_{r,i}|}{N - |D_r|}}{\frac{n_i - |D_{r,i}|}{N - |D_r|}} \right] \\ &= \sum_{i=1}^t w_{i,q} \times w_{i,j} \times \log \left[\frac{|D_{r,i}|}{|D_r| - |D_{r,i}|} \cdot \frac{N - |D_r| - n_i + |D_{r,i}|}{n_i - |D_{r,i}|} \right] \end{aligned}$$

Relevance Feedback for Probabilistic Model (cont.)

- Advantages
 - Feedback process is directly related to the derivation of new weights for query terms
 - The term reweighting is optimal under the assumptions of term independence and binary doc indexing
- Disadvantages
 - Document term weights are not taken into account
 - Weights of terms in previous query formulations are disregarded
 - No query expansion is used
 - The same set of index terms in the original query is reweighted over and over again

A Variant of Probabilistic Term Reweighting

Croft 1983

<http://ciir.cs.umass.edu/>

- **Differences**

- Distinct initial search assumptions
- Within-document frequency weight included

- **Initial search (assumptions)**

$$\text{sim}(d_j, q) \propto \sum_{i=1}^t w_{i,q} w_{i,j} F_{i,j,q}$$

$$F_{i,j,q} = (C + \text{idf}_i) \bar{f}_{i,j} \quad \bar{f}_{i,j} = K + (1 + K) \frac{f_{i,j}}{\max(f_{i,j})}$$

~ Inversed document frequency

~ Term frequency

(normalized with the maximum within-document frequency)

- C and K are adjusted with respect to the doc collection

A Variant of Probabilistic Term Reweighting (cont.)

- **Relevance feedback**

$$F_{i,j,q} = \left(C + \log \frac{P(k_i | R)}{1 - P(k_i | R)} + \log \frac{1 - P(k_i | \bar{R})}{P(k_i | \bar{R})} \right) \bar{f}_{i,j}$$

$$P(k_i | R) = \frac{|D_{r,i}| + 0.5}{|D_r| + 1}$$

$$P(k_i | \bar{R}) = \frac{n_i - |D_{r,i}| + 0.5}{N - |D_r| + 1}$$

A Variant of Probabilistic Term Reweighting (cont.)

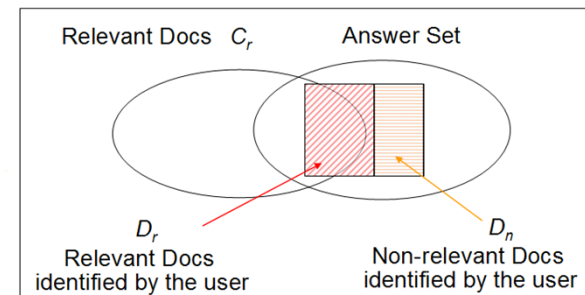
- Advantages
 - The *within-doc frequencies* are considered
 - A normalized version of these frequencies is adopted
 - Constants C and K are introduced for greater flexibility
- Disadvantages
 - More complex formulation
 - No query expansion (just reweighting of index terms)

Evaluation of Relevance Feedback Strategies

- Recall-precision figures of user reference feedback is unrealistic
 - Since the user has seen the docs during reference feedback
 - A significant part of the improvement results from the higher ranks assigned to the set R of **seen relevant docs**

$$\vec{q}_m = \alpha \cdot \vec{q} + \frac{\beta}{|D_r|} \cdot \sum_{\forall \vec{d}_i \in D_r} \vec{d}_i - \frac{\gamma}{|D_n|} \cdot \sum_{\forall \vec{d}_j \in D_n} \vec{d}_j$$

modified query original query



Doc Collection with size N

- The real gains in retrieval performance should be measured based on the docs **not seen** by the user yet

Evaluation of Relevance Feedback Strategies (cont.)

1. Recall-precision figures relative to the residual collection

- The residual collection is the set of all docs minus the set of feedback docs provided by the user
- Evaluate the retrieval performance of the modified query \vec{q}_m considering only the residual collection
- The recall-precision figures for \vec{q}_m tend to be lower than the figures for the original query \vec{q}
 - It's OK ! If we just want to compare the performance of different relevance feedback strategies

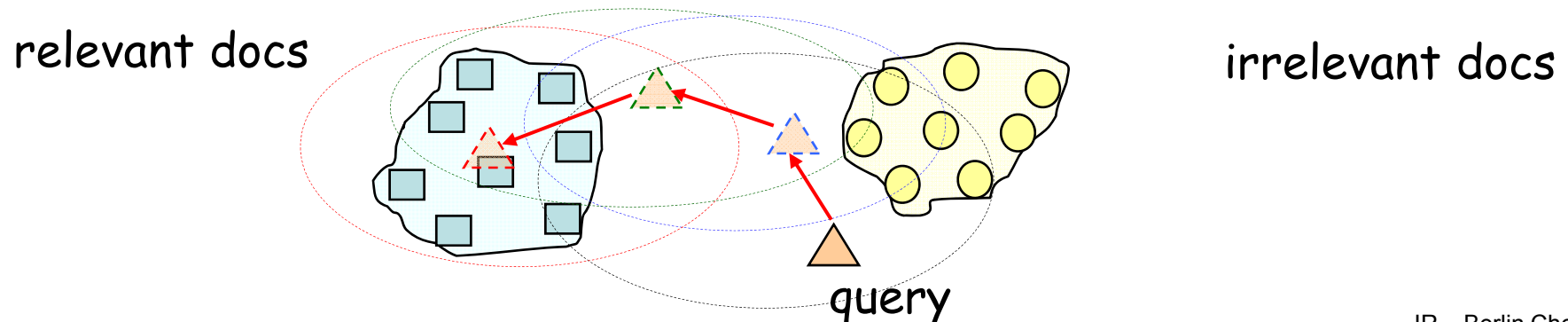
Evaluation of Relevance Feedback Strategies (cont.)

2. Or alternatively, perform a comparative evaluation of \vec{q} and \vec{q}_m on another collection
3. Or, the best evaluation of the utility of relevance feedback is to do user studies of its effectiveness in terms of how many documents a user find in a certain amount of time

Automatic Local/Global Analysis

- Remember that in user relevance feedback cycles
 - Top ranked docs separated into two classes
 - Relevant docs
 - Non-relevant docs
 - Terms in known relevant docs help describe a larger cluster of relevant docs
 - From a “clustering” perspective
 - Description of larger cluster of relevant docs is built iteratively **with assistance from the user**

Attar and Fraenkel 1977



Automatic Local/Global Analysis (cont.)

- Alternative approach: **automatically** obtain the description for a large cluster of relevant docs
 - Identify terms which are related to the query terms
 - Synonyms
 - Stemming variations
 - Terms are close each other in context

陳水扁 總統 李登輝 總統府 秘書長 陳師孟 一邊一國...

連戰 宋楚瑜 國民黨 一個中國 ...

Automatic Local/Global Analysis (cont.)

- Two strategies
 - Global analysis
 - All docs in collection are used to determine a global thesaurus-like structure for QE
 - Local analysis
 - Similar to relevance feedback but without user interference
 - Docs retrieved at query time are used to determine terms for QE
 - Local clustering, local context analysis

QE through Local Clustering

- QE through **Clustering**

- Build **global structures** such as **association matrices** to quantify term correlations
- Use the correlated terms for QE
- **But not always effective in general collections**

陳水扁 總統 呂秀蓮 綠色砂島 勇哥 吳淑珍 ...

陳水扁 視察 阿里山 小火車

- QE through **Local Clustering**

- Operate solely on the docs retrieved for the query
- **Not suitable for Web search**: time consuming
- **Suitable for intranets**

- Especially, as the assistance for search information in specialized doc collections like medical (patent) doc collections

QE through Local Clustering (cont.)

- Definition (Terminology)
 - Stem
 - $V(s)$: a non-empty subset of words which are grammatical variants of each other
 - E.g. {polish, polishing, polished}
 - A canonical form s of $V(s)$ is called a **stem**
 - e.g., $s = \text{polish}$
 - For a given query
 - Local doc set D_l : the set of documents retrieved
 - local vocabulary V_l : the set of all distinct words (stems) in the local document set
 - S_l : the set of all distinct stem derived from V_l

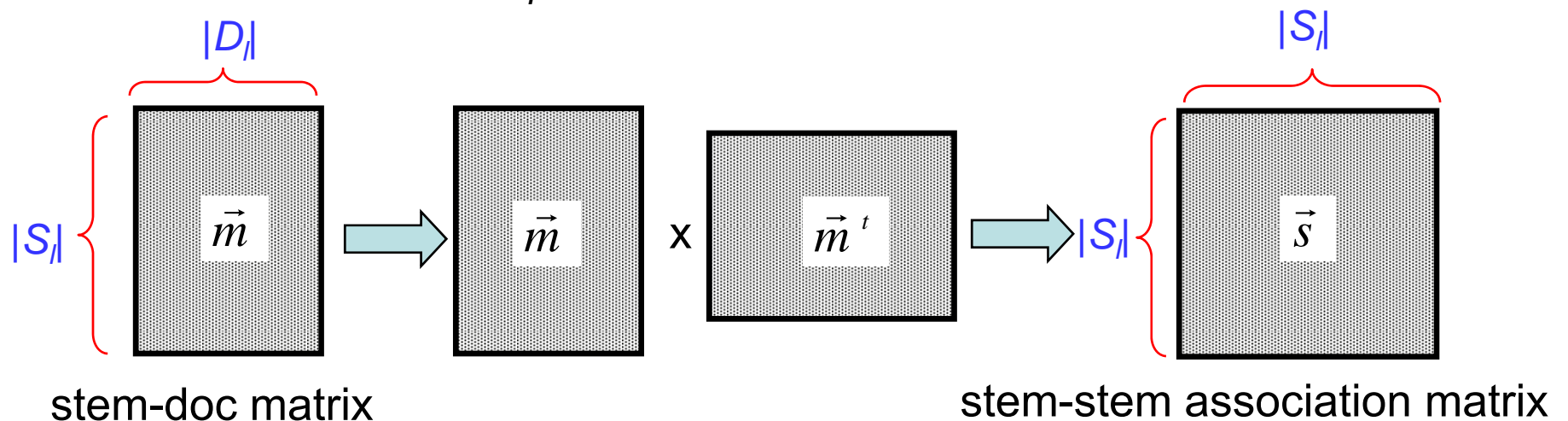
Strategies for Building Local Clusters

- **Association clusters**
 - Consider the **co-occurrence** of stems (terms) inside docs
- **Metric Clusters**
 - Consider the **distance** between two terms in a doc
- **Scalar Clusters**
 - Consider the **neighborhoods** of two terms
 - Do they have similar neighborhoods?

Strategies for Building Local Clusters (cont.)

- **Association clusters**

- Based on the **co-occurrence** of stems (terms) inside docs
 - Assumption: stems co-occurring frequently inside docs have a **synonymity** association
- An association matrix with $|S|$ rows and $|D|$ columns
 - Each entry $f_{s_i,j}$ the frequency of a stem s_i in a doc d_j



Strategies for Building Local Clusters (cont.)

- **Association clusters**

- Each entry in the stem-stem association matrix stands for **the correlation factor** between two stems

$$C_{u,v} = \sum_{d \ j \in D_l} f_{s_u, j} \times f_{s_v, j}$$

- The unnormalized form

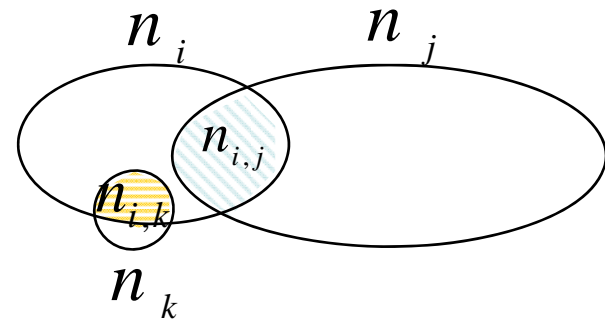
$$S_{u,v} = C_{u,v}$$

- Prefer terms with **high** frequency

- The normalized form (**ranged from 0 to 1**)

$$S_{u,v} = \frac{C_{u,v}}{C_{u,u} + C_{v,v} - C_{u,v}}$$

- Prefer terms with **low** frequency



Tanimoto coefficient

Strategies for Building Local Clusters (cont.)

- **Association clusters**

- The u -th row in the association matrix stands all the associations for the stem s_u
- A **local association cluster** $S_u(m)$
 - Defined as a set of stems s_v ($v \neq u$) with their respective values $s_{u,v}$ being the **top m** ones in the u -th row of the association matrix
- Given a query, only **the association clusters of query terms** are calculated
 - The stems (terms) belong to the association clusters are selected and added the query formulation

Strategies for Building Local Clusters (cont.)

- **Association clusters**

- Other measures for term association

- Dice coefficient

$$s_{u,v} = \frac{2 \times c_{u,v}}{c_{u,u} + c_{v,v}}$$

- Mutual information

$$s_{u,v} = MI(k_u, k_v) = \log \frac{P(k_u, k_v)}{P(k_u)P(k_v)} = \log \frac{\frac{n_{u,v}}{N}}{\frac{n_u}{N} \times \frac{n_v}{N}}$$

Strategies for Building Local Clusters (cont.)

- **Metric Clusters**

- Key idea

- Association clusters are simply based on the frequency of co-occurrence of pairs of terms in documents and do not take into account *where* the terms occur in a document
 - Two terms occurring in the same sentence seem more correlated than two terms occurring far apart in a document
 - It would be worthwhile to factor in the distance between two terms in the computation of their correlation factor

Strategies for Building Local Clusters (cont.)

- **Metric Clusters**

- Take into consideration the **distance** between two terms in a doc while computing their correlation factor

$$C_{u,v} = \sum_{d, j \in D_l} \sum_n \sum_m \frac{1}{r(k_u(n, j), k_v(n, j))}$$

- The n -th occurrence of term k_v in doc j
- $r(.)$ is a function computing the distance (in terms of the number of words) between k_u and k_v

- The entry of **local stem-stem metric correlation** matrix \vec{s} can be expressed as

- The unnormalized form

$$S_{u,v} = C_{u,v}$$

- The normalized form

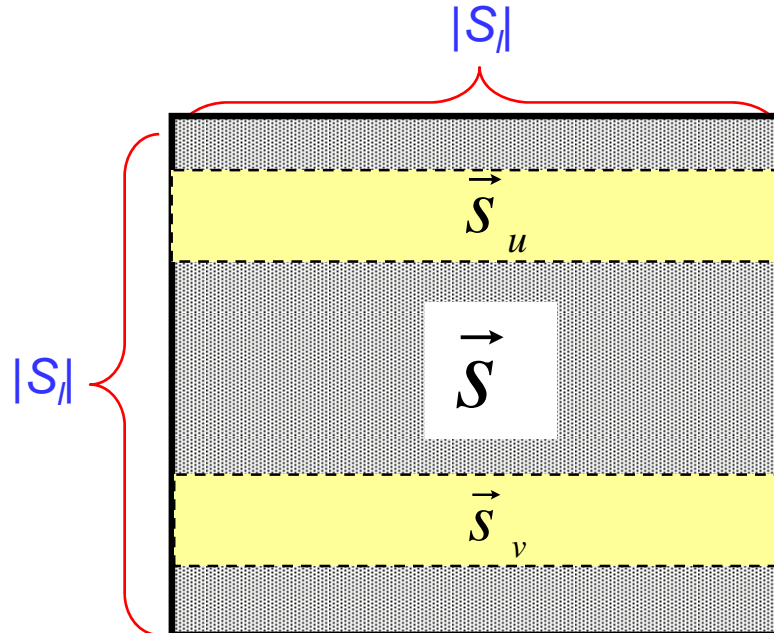
$$S_{u,v} = \frac{C_{u,v}}{\text{total number of } [k_u, k_v] \text{ pairs considered}}$$

The local association clusters of stems can be similarly defined

Strategies for Building Local Clusters (cont.)

- **Scalar Clusters**

- **Idea:** two stems (terms) with similar neighborhoods have some synonymy relationship
- Derive the synonymy relationship between two stems by comparing the sets $S_u(m)$ and $S_v(m)$



$$\longrightarrow s_{u,v} = \frac{\vec{s}_u \cdot \vec{s}_v}{|\vec{s}_u| \times |\vec{s}_v|}$$

↓

Use Cosine measure to derive a new scalar association matrix

The stem-stem association matrix achieved before

QE with Neighbor Terms

- Terms that belong to clusters associated to the query terms can be used to expand the original query
- Such terms are called neighbors of the query terms and are characterized as follows
- A term k_v that belongs to a cluster $C_u(n)$, associated with another term k_u , is said to be a **neighbor** of k_u
- Often, neighbor terms represent distinct keywords that are correlated by the current query context

QE with Neighbor Terms (cont.)

- Consider the problem of expanding a given user query q with neighbor terms
- One possibility is to expand the query as follows
- For each term $k_u \in q$, select m neighbor terms from the cluster $C_u(n)$ and add them to the query
- This can be expressed as follows:

$$q_m = q \cup \{k_v \mid k_v \in C_u(n), k_u \in q\}$$

- Hopefully, the additional neighbor terms k_v will retrieve new relevant documents

QE with Neighbor Terms (cont.)

- The set $C_u(n)$ might be composed of terms obtained using correlation factors normalized and un-normalized
- Query expansion is important because it tends to improve recall
- However, the larger number of documents to rank also tends to lower precision
- Thus, query expansion needs to be exercised with great care and fine tuned for the collection at hand

Local Context Analysis

Local context analysis combines features from both

- Local Analysis

Calculation of term correlations at query time

- Based on the set of docs retrieved for the original query
- Based on term (stem) correlation inside docs
- Terms are neighbors of **each query terms** are used to expand the query

- Global Analysis

Pre-calculation of term correlations

- Based on the whole doc collection
- The thesaurus for term relationships are built by considering small contexts (e.g. passages) and phrase structures instead of the context of the whole doc
- Terms closest to **the whole query** are selected for query expansion

Local Context Analysis (cont.)

Xu and Croft 1996

- Operations of local context analysis
 - **Document concepts**: Noun groups (named **concept** here) from retrieved docs as the units for QE instead of single keywords
 - **Concepts** selected from the top ranked passages (instead of docs) based on their co-occurrence with the whole set of query terms (no stemming)

QE through Local Context Analysis

- The operations can be further described in three steps
 - Retrieve the top n ranked passages using the original query (a doc is segmented into several passages)
 - For each concept c in the top ranked passages, the similarity $sim(q,c)$ between the whole query q and the concept c is computed using a variant of *tf-idf* ranking
 - The top m ranked concepts are added to the original query q and appropriately weighted, e.g.
 - Each concept is assigned a weight $1-0.9 \times i/m$ (i : the position in rank)
 - Original query terms are stressed by a weight of 2

QE through Local Context Analysis (cont.)

- The similarity between a concept and a query

$$sim(q, c) = \prod_{k_i \in q} \left(\delta + \frac{\log(f(c, k_i) \times idf_c)}{\log n} \right)^{idf_i}$$

δ : Set to 0.1 to avoid zero
 n : the no. of top ranked passages considered
 idf_i : emphasize the infrequent terms

$$f(c, k_i) = \sum_{j=1}^n pf_{i,j} \times pf_{c,j}$$

$$idf_c = \max \left(1, \frac{\log_{10} N / np_c}{5} \right)$$

N : the no. of passages in the collection
 np_c : Frequency of the concept c in passage j

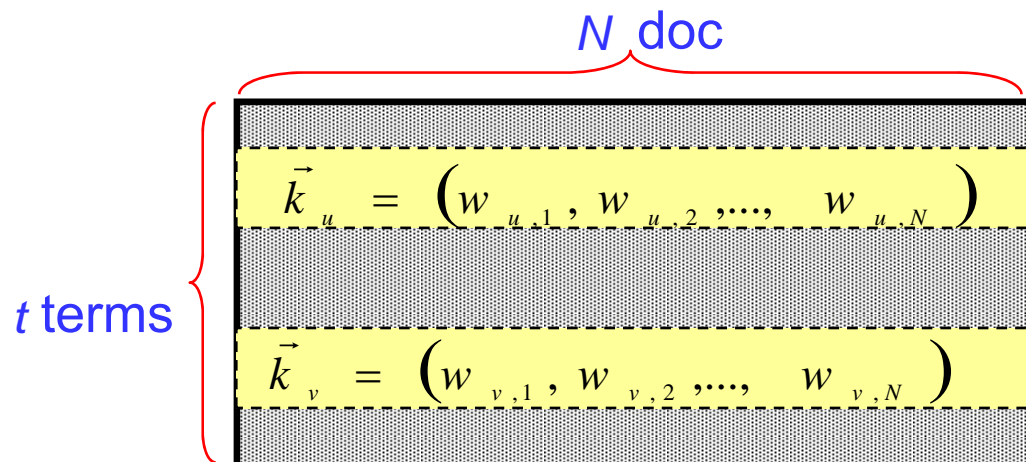
$$idf_i = \max \left(1, \frac{\log_{10} N / np_i}{5} \right)$$

np_i : the no. of passages containing concept c

QE based on a Similarity Thesaurus

Qiu and Frei 1993

- Belongs to Global Analysis
- How to construct the similarity thesaurus
 - Term to term relationships rather than term co-occurrences are considered
- How to select term for query expansion
 - Terms for query expansion are selected based on their similarity to the whole query rather the similarities to individual terms



Docs are interpreted as indexing elements here

- Doc frequency within the term vector
- Inverse term frequency

term-doc matrix

QE based on a Similarity Thesaurus (cont.)

- Definition

- $f_{u,j}$: the frequency of term k_u in document d_j
- t_j : the number of distinct index terms in document d_j
- Inverse term frequency (*ITF*)

$$ITF_j = \log \left(\frac{t}{t_j} \right) \quad (\text{doc containing more distinct terms is less important})$$

- The weight associated with each entry in the term-doc matrix

$$w_{u,j} = \frac{\left(0.5 + 0.5 \frac{f_{u,j}}{\max_g f_{u,g}} \right) \times ITF_j}{\sqrt{\sum_{l=1}^N \left[\left(0.5 + 0.5 \frac{f_{u,l}}{\max_g f_{u,g}} \right) \times ITF_l \right]^2}}$$

The importance of the doc d_j to a term k_u

Let term vector have a unit norm

QE based on a Similarity Thesaurus (cont.)

- The relationship between two terms k_u and k_v

$$c_{u,v} = \vec{k}_u \cdot \vec{k}_v = \sum_{\forall d_j} w_{u,j} \times w_{v,j}$$

is just a cosine measure?

ranged from 0 to 1

- The vector representations are normalized
- The computation is computationally expensive
 - There may be several hundred thousands of docs

QE based on a Similarity Thesaurus (cont.)

Concept-based QE

- Steps for QE based on a similarity thesaurus

1. Represent the query in the term-concept space

$$\vec{q} = \sum_{k_u \in q} w_{u,q} \times \vec{k}_u$$

2. Based on the global thesaurus, compute a similarity between the each term k_v and the whole query q

$$\text{sim}(q, k_v) = \left(\sum_{k_u \in q} w_{u,q} \times \vec{k}_u \right) \cdot \vec{k}_v = \sum_{k_u \in q} w_{u,q} \times c_{u,v}$$

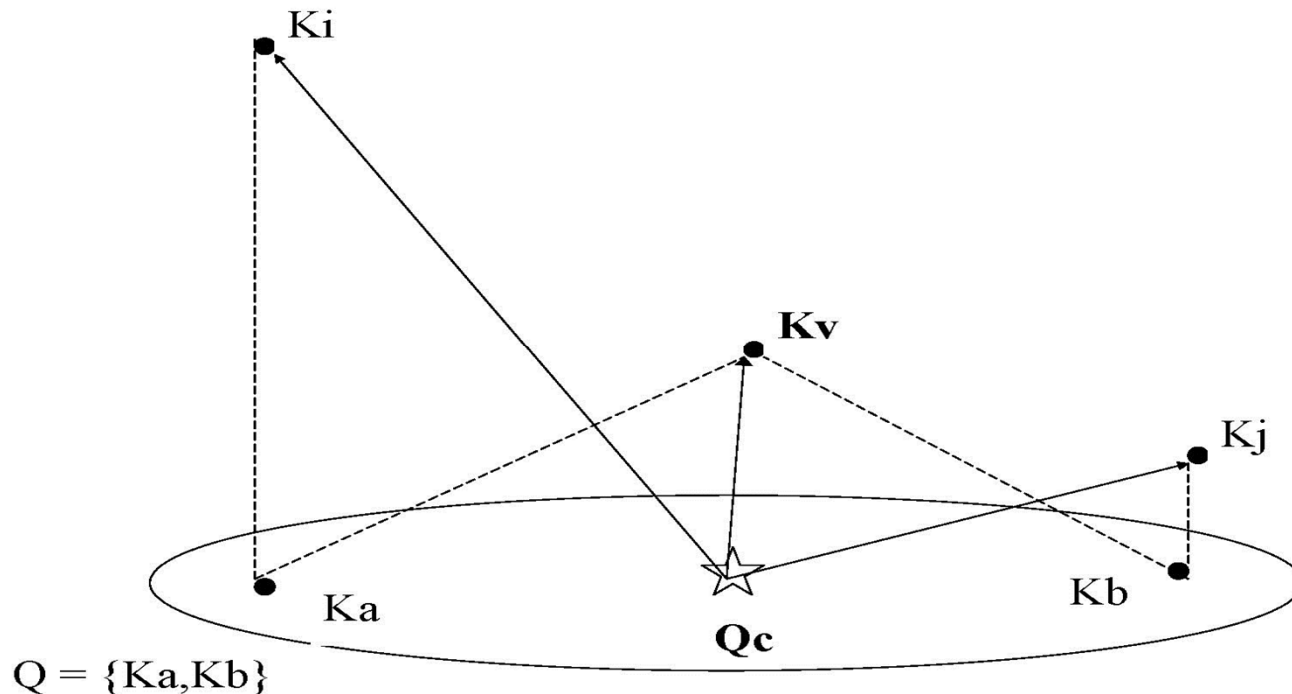
3. Expand the query with the top r ranked terms according to $\text{sim}(q, k_v)$

- The weight assigned to the expansion term

$$w_{v,q'} = \frac{\text{sim}(q, k_v)}{\sum_{k_u \in q} w_{u,q}} = \frac{\sum_{k_u \in q} w_{u,q} \times c_{u,v}}{\sum_{k_u \in q} w_{u,q}} \quad \text{ranged from 0 to 1?}$$

QE based on a Similarity Thesaurus (cont.)

- The term k_v selected for query expansion might be quite close to the whole query while its distances to individual query terms are larger



QE based on a Similarity Thesaurus (cont.)

- The similarity between query and doc measured in the term-concept space

- Doc is first represented in the **term-concept space**

$$\vec{d}_j = \sum_{k_v \in d_j} w_{v,j} \times \vec{k}_v$$

- Similarity measure

$$\text{sim}(q, d_j) \propto \sum_{k_v \in d_j} \sum_{k_u \in q} w_{v,j} \times w_{u,q} \times C_{u,v}$$

- Analogous to the formula for query-doc similarity in the generalized vector space model

- Differences

- » Weight computation

- » Only the top r ranked terms are used here

QE based on a Statistical Thesaurus

- Belongs to Global Analysis
- Global thesaurus is composed of classes which group correlated terms in the context of the whole collection
- Such correlated terms can then be used to expand the original user query
 - The terms selected must be **low frequency terms**
 - With high discrimination values

QE based on a Statistical Thesaurus (cont.)

- However, it is difficult to cluster **low frequency terms**
 - To circumvent this problem, we **cluster docs into classes instead and use the low frequency terms in these docs to define our thesaurus classes**
 - This algorithm must produce **small and tight clusters**
 - Depend on the cluster algorithm

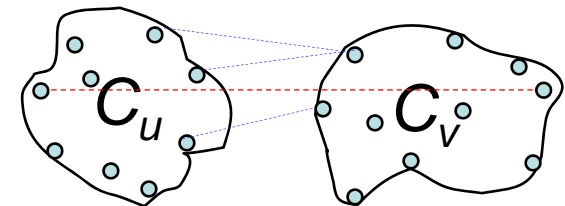
QE based on a Statistical Thesaurus (cont.)

- Complete Link Algorithm

- Place each doc in a distinct cluster
- Compute the similarity between all pairs of clusters
- Determine the pair of clusters $[C_u, C_v]$ with the highest inter-cluster similarity (using the cosine formula)
- Merge the clusters C_u and C_v
- Verify a stop criterion. If this criterion is not met then go back to step 2
- Return a hierarchy of clusters

- Similarity between two clusters is defined as

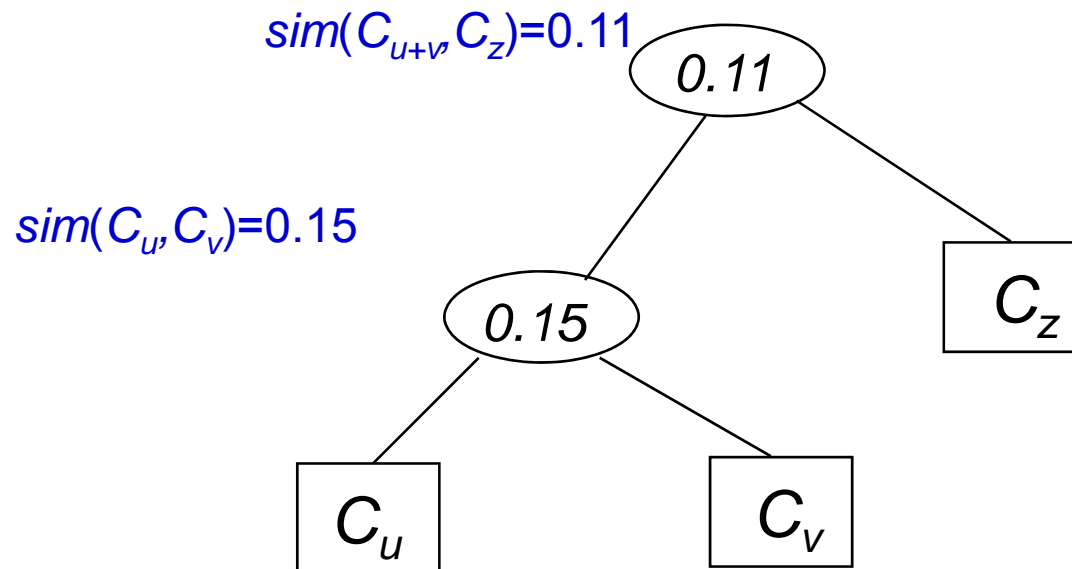
- The **minimum of similarities between all pairs** of inter-cluster docs



Cosine formula of the vector model is used

QE based on a Statistical Thesaurus (cont.)

- Example: hierarchy of three clusters



- Higher level clusters represent a looser grouping
 - Similarities decrease as moving up in the hierarchy

QE based on a Statistical Thesaurus (cont.)

- Given the doc cluster hierarchy for the whole collection, the terms that compose each class of the global thesaurus are selected as follows
 - Three parameters obtained from the user
 - *TC*: Threshold class
 - *NDC*: Number of docs in class
 - *MIDF*: Minimum inverse doc frequency

QE based on a Statistical Thesaurus (cont.)

- Use the parameter TC as threshold value for determining the doc clusters that will be used to generate thesaurus classes
 - It has to be surpassed by $\text{sim}(C_u, C_v)$ if the docs in the clusters C_u and C_v are to be selected as sources of terms for a thesaurus class
- Use the parameter NDC as a limit on the size of clusters (number of docs) to be considered
 - A low value of NDC might restrict the selection to the smaller clusters

QE based on a Statistical Thesaurus (cont.)

- Consider the set of docs in each doc cluster pre-selected above
 - Only **the lower frequency terms** are used as sources of terms for the thesaurus classes
 - The parameter *MIDF* defines the minimum value of **inverse doc frequency** for any term which is selected to participate in a thesaurus class
- Given the thesaurus classes have been built, they can be to query expansion

QE based on a Statistical Thesaurus (cont.)

- Example

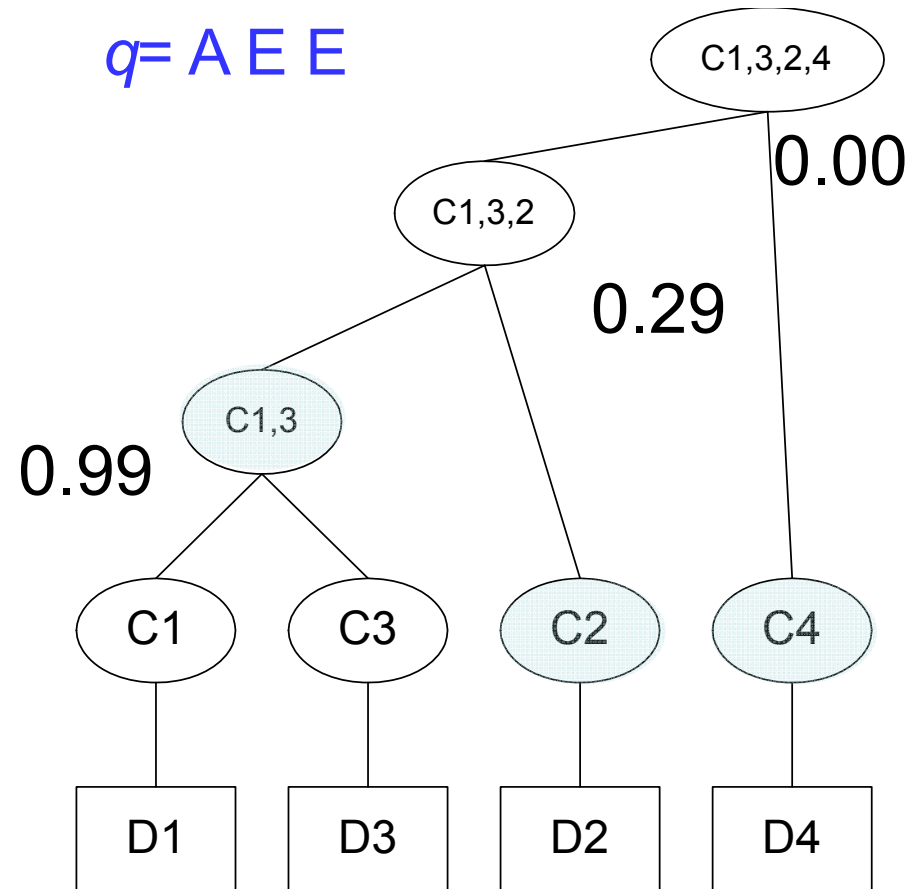
Doc1 = D, D, A, B, C, A, B, C
 Doc2 = E, C, E, A, A, D
 Doc3 = D, C, B, B, D, A, B, C, A
 Doc4 = A

sim(1,3) = 0.99
 sim(1,2) = 0.40
 sim(2,3) = 0.29
 sim(4,1) = 0.00
 sim(4,2) = 0.00
 sim(4,3) = 0.00

Cosine formula
 with *tf-idf* weighting

idf A = 0.0
 idf B = 0.3
 idf C = 0.12
 idf D = 0.12
 idf E = 0.60

$q = A E E$



- $TC = 0.90$ $NDC = 2.00$ $MIDF = 0.2$

$q' = A B E E$

QE based on a Statistical Thesaurus (cont.)

- Problems
 - Initialization of parameters TC , NDC and $MIDF$
 - TC depends on the collection
 - Inspection of the cluster hierarchy is almost always necessary for assisting with the setting of TC
 - A high value of TC might yield classes with too few terms
 - While a low value of TC yields **too few classes**

Explicit Feedback Through Clicks

- Web search engine users not only inspect the answers to their queries, they also click on them
- The **clicks** reflect preferences for particular answers **in the context of a given query**
- They can be collected in large numbers without interfering with the user actions
- The immediate question is whether they also reflect relevance judgments on the answers
- Under certain restrictions, the answer is affirmative as we now discuss

Eye Tracking

- Clickthrough data provides limited information on the user behavior
- One approach to complement information on user behavior is to use **eye tracking devices**
 - Such commercially available devices can be used to determine the area of the screen the user is focused in
 - The approach **allows correctly detecting the area of the screen of interest to the user in 60-90% of the cases**
 - Further, the cases for which the method does not work can be determined

Eye Tracking (cont.)

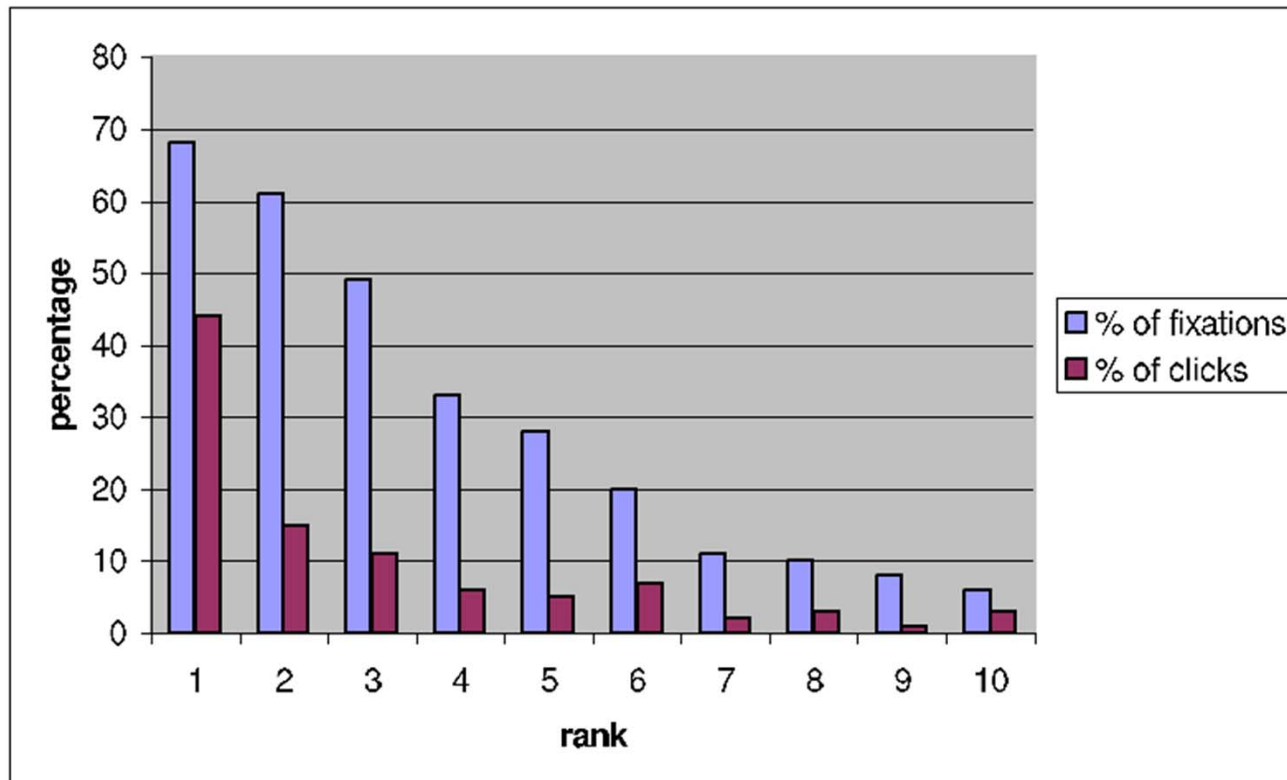
- Eye movements can be classified in four types: **fixations**, **saccades**, **pupil dilation**, and **scan paths**
 - **Fixations** are a gaze at a particular area of the screen lasting for 200-300 milliseconds
 - This time interval is large enough to allow effective brain capture and interpretation of the image displayed
 - Fixations are the ocular activity normally associated with visual information acquisition and processing
 - That is, **fixations are key to interpreting user behavior**

User Behavior

- Eye tracking experiments have shown that users scan the query results from top to bottom
- The users inspect the first and second results right away, within the second or third fixation
- Further, they tend to scan the top 5 or top 6 answers thoroughly, before scrolling down to see other answers

User Behavior (cont.)

- Percentage of times each one of the top results was viewed and clicked on by a user, for 10 test tasks and 29 subjects (Joachims et al.)



User Behavior (cont.)

- We notice that the users inspect the top 2 answers almost equally, but they click three times more in the first
- This might be indicative of a user bias towards the search engine
 - That is, that the users tend to trust the search engine in recommending a top result that is relevant

User Behavior (cont.)

- This can be better understood by presenting test subjects with two distinct result sets:
 - The normal ranking returned by the search engine
 - And, a modified ranking in which the top 2 results have their positions swapped
- Analysis suggests that the user displays a *trust bias* in the search engine that favors the top result
 - That is, the position of the result has great influence on the user's decision to click on it

Clicks as a Metric of Preferences

- Thus, it is clear that interpreting clicks as a direct indicative of relevance is not the best approach
- More promising is to interpret clicks as a metric of **user preferences**
 - For instance, a user can look at a result and decide to skip it to click on a result that appears lower
 - In this case, we say that the user prefers the result clicked on to the result shown upper in the ranking
 - This type of preference relation takes into account:
 - The results clicked on by the user
 - The results that were inspected and not clicked on
- More discussion on this issue is given in Ch. 11

Trends and Research Issues (1/3)

- Visual display
 - Graphical interfaces (2D or 3D) for relevance feedback
 - Quickly identify (by visual inspection) relationships among doc in the answer set

Allow users to visually explore the document space!

The screenshot displays a web-based interface for news retrieval. At the top left, a menu lists various news categories such as 'International Political News', 'Local Political News', 'International Business', 'Local Business', 'International Entertainment', 'Local Entertainment', 'International Sports', and 'Local Sports'. Each category has a 'Topic Map' link. Below the menu is a list of news items with checkboxes and timestamps. To the right, a large grid of colored boxes represents a topic map. The boxes are labeled with terms like '伊拉克 巴格達', '以色列 阿拉法特', '美軍 陸戰隊', '巴勒斯坦 迦薩市', '國土安全部 民航機', '聯合國 安理會', '蓋達組織 中情局', '武檢人員 武器', '阿拉法特 阿巴斯', '以色列 夏隆', '雷馬拉 任命', '約旦河 英國', '中東 鮑爾', '和平 路線', '巴格達 炸彈', and '自殺 巴士'. Red arrows indicate navigation between different levels of the hierarchy, with labels like 'go to Level-1' and 'go to Level-2'.

Lee and Chen, "Spoken document understanding and organization," *IEEE Signal Processing Magazine* 22 (5), Sept. 2005

- Utilization of local and global analysis techniques to the Web environments
 - How to alleviate the computational burden imposed on the search engine?

Trends and Research Issues (2/3)

- Yahoo! uses manually built hierarchy of concepts to assist the user with forming the query

The screenshot displays the Yahoo! search engine interface. At the top, the 'YAHOO! SEARCH' logo is visible. Below it, there are navigation links for 'Web', 'Images', 'Video', 'Audio', 'Directory', 'Local', 'News', 'Shopping', and 'More'. A search bar contains the text 'palm' and a 'Search' button. Below the search bar, there are links for 'Answers', 'My Web', 'Search Services', 'Advanced Search', and 'Preferences'. The search results section shows 'Search Results' for 'palm' with 1-10 of about 160,000,000 results. Below this, there are suggestions: 'Also try: palm springs, palm pilot, palm trees, palm reading More...'. The main search results are divided into two columns. The left column contains sponsored results for 'Official Palm Store' and 'Palms Hotel - Best Rate Guarantee', and a result for 'Palm Pilots - Palm Downloads'. The right column contains sponsored results for 'Palm Memory', 'The Palms, Turks and Caicos Islands', and 'The Palms Casino Resort, Las Vegas'.

► **Figure 9.6** An example of query expansion in the interface of the Yahoo! web search engine in 2006. The expanded query suggestions appear just below the “Search Results” bar.

Trends and Research Issues (3/3)

- Building relationships between named entities
 - Renlifang (人立方) of msra



<http://renlifang.msra.cn/>