

CS3245

# Information Retrieval

# 10

Lecture 10: XML IR,  
Relevance Feedback and Query Expansion



# Last Time

---

## Evaluating a search engine

- Benchmarks: 3 components  
Queries, documents and relevance judgments
- Precision-recall curves
- Composite, single number summaries
  
- A/B Testing
  - Test the appropriate overall evaluation criterion (OEC)



# Today

---

## Chapter 10

### 1. XML

- Basic XML concepts
- Challenges in XML IR
- Vector space model for XML IR
  
- Evaluation of XML IR

## Chapter 9

### 2. Relevance Feedback

#### *Document Level*

- Explicit RF – Rocchio (1971)
  - When does it work?
- Variants: Implicit and Blind

### 3. Query Expansion

#### *Term Level*

- Controlled Vocabularies
- WordNet
- Automatic Thesaurus Generation



# XML RETRIEVAL



# Structured retrieval



Premise: queries are structured or unstructured; documents are structured.

## Applications of structured retrieval

Digital libraries, patent databases, blogs, tagged text with entities like persons and locations (named entity tagging)

## Example

- Digital libraries: *give me a full-length article on fast fourier transforms*
- Patents: *give me patents whose claims mention RSA public key encryption and that cite US patent 4,405,829*
- Entity-tagged text: *give me articles about sightseeing tours of the Vatican and the Coliseum*



# Why RDB is not suitable in this case

## Three main problems

1. An unranked system (like a DB) can return a large set leading to information overload
2. Users often don't precisely state structural constraints – may not know possible structure elements are supported
  - *tours AND (COUNTRY: Vatican OR LANDMARK: Coliseum)?*
  - *tours AND (STATE: Vatican OR BUILDING: Coliseum)?*
3. Users may be unfamiliar with structured search and the necessary advanced search interfaces or syntax

**Solution:** adapt ranked retrieval to structured documents



# Structured Retrieval

## RDB search, Unstructured IR, Structured IR

	RDB search	unstructured retrieval	structured retrieval
objects	records	unstructured docs	trees with text at leaves
main data structure	table	inverted index	?
model	relational model	vector space & others	?
queries	SQL	free text queries	?

- Standard for encoding structured documents: Extensible Markup Language (XML)
  - structured IR → XML IR
  - also applicable to other types of markup (HTML, SGML, ...)



# XML Document

The **leaf nodes**

consist of text

The internal nodes

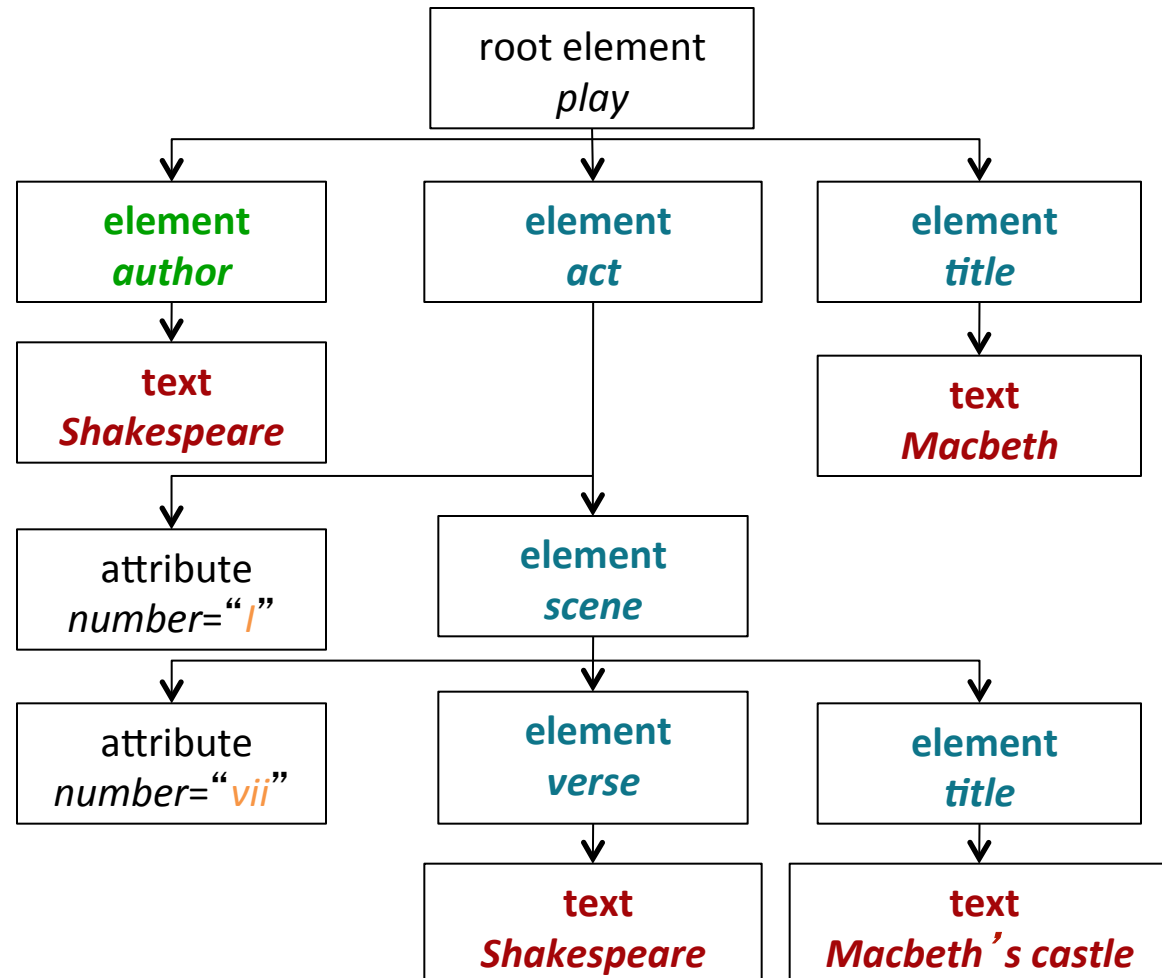
encode

**document structure**

or **metadata**

An element can  
have one or more

XML **attributes**





# CHALLENGES IN XML RETRIEVAL

# First challenge: Document parts to retrieve



- Structured or XML retrieval: users want parts of documents (i.e., XML elements), not the entire thing.

## Example

If we query Shakespeare's plays for *Macbeth's castle*, should we return the scene, the act or the entire play?

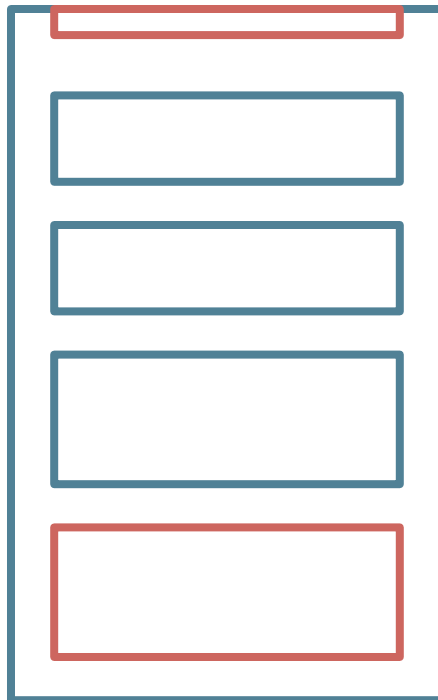
- In this case, the user is probably looking for the scene.
- However, an otherwise unspecified search for *Macbeth* should return the play of this name, not a subunit.
- Solution: structured document retrieval principle

# Structured document retrieval principle



## Structured document retrieval principle

*A system should always retrieve the most specific part of a document that answers the query.*



- Hard to implement this principle algorithmically. E.g. query: `title:Macbeth` can match both the title of the tragedy, *Macbeth*, and the title of Act I, Scene vii, *Macbeth's castle*.

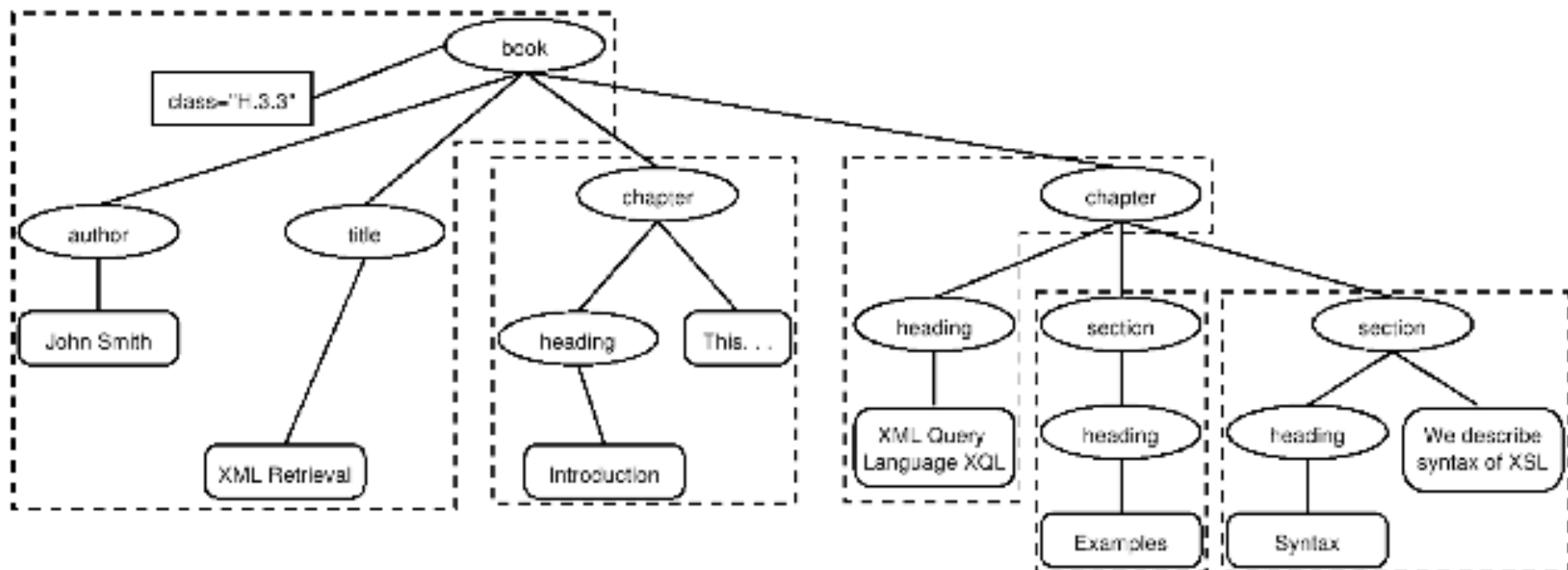
# Second challenge: Indexing Unit



- In unstructured retrieval, this is usually straightforward: files on your desktop, email messages, web pages, etc.
- In structured retrieval not so obvious what are document boundaries. 4 main methods:
  1. Non-overlapping pseudo-documents
  2. Top down
  3. Bottom up
  4. All units

# 1) Non-overlapping pseudodocuments

Group nodes into non-overlapping subtrees

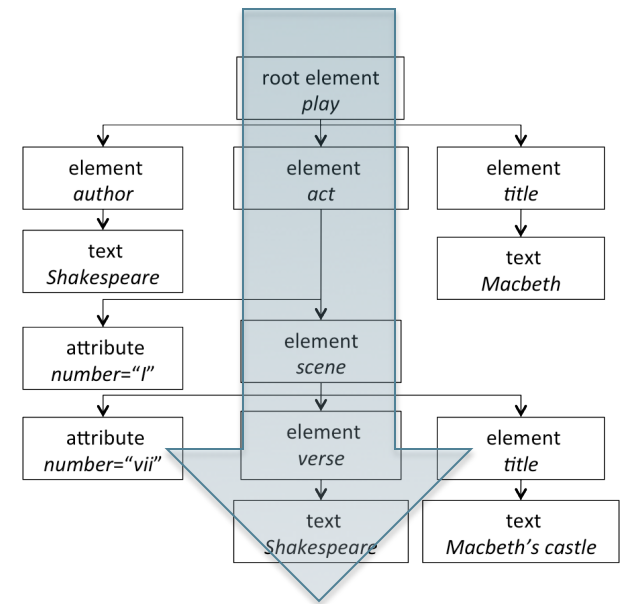


- Indexing units: books, chapters, section, but without overlap.
- Disadvantage: pseudodocuments may not make sense to the user because they are not coherent units.

## 2) Top down

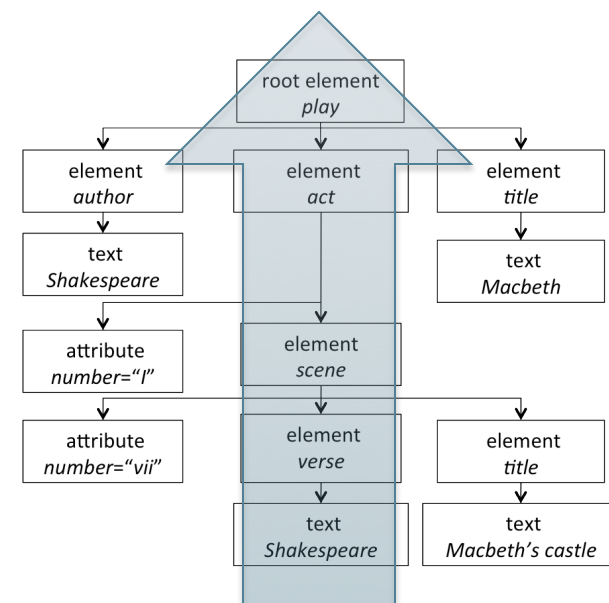


- A 2-stage process:
  1. Start with one of the largest elements as the indexing unit, e.g. the **<book>** element in a collection of books
  2. Then postprocess search results to find for each book the subelement that is the best hit.
- This two-stage process often fails to return the best subelement
  - the relevance of a whole book is often **not** a good predictor of the relevance of subelements within it.



### 3) Bottom Up

- We can search all **leaves**, select the most relevant ones and then extend them to larger units in postprocessing (bottom up).
- Similar problem as top down: the relevance of a **leaf** element is often **not** a good predictor of the relevance of elements it is contained in.



## 4) Index all elements



- The least restrictive approach, but also problematic:
- Many XML elements are not meaningful search results, e.g., an ISBN number, bolded text
- Indexing all elements means that search results will be highly redundant, due to **nested elements**.

### Example

For the query *Macbeth's castle* we would return all of the *play*, *act*, *scene* and *title* elements on the path between the root node and *Macbeth's castle*. The leaf node would then occur 4 times in the result set: 1 directly and 3 as part of other elements.

# Third challenge: Nested elements



Due to the redundancy of nested elements, it is common to restrict the set of elements eligible for retrieval.

Restriction strategies include:

- discard all small elements
- discard all elements that **users** do not look at (from examining retrieval system logs)
- discard all elements that **assessors** generally do not judge to be relevant (when relevance assessments are available)
- only keep elements that a **system designer or librarian** has deemed to be useful

In most of these approaches, result sets will still contain nested elements.

# Third challenge: Nested elements



Further techniques:

- remove nested elements in a postprocessing step to reduce redundancy.
- collapse several nested elements in the results list and use highlighting of query terms to draw the user's attention to the relevant passages.

## Highlighting

- Gain 1: enables users to scan medium-sized elements (e.g., a section); thus, if the section and the paragraph both occur in the results list, it is sufficient to show the section.
- Gain 2: paragraphs are presented in-context (i.e., their embedding section). This context may be helpful in interpreting the paragraph.

# Nested elements and term statistics

- Further challenge related to nesting: we may need to distinguish different contexts of a term when we compute term statistics for ranking, in particular inverse document frequency (idf).

## Example

The term *Gates* under the node *author* is unrelated to an occurrence under a content node like *section* if used to refer to the plural of *gate*. It makes little sense to compute a single document frequency for *Gates* in this example.

- Solution: compute idf for XML-context term pairs.
- Sparse data problems (many XML-context pairs occur too rarely to reliably estimate df)
- Compromise: consider the parent node **x** of the term and not the rest of the path from the root to **x** to distinguish contexts.

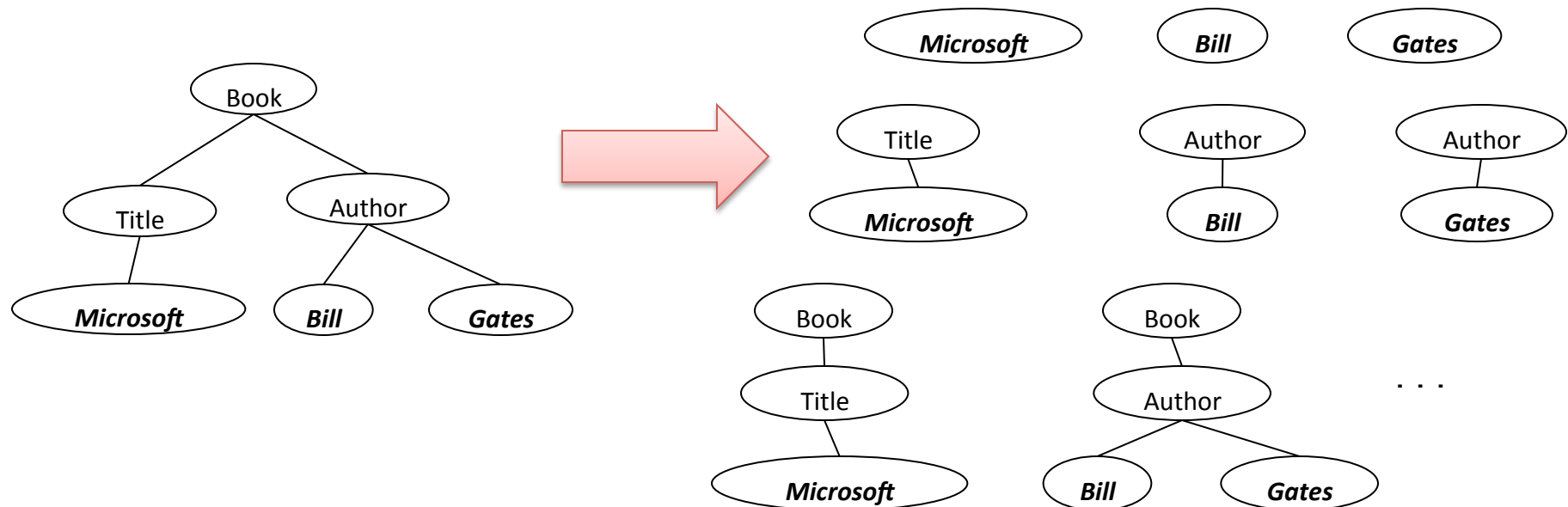


# **VECTOR SPACE MODEL FOR XML IR**

# Main idea: lexicalized subtrees

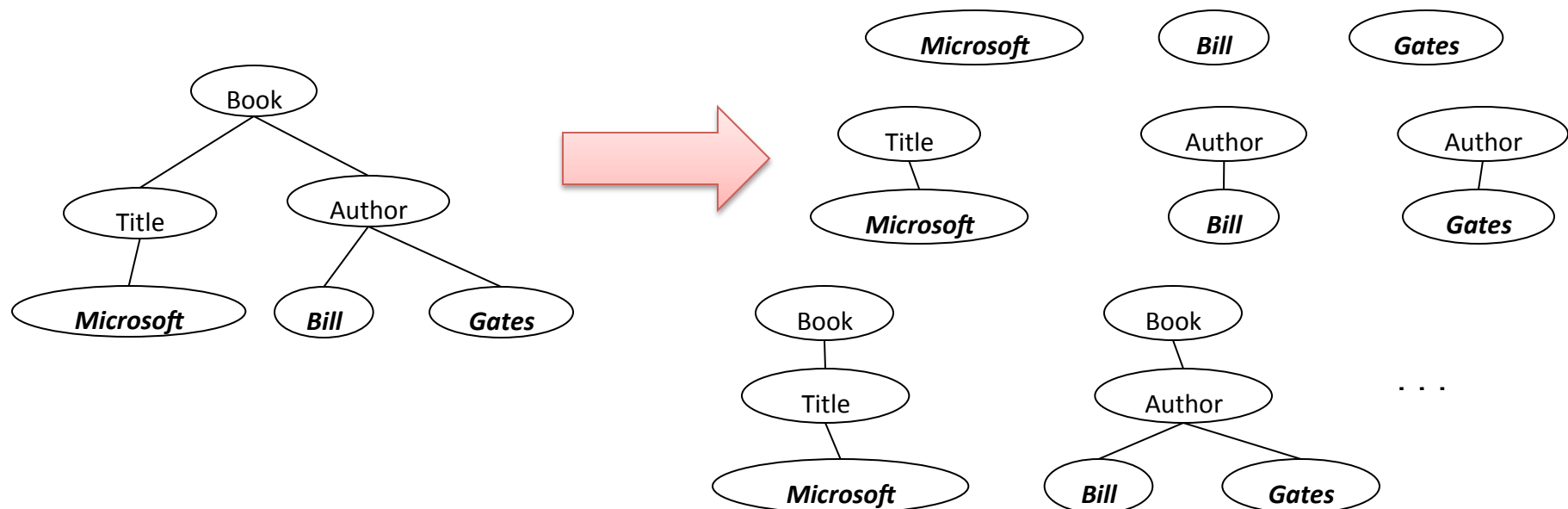
- Aim: to have each dimension of the vector space encode a word together with its position within the XML tree.
- How: Map XML documents to lexicalized subtrees.

“With words”



# Creating lexicalized subtrees

- Take each text node (leaf) and break it into multiple nodes, one for each word. E.g. split **Bill Gates** into **Bill** and **Gates**
- Define the dimensions of the vector space to be lexicalized subtrees of documents – subtrees that contain at least one vocabulary term.





# Lexicalized subtrees

- We can now represent queries and documents as vectors in this space of lexicalized subtrees and compute matches between them,
- e.g. using the vector space formalism.

## Vector space formalism in unstructured vs. structured IR

The **main difference** is that the dimensions of vector space in unstructured retrieval are **vocabulary terms** whereas they are **lexicalized subtrees** in XML retrieval.

# Structural term

- There is a tradeoff between the dimensionality of the space and the accuracy of query results.
  - If we restrict dimensions to vocabulary terms, then the VSM retrieval system will retrieve many documents that do not match the structure of the query (e.g., Gates in the title as opposed to the author element).
  - If we create a separate dimension for each lexicalized subtree in the collection, the dimensionality becomes too large.
- **Feast or Famine** Compromise: index all paths that end in a single vocabulary term (i.e., all XML-context term pairs). We call such an XML-context term pair a structural term and denote it by  $\langle c, t \rangle$ : a pair of XML-context  $c$  and vocabulary term  $t$ .



# Context resemblance

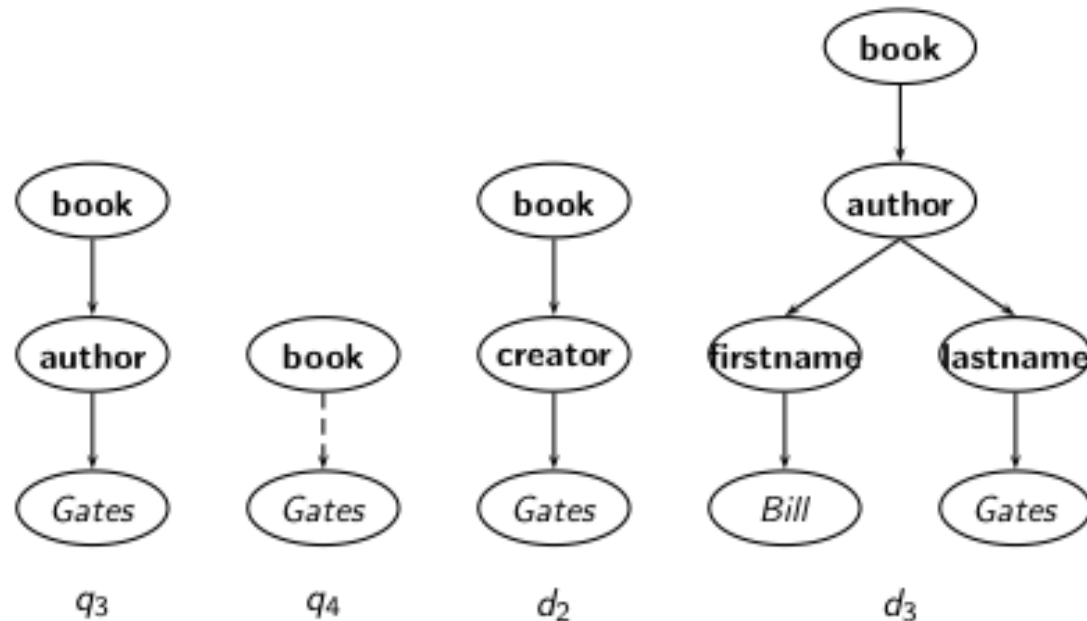
- A simple measure of the similarity of a path  $c_q$  in a query and a path  $c_d$  in a document is the following **context resemblance** function CR:

$$\text{CR}(c_q, c_d) = \begin{cases} \frac{1+|c_q|}{1+|c_d|} & \text{if } c_q \text{ matches } c_d \\ 0 & \text{if } c_q \text{ does not match } c_d \end{cases}$$

$|c_q|$  and  $|c_d|$  are the number of nodes in the query path and document path, respectively

- $c_q$  matches  $c_d$  **iff** we can transform  $c_q$  into  $c_d$  by inserting additional nodes.

# Context resemblance example

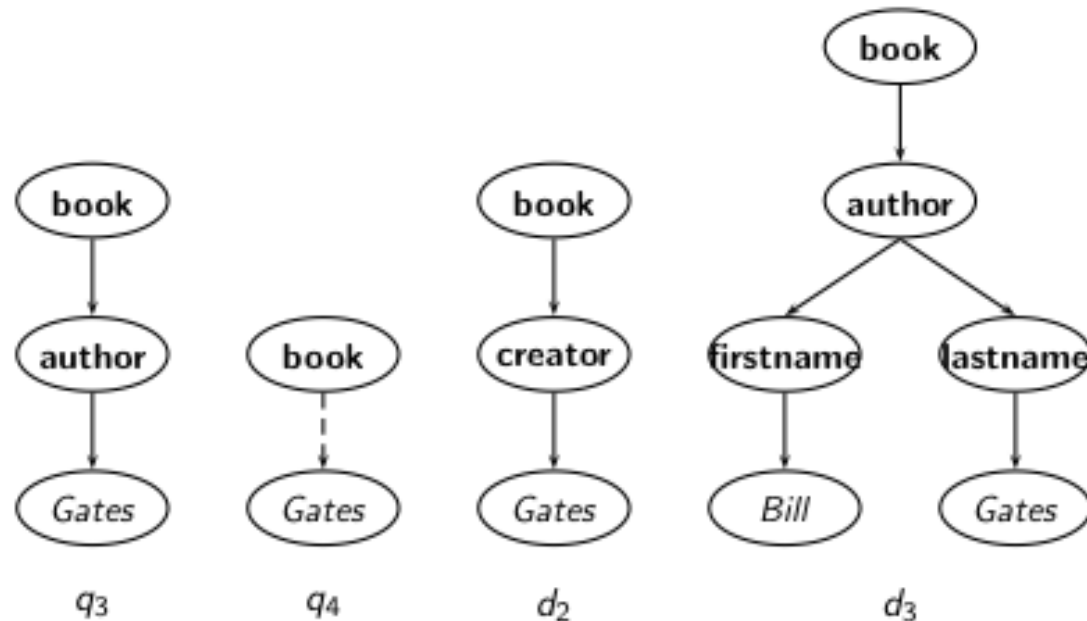


$$CR(c_q, c_d) = \begin{cases} \frac{1+|c_q|}{1+|c_d|} & \text{if } c_q \text{ matches } c_d \\ 0 & \text{if } c_q \text{ does not match } c_d \end{cases}$$

- $CR(c_{q4}, c_{d2}) = 3/4 = 0.75$ .  
The value of  $CR(c_q, c_d)$  is 1.0 if  $q$  and  $d$  are identical.

Blanks on slides, you may want to fill in

# Context resemblance example



$$CR(c_q, c_d) = \begin{cases} \frac{1+|c_q|}{1+|c_d|} & \text{if } c_q \text{ matches } c_d \\ 0 & \text{if } c_q \text{ does not match } c_d \end{cases}$$

- $CR(c_{q?}, c_{d?}) = CR(c_q, c_d) = 3/5 = 0.6.$

# Document similarity measure

- The final score for a document is computed as a variant of the cosine measure, which we call SIMNOMERGE.
- $\text{SIMNOMERGE}(q, d) =$

$$\sum_{c_k \in B} \sum_{c_l \in B} \text{CR}(c_k, c_l) \sum_{t \in V} \text{weight}(q, t, c_k) \frac{\text{weight}(d, t, c_l)}{\sqrt{\sum_{c \in B, t \in V} \text{weight}^2(d, t, c)}}$$

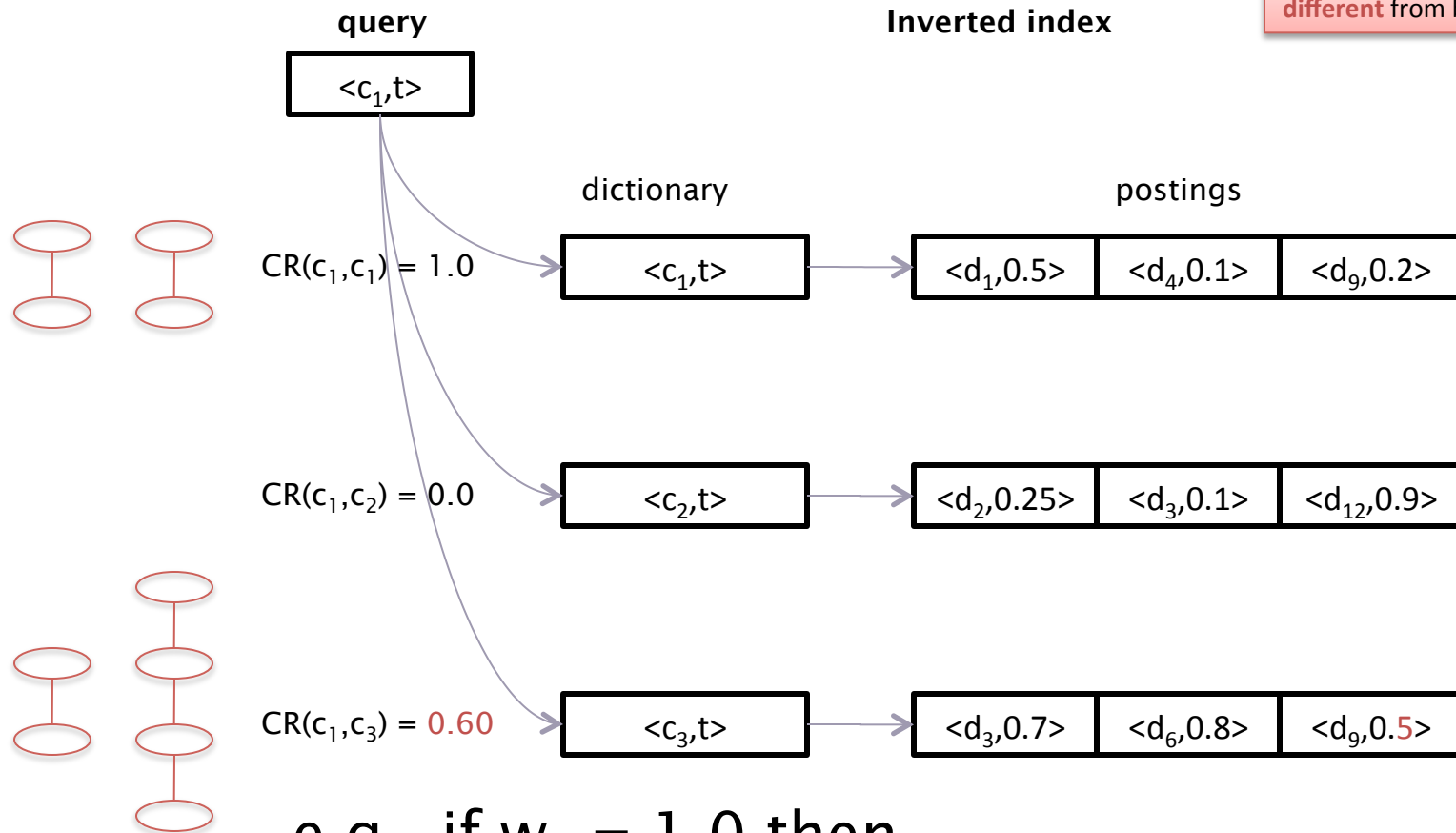
- $V$  is the vocabulary of non-structural terms
- $B$  is the set of all XML contexts
- $\text{weight}(q, t, c)$ ,  $\text{weight}(d, t, c)$  are the weights of term  $t$  in XML context  $c$  in query  $q$  and document  $d$ , resp. (standard weighting e.g.  $\text{idf}_t \times \text{wf}_{t,d}$ , where  $\text{idf}_t$  depends on which elements we use to compute  $\text{df}_t$ .)
- $\text{SIMNOMERGE}(q, d)$  is not a true cosine measure since its value can be larger than 1.0.

Blanks on slides, you may want to fill in

# SIMNOMERGE example



example **slightly different** from book



e.g., if  $w_q = 1.0$  then  
 $\text{sim}(q, d_9) = 1.0 ((1.0 \times 0.2) + (0.6 \times 0.5)) = .5$



# **XML IR EVALUATION**



## Initiative for the Evaluation of XML retrieval (INEX)

INEX: standard benchmark evaluation (yearly) that has produced test collections (documents, sets of queries, and relevance judgments). Based on IEEE journal collection (since 2006 INEX uses the much larger English Wikipedia test collection). The relevance of documents is judged by human assessors.

### INEX 2002 collection statistics

12,107	number of documents
494 MB	size
1995—2002	time of publication of articles
1,532	average number of XML nodes per document
6.9	average depth of a node
30	number of CAS topics
30	number of CO topics

# INEX Topics



- Two types:
  1. content-only or **CO topics**: regular keyword queries as in unstructured information retrieval
  2. content-and-structure or **CAS topics**: have structural constraints in addition to keywords
- Since CAS queries have both structural and content criteria, relevance assessments are more complicated than in unstructured retrieval

# INEX relevance assessments



- INEX 2002 defined component coverage and topical relevance as orthogonal dimensions of relevance.

## Component coverage

Evaluates whether the element retrieved is “structurally” correct, i.e., neither too low nor too high in the tree.

- We distinguish four cases:
  1. Exact coverage (E): The information sought is the main topic of the component and the component is a meaningful unit of information.
  2. Too small (S): The information sought is the main topic of the component, but the component is not a meaningful (self-contained) unit of information.
  3. Too large (L): The information sought is present in the component, but is not the main topic.
  4. No coverage (N): The information sought is not a topic of the component.

# INEX relevance assessments



- The **topical relevance** dimension also has four levels: highly relevant (3), fairly relevant (2), marginally relevant (1) and nonrelevant (0).

## Combining the relevance dimensions

Components are judged on both dimensions and the judgments are then combined into a digit-letter code, e.g. **2S** is a fairly relevant component that is too small. In theory, there are 16 combinations of coverage and relevance, but many cannot occur. For example, a nonrelevant component cannot have exact coverage, so the combination **3N** is not possible.

# INEX relevance assessments



- The relevance-coverage combinations are quantized as follows:

$$Q(rel, cov) = \begin{cases} 1.00 & \text{if } (rel, cov) = 3E \\ 0.75 & \text{if } (rel, cov) \in \{2E, 3L\} \\ 0.50 & \text{if } (rel, cov) \in \{1E, 2L, 2S\} \\ 0.25 & \text{if } (rel, cov) \in \{1S, 1L\} \\ 0.00 & \text{if } (rel, cov) = 0N \end{cases}$$

- This evaluation scheme takes account of the fact that binary relevance judgments are not appropriate for XML retrieval. The quantization function  $Q$  instead allows us to grade each component as partially relevant. The number of relevant components in a retrieved set  $A$  of components can then be computed as:

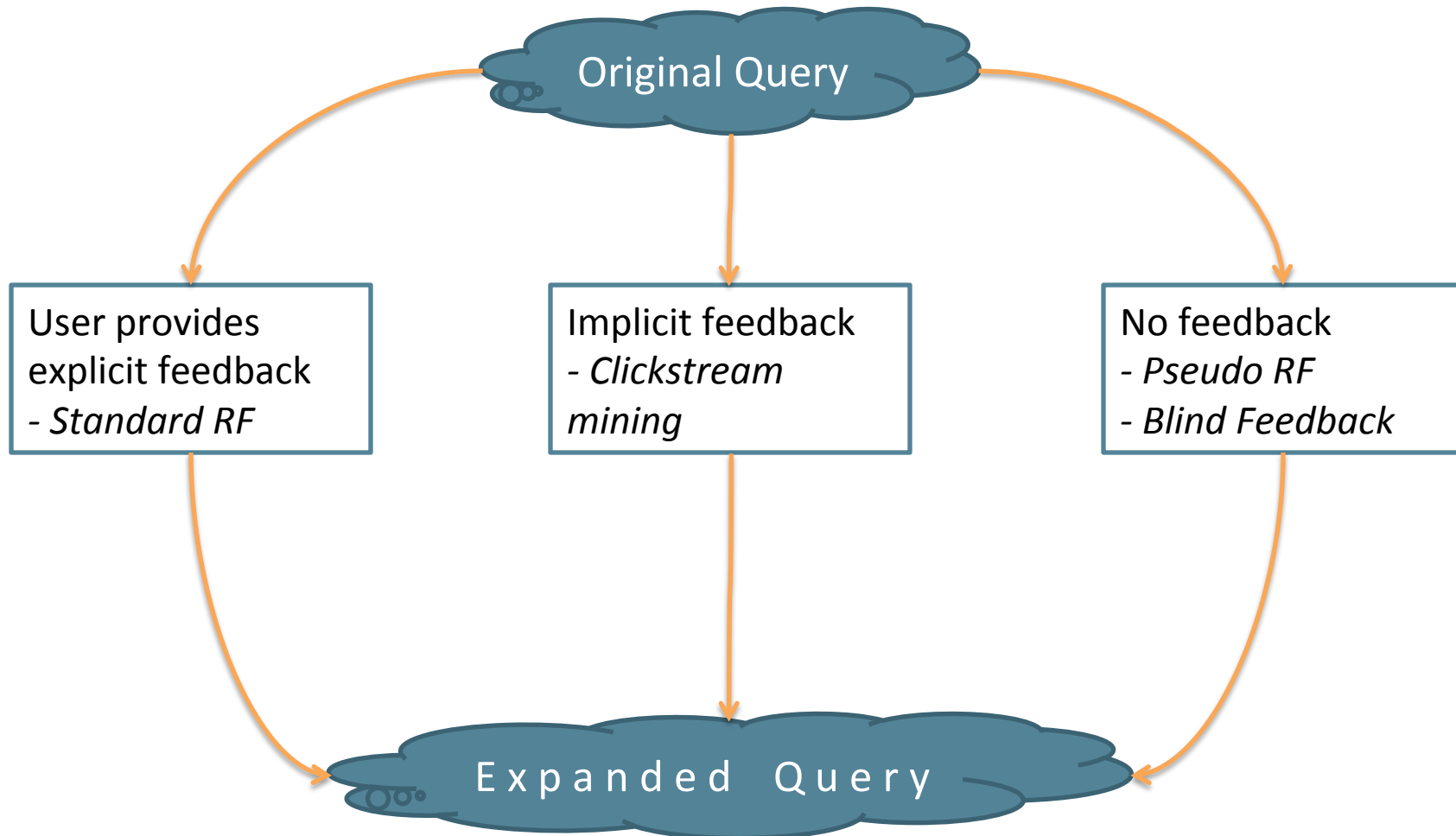
$$\#(\text{relevant items retrieved}) = \sum_{c \in A} Q(rel(c), cov(c))$$



# RELEVANCE FEEDBACK



# Relevance Feedback



# Similar pages



## Google Similar Pages beta (by Google)

★★★★★ (792) | [Fun](#) | from chrome.google.com | 162,414 users

AVAILABLE ON CHROME



OVERVIEW

DETAILS

REVIEWS

RELATED

632

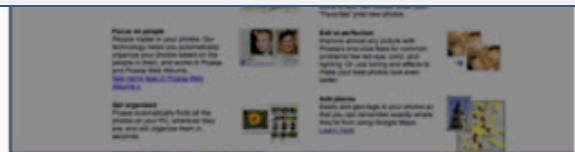


**A Google User** Jan 8, 2010

very useful..especially when you are say doing some research on a product , company, service etc. , you get objects operating in the same space. For e.g. when I open Forrester, it shows the landing pages for other companies in the same space like Gartner, IDC , etc. and this way it serves 2 purposes - you don't have to save the multiple bookmarks and more importantly, you get similar content offered by other sites

Was this review helpful?  Yes  No **Mark as spam**

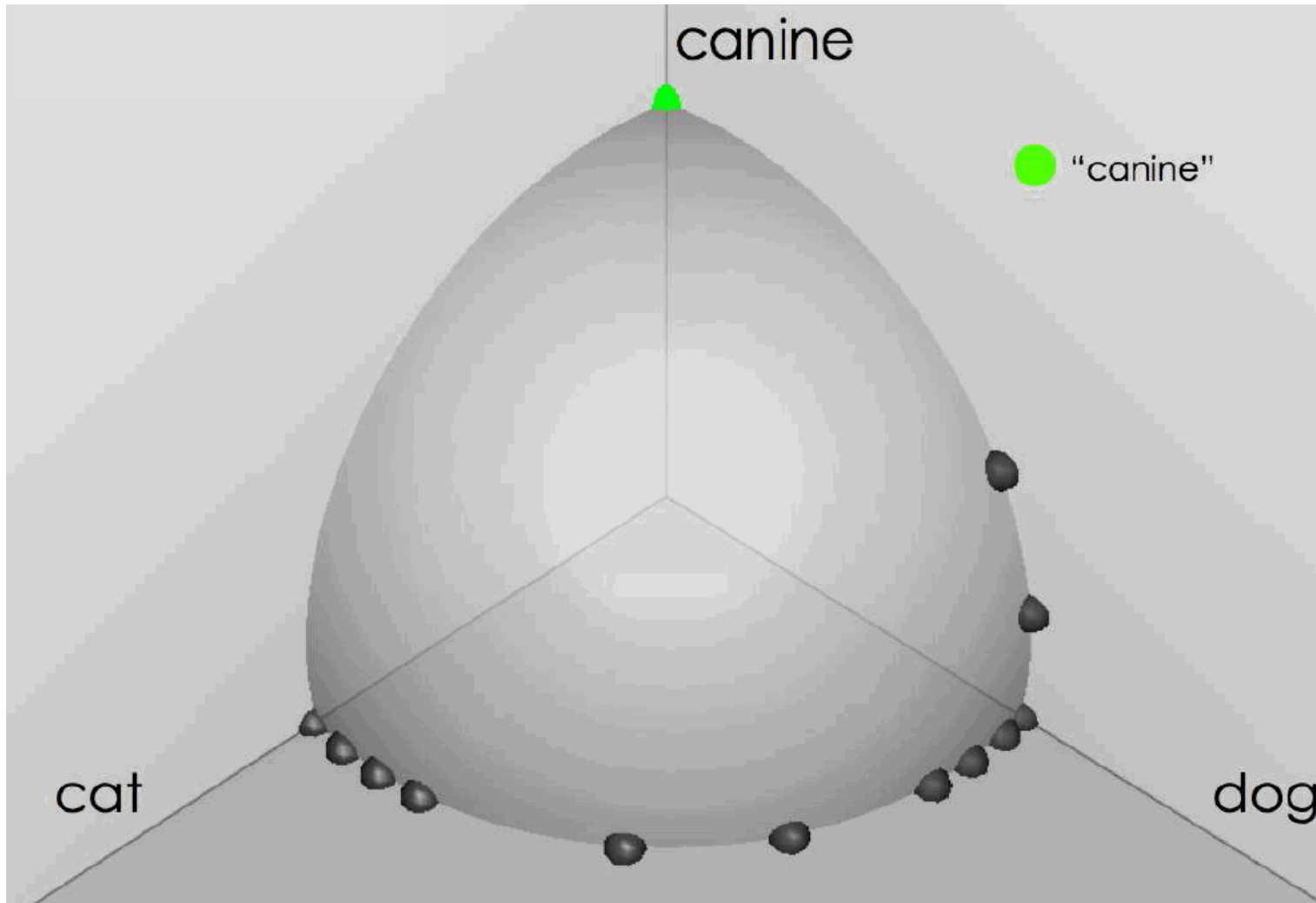
*5 out of 6 found this review helpful.*



for the page you are viewing. The data related to the query will be handled as described in Google's privacy policy (<http://www.google.com/privacypolicy.html>).

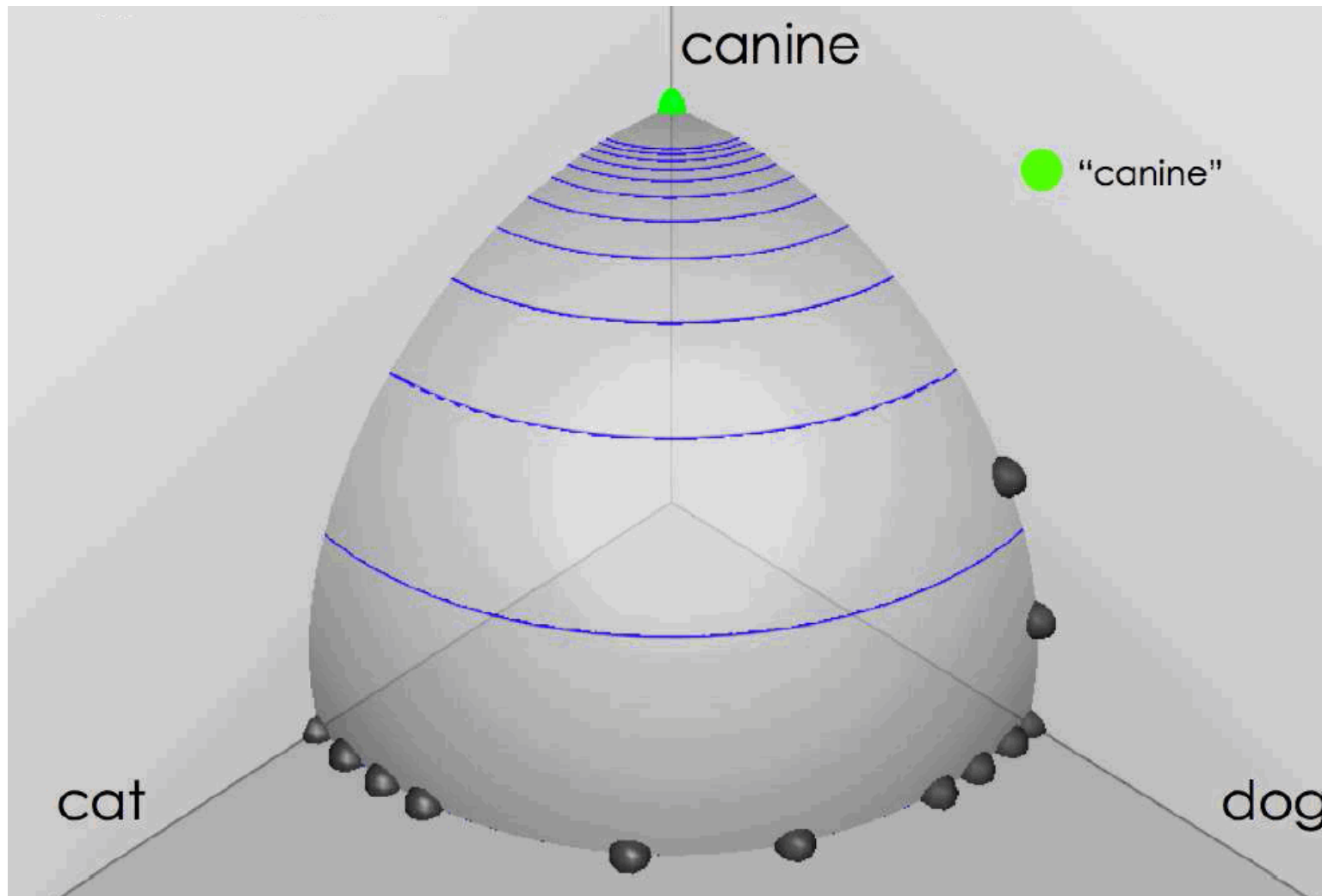
# Initial results for query *canine*

source: Fernando Diaz



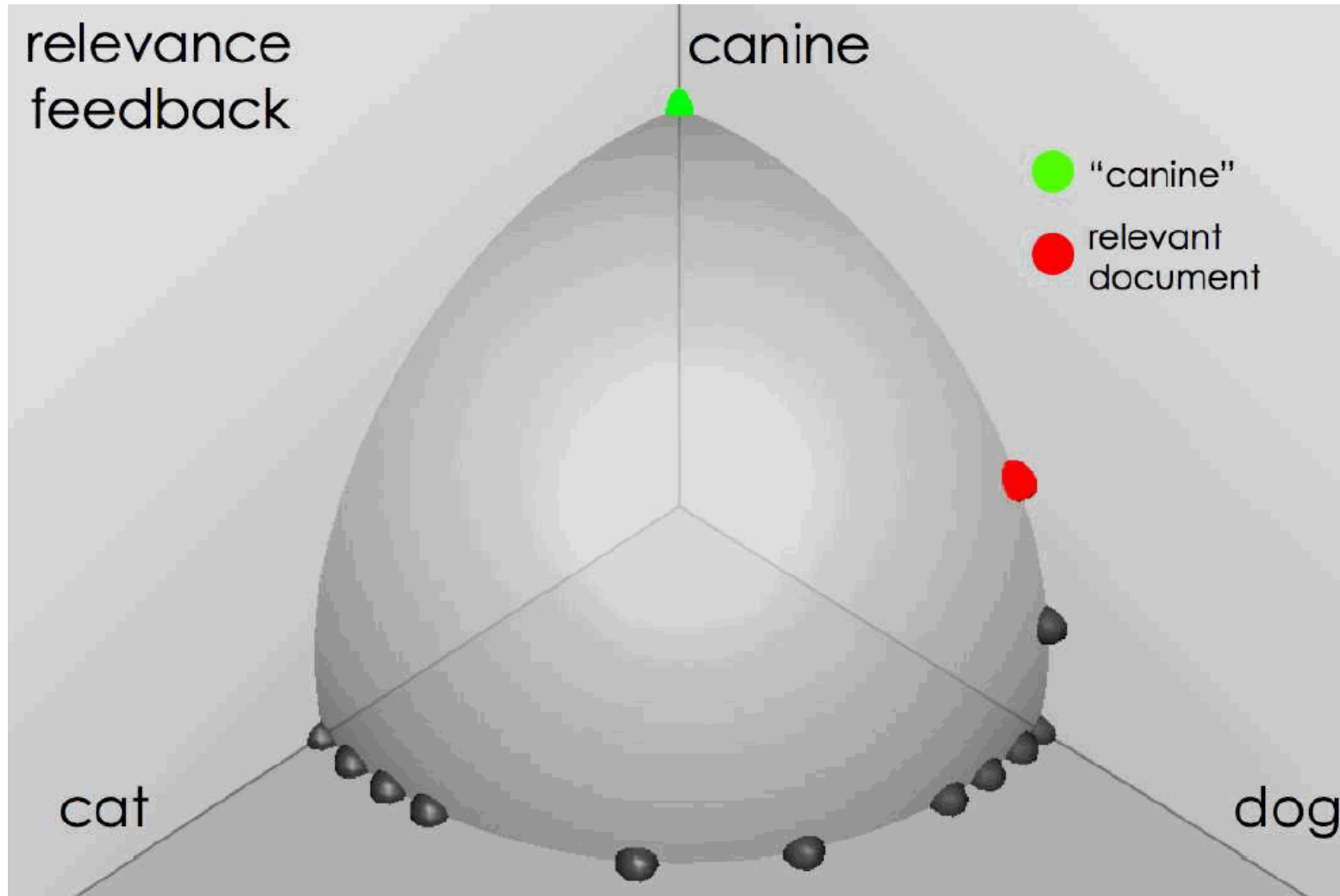
# Ad hoc results for query *canine*

source: Fernando Diaz



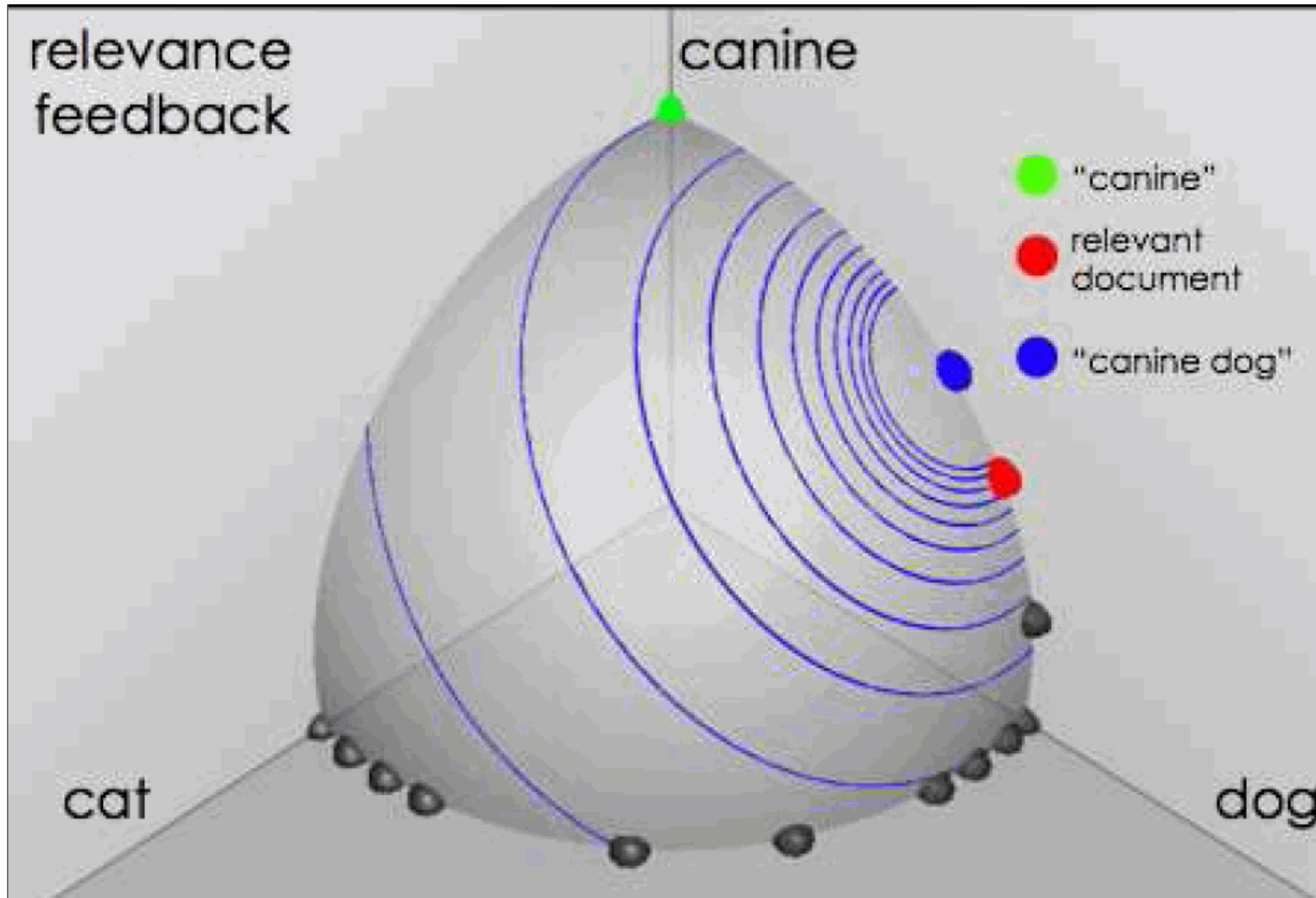
# User feedback: Select what is relevant

source: Fernando Diaz



# Results after relevance feedback

source: Fernando Diaz



# Initial query/results



Initial query: *New space satellite applications*

User marks  
relevant  
items

- + 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](#)
- 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, [Arianespace Receives Satellite Launch Pact From Telesat Canada](#)
- + 8. 0.509, 12/02/87, [Telecommunications Tale of Two Companies](#)

Assume  
others as  
nonrelevant



## Expanded query after relevance feedback

---

2.074 new	15.10 space
30.81 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



# Results for the expanded query

- 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](#)



Originally  
Marked  
Relevant  
Documents

# Key concept: Centroid

---



- The **centroid** is the center of mass of a set of points.

Definition: Centroid

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

Where  $D$  is a set of documents.

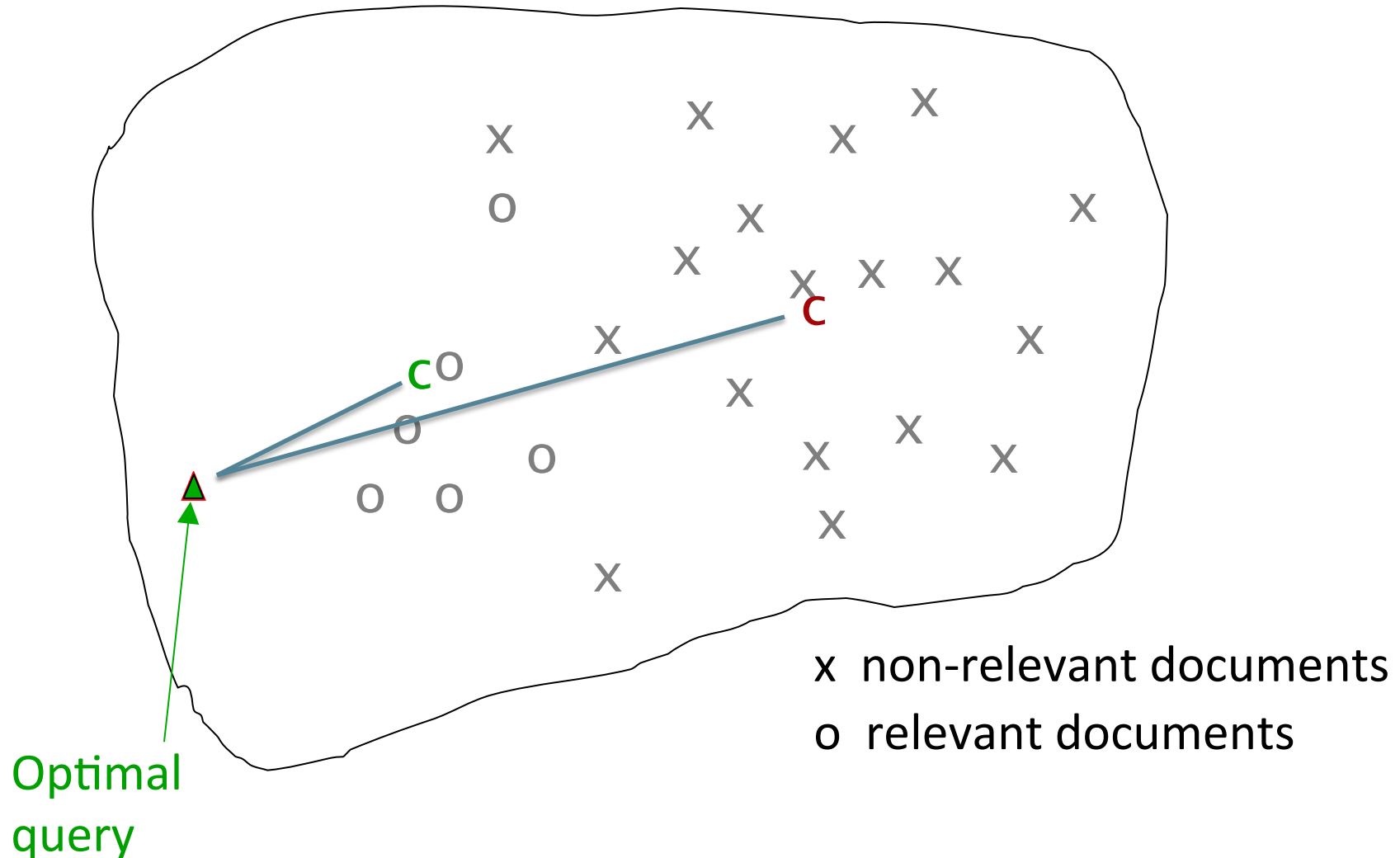


# Rocchio Algorithm

- Intuitively, we want to separate docs marked as relevant and non-relevant from each other
- The Rocchio algorithm uses the vector space model to pick a new query

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

# The Theoretically Best Query



# Rocchio (1971)



*Popularized in the SMART system (Salton)*

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 known relevant doc vectors
- $D_{nr}$  = set of known irrelevant doc vectors
- ⚠ Different from  $C_r$  and  $C_{nr}$  as we only get judgments from a few documents
- $\{\alpha, \beta, \gamma\}$  = weights (hand-chosen or set empirically)

# Weighting



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

- Tradeoff  $\alpha$  vs.  $\beta/\gamma$  : What if we have only a few judged documents?
- $\beta$  vs.  $\gamma$ : Which is more valuable?
  - Many systems only allow positive feedback ( $\gamma=0$ ). Why?
- Some weights in the query vector can go negative
  - So negative term weights are ignored (set to 0)



## When does RF work?

---

Empirically, a round of RF is often very useful. Two rounds is sometimes marginally useful.

When does it work? When two assumptions hold:

1. User's initial query at least partially works.
2. (Non)-relevant documents are similar.  
*or* term distribution in non-relevant documents are sufficiently distinct from relevant documents

# Violation of Assumption 1

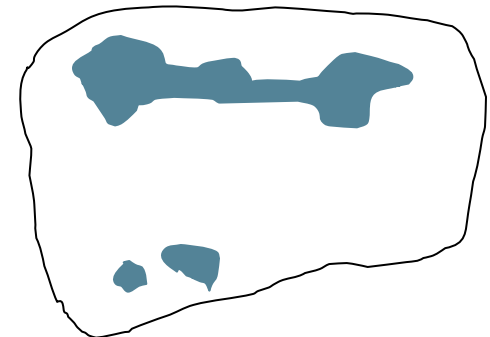
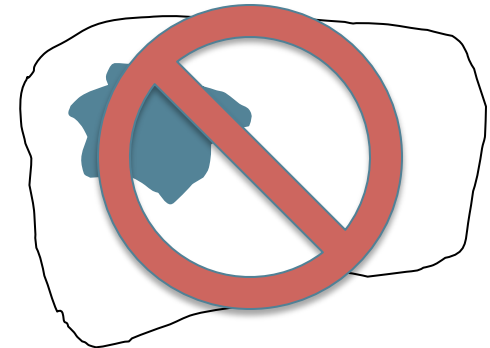


- User does not have sufficient initial knowledge.
- Examples:
  - Misspellings (but not *Brittany Speers*).
  - Cross-language information retrieval (hígado).
  - Mismatch of searcher’s vocabulary vs. collection vocabulary
    - Q: “laptop” but collection all uses “notebook”

# Violation of Assumption 2



- There are several relevance prototypes.
- Examples:
  - Burma/Myanmar: change of name
  - Instances of a general concept
  - Pop stars that worked at Burger King



# Relevance Feedback: Problems

---

- Long queries are inefficient for typical IR engine.
  - Long response times for user, as it deals with long queries.
  - Hack: reweight only a # of prominent terms, e.g., top 20.
- Users reluctant to provide explicit feedback
- Harder to understand why particular document was retrieved after RF

# Evaluation of relevance feedback strategies



Use  $q_m$  and compute precision recall graph

1. Assess on all documents in the collection
    - Spectacular improvements, but ... it's cheating!
    - Must evaluate with respect to documents not seen by user
  2. 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
- Best: use two collections each with their own relevance assessments
    - $q_o$  and user feedback from first collection
    - $q_m$  run on second collection and measured

# RF in Web search



- True evaluation of RF must also account for usability and time.
- Alternative: User revises and resubmits query.
- Users may prefer revision/resubmission to having to judge relevance of documents (more transparent)
- Some search engines offer a similar/related pages
  - Google (link-based), Altavista, Stanford WebBase
- Some don't use RF because it's hard to explain:
  - Alltheweb, Bing, Yahoo!
- Excite initially had true RF, but **abandoned it** due to lack of use.



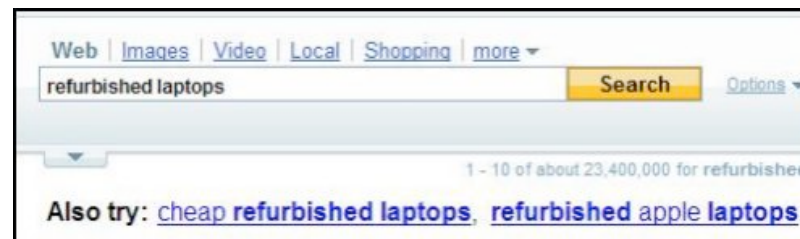
# Pseudo relevance feedback (PRF)

- **Blind feedback** automates the “manual” part of true RF, by assuming the top  $k$  is actually relevant.
- Algorithm:
  - Retrieve a ranked list of hits for the user’s query
  - Assume that the top  $k$  documents are relevant.
  - Do relevance feedback
- Works very well on average
- But can go horribly wrong for some queries
- Several iterations can cause **query drift**

# Query Expansion vs. Relevance Feedback



- 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**



# 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
  - Refinements based on query log mining



# 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
- Generally increases recall, but may decrease precision when terms are ambiguous.  
E.g., “interest rate” → “interest rate fascinate evaluate”



# An example of thesaurii: MeSH

NCBI Resources How To

PubMed.gov  
US National Library of Medicine  
National Institutes of Health

PubMed ("neopla")

RSS

Show additional filters

Article types  
Clinical Trial  
Review  
more ...

Text availability  
Abstract available  
Free full text available  
Full text available

Publication dates  
5 years  
10 years

Display Settings: [v] Summary

Results: 1 to 20 of 2

[\[Rectal cancer: im](#)

1. Krome S.  
Dtsch Med Wochensc  
PMID: 23520620 [Pub

[Isolation of low-mo](#)

2. Galbas M, Porzuce  
Acta Biochim Pol. 201  
PMID: 23520576 [Pub

MeSH Tree Structures - 2013

[Return to Entry Page](#)

1. - Anatomy [A]

- o [Body Regions \[A01\] +](#)
- o [Musculoskeletal System \[A02\] +](#)
- o [Digestive System \[A03\] +](#)
- o [Respiratory System \[A04\] +](#)
- o [Urogenital System \[A05\] +](#)
- o [Endocrine System \[A06\] +](#)
- o [Cardiovascular System \[A07\] +](#)
- o [Nervous System \[A08\] +](#)
- o [Sense Organs \[A09\] +](#)
- o [Tissues \[A10\] +](#)
- o [Cells \[A11\] +](#)
- o [Fluids and Secretions \[A12\] +](#)
- o [Animal Structures \[A13\] +](#)
- o [Stomatognathic System \[A14\] +](#)
- o [Hemic and Immune Systems \[A15\] +](#)
- o [Embryonic Structures \[A16\] +](#)
- o [Integumentary System \[A17\] +](#)
- o [Plant Structure \[A18\] +](#)
- o [Fungal Structure \[A19\] +](#)
- o [Bacterial Structure \[A20\] +](#)
- o [Viral Structure \[A21\] +](#)

2. + Organisms [B]

3. + Diseases [C]

4. + Chemicals and Drugs [D]

5. + Analytical, Diagnostic and Therapeutic Techniques and Equipment [E]

# Princeton's WordNet



## WordNet Search - 3.1

- [WordNet home page](#) - [Glossary](#) - [Help](#)

Word to search for:

Display Options:

Key: "S:" = Show Synset (semantic) relations, "W:" = Show Word (lexical) relations  
 Display options for sense: (gloss) "an example sentence"

### Noun

- **S: (n)** [washer](#), [automatic washer](#), **washing machine** (a home appliance for washing clothes and linens automatically)
  - [direct hypernym](#) / [inherited hypernym](#) / [sister term](#)
    - **S: (n)** [white goods](#) (large electrical home appliances (refrigerators or washing machines etc.) that are typically finished in white enamel)
      - **S: (n)** [home appliance](#), [household appliance](#) (an appliance that does a particular job in the home)
        - **S: (n)** [appliance](#) (durable goods for home or office use)
          - **S: (n)** [durables](#), [durable goods](#), [consumer durables](#) (consumer goods that are not destroyed by use)
            - **S: (n)** [consumer goods](#) (goods (as food or clothing) intended for direct use or

```
from nltk.corpus import wordnet as wn
```

```
wn.synsets("motorcar")
```

```
wn.synsets("car.n.01").lemma_names
```



# Automatic Thesaurus Generation

You shall know a word by the company it keeps  
– John R. Firth

- You can “harvest”, “peel”, “eat” and “prepare” **apples** and **pears**, so **apples** and **pears** must be similar
- Generate a thesaurus by analyzing the documents
- Assumption: distributional similarity
- I.e., Two words are similar if they co-occur / share same grammatical relations with similar words.

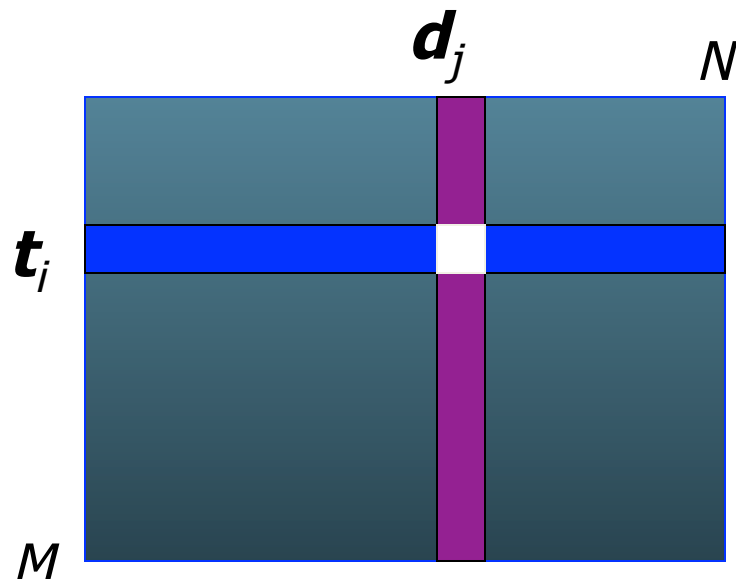
Co-occurrences are more robust; grammatical relations are more accurate. Why?

# 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  $(t_i, d_j)$



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

In NLTK. Did you forget?

A concordance permits us to see words in context. For example, we saw that then inserting the relevant word in parentheses:

```
>>> text1.similar("monstrous")
Building word-context index...
subtly impalpable pitiable curious imperial perilous trust
abundant untoward singular lamentable few maddens horrible
mystifying christian exasperate puzzled
>>> text2.similar("monstrous")
Building word-context index...
very exceedingly so heartily a great good amazingly as swee
remarkably extremely vast
>>>
```

Observe that we get different results for different texts. Austen uses this word

The term `common_contexts` allows us to examine just the contexts that are sh

```
>>> text2.common_contexts(["monstrous", "very"])
be_glad am_glad a_pretty is_pretty a_lucky
>>>
```

# Automatic Thesaurus Generation: Problems



- Term ambiguity may introduce irrelevant statistically correlated terms.
  - “Apple computer” → “Apple red fruit computer”
- Problems:
  - False positives: Words deemed similar that are not (Especially opposites)
  - False negatives: Words deemed dissimilar that are similar
- Since terms are highly correlated anyway, expansion may not retrieve many additional documents.



# Summary

---

- Chapter 9 of IIR
  1. Relevance Feedback “Documents”
  2. Query Expansion “Terms”

Rocchio: Intuition: Maximize similarity with relevant and difference from non-relevant.

In the web context, clickstream and implicit feedback common

- Resources
  - MG Ch. 4.7 and MIR Ch. 5.2 – 5.4