



Preliminaries: Information Retrieval



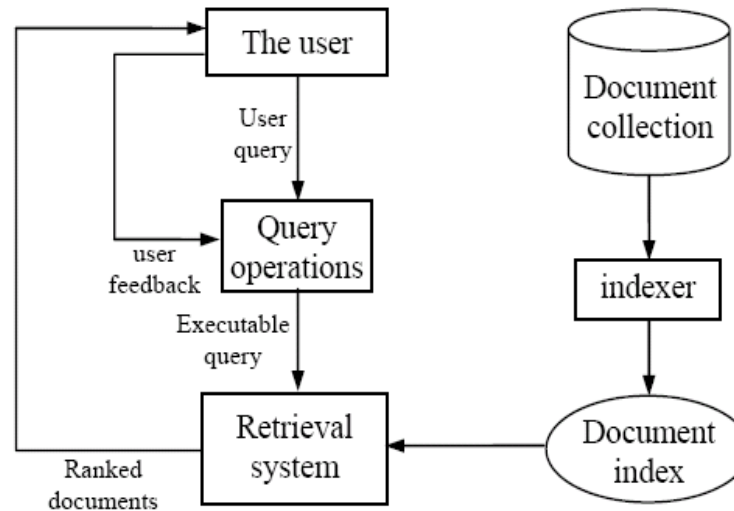
Introduction

- ⌘ Text mining refers to data mining using text documents as data.
- ⌘ Most text mining tasks use **Information Retrieval** (IR) methods to pre-process text documents.
- ⌘ These methods are quite different from traditional data pre-processing methods used for relational tables.
- ⌘ Web search also has its root in IR.

Information Retrieval (IR)

- ∞ Conceptually, IR is the study of finding needed information. I.e., IR helps users find information that matches their information needs.
 - Expressed as queries
- ∞ Historically, IR is about document retrieval, emphasizing document as the basic unit.
 - Finding documents relevant to user queries
- ∞ Technically, IR studies the acquisition, organization, storage, retrieval, and distribution of information.

IR architecture



IR queries

- ∞ Keyword queries
- ∞ Boolean queries (using AND, OR, NOT)
- ∞ Phrase queries
- ∞ Proximity queries
- ∞ Full document queries
- ∞ Natural language questions

Information retrieval models

- ∞ An IR model governs how a document and a query are represented and how the relevance of a document to a user query is defined.
- ∞ Main models:
 - Boolean model
 - Vector space model
 - etc

Boolean model

- Each document or query is treated as a “**bag**” of words or **terms**. Word sequence is not considered.
- Given a collection of documents D , let $V = \{t_1, t_2, \dots, t_{|V|}\}$ be the set of distinctive words/terms in the collection. V is called the **vocabulary**.
- A weight $w_{ij} > 0$ is associated with each term t_i of a document $\mathbf{d}_j \in D$. For a term that does not appear in document \mathbf{d}_j , $w_{ij} = 0$.
$$\mathbf{d}_j = (w_{1j}, w_{2j}, \dots, w_{|V|j}),$$
- Boolean model: weight is Boolean

Boolean model (contd)

- ∞ Query terms are combined logically using the Boolean operators **AND**, **OR**, and **NOT**.
 - E.g., ((*data* AND *mining*) AND (NOT *text*))
- ∞ Retrieval
 - Given a Boolean query, the system retrieves every document that makes the query logically true.
 - Called **exact match**.
- ∞ The retrieval results are usually quite poor because term frequency is not considered.

Vector space model

- Documents are also treated as a “bag” of words or terms.
- Each document is represented as a vector.
- However, the term weights are no longer 0 or 1. Each term weight is computed based on some variations of **TF** or **TF-IDF** scheme.
- Term Frequency (TF) Scheme:** The weight of a term t_i in document \mathbf{d}_j is the number of times that t_i appears in \mathbf{d}_j , denoted by f_{ij} . Normalization may also be applied.

TF-IDF term weighting scheme

∞ The most well known weighting scheme

- TF: still term frequency
- IDF: inverse document frequency.

N : total number of docs

df_i : the number of docs that t_i appears.

∞ The final TF-IDF term weight is:

$$tf_{ij} = \frac{f_{ij}}{\max\{f_{1j}, f_{2j}, \dots, f_{|V|j}\}}$$

$$idf_i = \log \frac{N}{df_i} \quad H(X) = - \sum_{i=1}^n p(x_i) \log_b p(x_i)$$

$$w_{ij} = tf_{ij} \times idf_i.$$

Retrieval in vector space model

- ∞ Query \mathbf{q} is represented in the same way
- ∞ **Relevance of \mathbf{d}_j to \mathbf{q}** : Compare the similarity of query \mathbf{q} and document \mathbf{d}_j .
- ∞ Cosine similarity (the cosine of the angle between the two vectors)

$$\text{cosine}(\mathbf{d}_j, \mathbf{q}) = \frac{\langle \mathbf{d}_j \bullet \mathbf{q} \rangle}{\|\mathbf{d}_j\| \times \|\mathbf{q}\|} = \frac{\sum_{i=1}^{|\mathcal{V}|} w_{ij} \times w_{iq}}{\sqrt{\sum_{i=1}^{|\mathcal{V}|} w_{ij}^2} \times \sqrt{\sum_{i=1}^{|\mathcal{V}|} w_{iq}^2}}$$

- ∞ Cosine is also commonly used in text clustering

An Example

- ⌘ A document space is defined by three terms:
 - hardware, software, users
 - the vocabulary
- ⌘ A set of documents are defined as:
 - $A1=(1, 0, 0)$, $A2=(0, 1, 0)$, $A3=(0, 0, 1)$
 - $A4=(1, 1, 0)$, $A5=(1, 0, 1)$, $A6=(0, 1, 1)$
 - $A7=(1, 1, 1)$ $A8=(1, 0, 1)$. $A9=(0, 1, 1)$
- ⌘ If the Query is “hardware and software”
- ⌘ what documents should be retrieved?

An Example (cont.)

∞ In Boolean query matching:

- document A4, A7 will be retrieved (“AND”)
- retrieved: A1, A2, A4, A5, A6, A7, A8, A9 (“OR”)

∞ In similarity matching (cosine):

- $q=(1, 1, 0)$

- $S(q, A1)=0.71$, $S(q, A2)=0.71$, $S(q, A3)=0$
- $S(q, A4)=1$ $S(q, A5)=0.5$ $S(q, A6)=0.5$
- $S(q, A7)=0.82$ $S(q, A8)=0.5$ $S(q, A9)=0.5$

- Document retrieved set (with ranking)=
 - {A4, A7, A1, A2, A5, A6, A8, A9}

Okapi relevance method

- Another way to assess the degree of relevance is to directly compute a relevance score for each document to the query.
- The **Okapi** method and its variations are popular techniques in this setting.

$$okapi(d_j, q) = \sum_{t_i \in q, d_j} \ln \frac{N - df_i + 0.5}{df_i + 0.5} \times \frac{(k_1 + 1)f_{ij}}{k_1(1 - b + b \frac{dl_j}{avdl}) + f_{ij}} \times \frac{(k_2 + 1)f_{iq}}{k_2 + f_{iq}},$$

df: document frequency, dl: document length

k1, b, k2: parameters in ranges (1,2), ~0.75, (1,1000), respectively

Relevance feedback

- ∞ Relevance feedback is one of the techniques for improving retrieval effectiveness. The steps:
 - the user first identifies some relevant (D_r) and irrelevant documents (D_{ir}) in the initial list of retrieved documents
 - the system expands the query \mathbf{q} by extracting some additional terms from the sample relevant and irrelevant documents to produce \mathbf{q}_e
 - Perform a second round of retrieval.
- ∞ **Rocchio method** (α , β and γ are parameters)

$$\mathbf{q}_e = \alpha \mathbf{q} + \frac{\beta}{|D_r|} \sum_{\mathbf{d}_r \in D_r} \mathbf{d}_r - \frac{\gamma}{|D_{ir}|} \sum_{\mathbf{d}_{ir} \in D_{ir}} \mathbf{d}_{ir}$$

Text pre-processing

- ∞ Word (term) extraction: easy
- ∞ Stopwords removal
- ∞ Stemming
- ∞ Frequency counts and computing TF-IDF term weights.

Stopwords removal

- ∞ Many of the most frequently used words in English are useless in IR and text mining – these words are called *stop words*.
 - the, of, and, to,
 - Typically about 400 to 500 such words
 - For an application, an additional domain specific stopwords list may be constructed
- ∞ Why do we need to remove stopwords?
 - Reduce indexing (or data) file size
 - stopwords accounts 20-30% of total word counts.
 - Improve efficiency and effectiveness
 - stopwords are not useful for searching or text mining
 - they may also confuse the retrieval system.

Stemming

Techniques used to find out the root/stem of a word.
E.g.,

- | | |
|---------|-------------|
| • user | engineering |
| • users | engineered |
| • used | engineer |
| • using | |

stem: use engineer

Usefulness:

improving effectiveness of IR and text mining

- matching similar words
- Mainly improve recall

reducing indexing size

- combining words with same roots may reduce indexing size as much as 40-50%.

Basic stemming methods

Using a set of rules. E.g., Porter Stemmer

∞ remove ending

- if a word ends with a consonant other than s, followed by an s, then delete s.
- if a word ends in es, drop the s.
- if a word ends in ing, delete the ing unless the remaining word consists only of one letter or of th.
- If a word ends with ed, preceded by a consonant, delete the ed unless this leaves only a single letter.
-

∞ transform words

- if a word ends with “ies” but not “eies” or “aies” then “ies → y.”

Evaluation: Precision and Recall

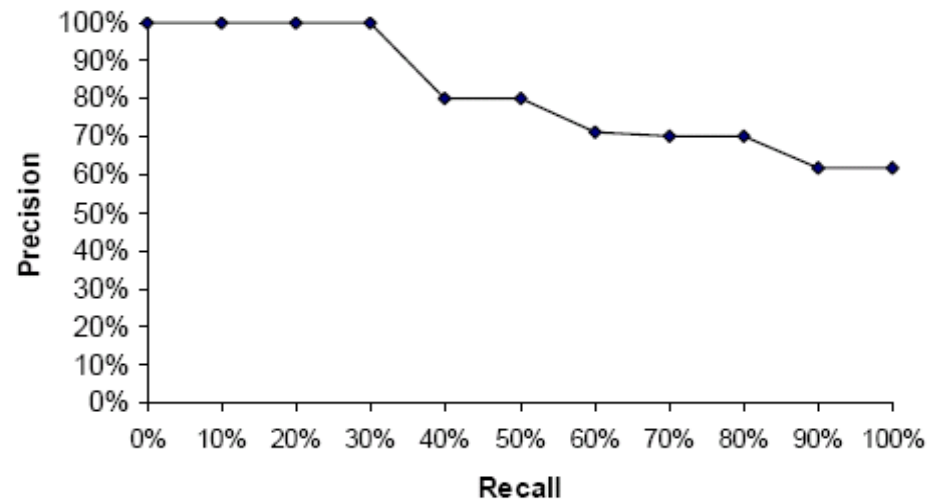
∞ Given a query:

- Are all retrieved documents relevant?
- Have all the relevant documents been retrieved?

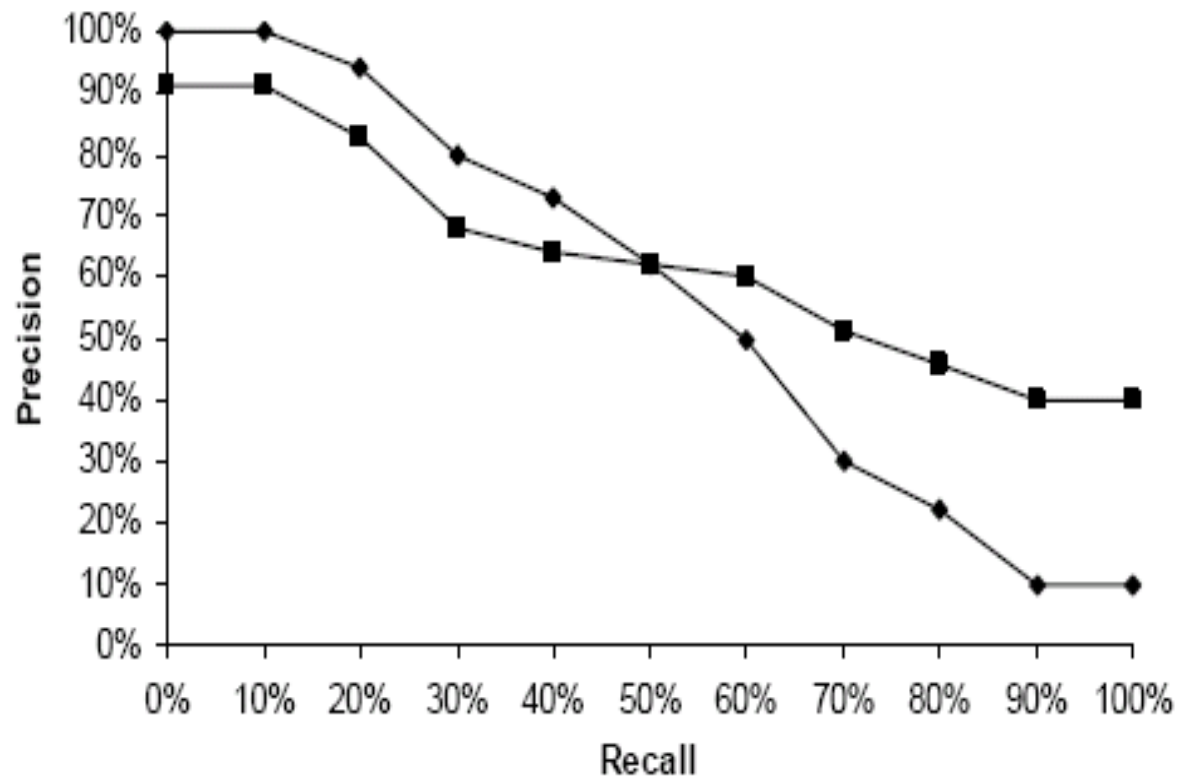
∞ Measures for system performance:

- The first question is about the **precision** of the search
- The second is about the completeness (**recall**) of the search.

Precision-recall curve



Compare different retrieval algorithms



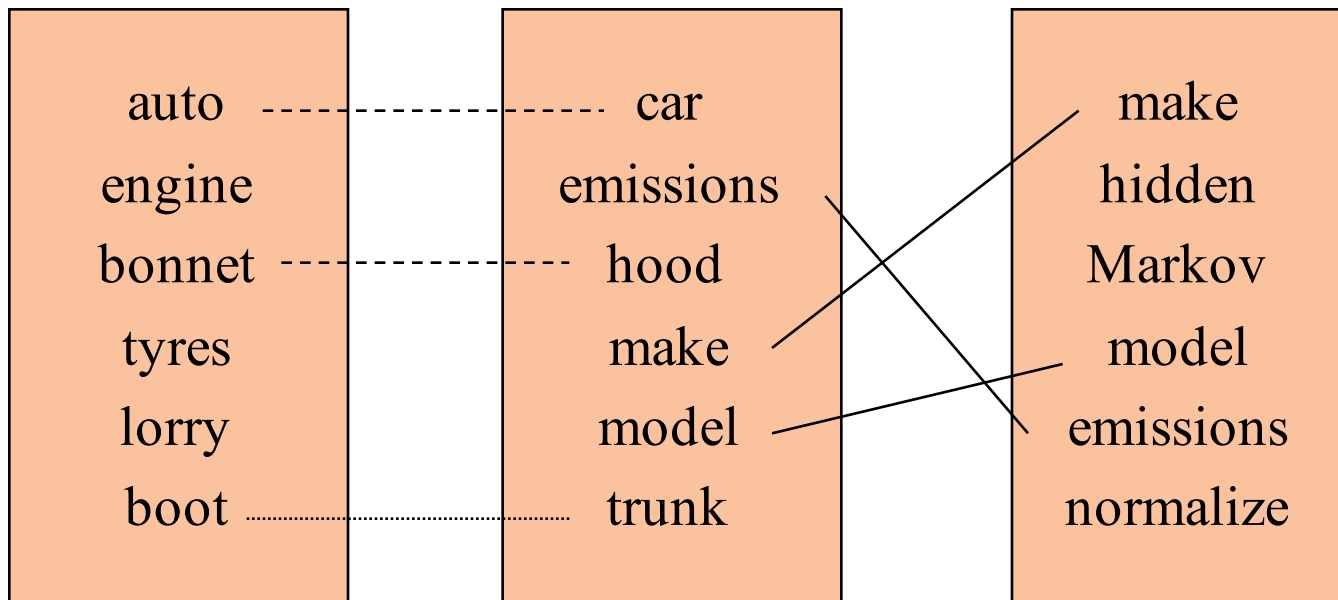
Rank precision

- ✎ Compute the precision values at some selected rank positions.
- ✎ Mainly used in Web search evaluation.
- ✎ For a Web search engine, we can compute precisions for the top 5, 10, 15, 20, 25 and 30 returned pages
 - as the user seldom looks at more than 30 pages.
- ✎ Recall is not very meaningful in Web search.
 - Why?

The Problem

∞ Example: Vector Space Model

- (from Lillian Lee)



Synonymy

Polysemy

Example

Technical Memo Titles

- c1: *Human machine interface for ABC computer applications*
 - c2: *A survey of user opinion of computer system response time*
 - c3: *The EPS user interface management system*
 - c4: *System and human system engineering testing of EPS*
 - c5: *Relation of user perceived response time to error measurement*
-
- m1: *The generation of random, binary, ordered trees*
 - m2: *The intersection graph of paths in trees*
 - m3: *Graph minors IV: Widths of trees and well-quasi-ordering*
 - m4: *Graph minors: A survey*

Term-by-Document Matrix

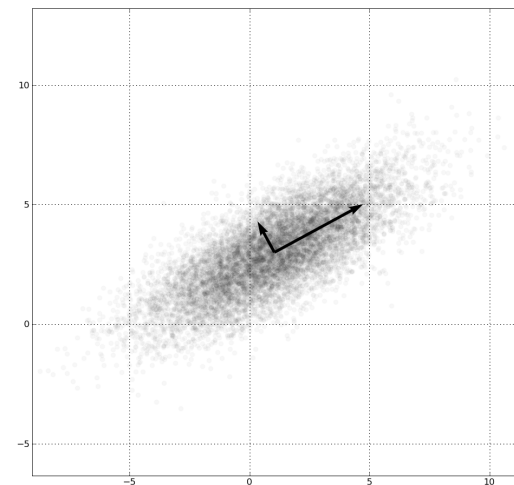
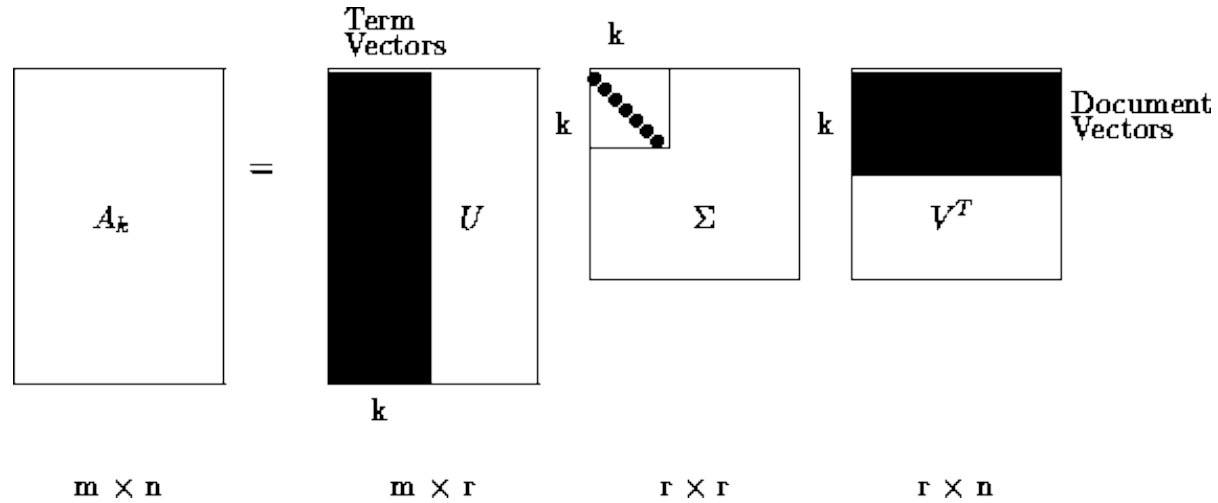
	c1	c2	c3	c4	c5	m1	m2	m3	m4
human	1	0	0	1	0	0	0	0	0
interface	1	0	1	0	0	0	0	0	0
computer	1	1	0	0	0	0	0	0	0
user	0	1	1	0	1	0	0	0	0
system	0	1	1	2	0	0	0	0	0
response	0	1	0	0	1	0	0	0	0
time	0	1	0	0	1	0	0	0	0
EPS	0	0	1	1	0	0	0	0	0
survey	0	1	0	0	0	0	0	0	1
trees	0	0	0	0	0	1	1	1	0
graph	0	0	0	0	0	0	1	1	1
minors	0	0	0	0	0	0	0	1	1

-.38

-.29

Deerwester, S., Dumais, S.T., Landauer, T.K., Furnas, G.W. and Harshman, R.A. (1990). "Indexing by latent semantic analysis." *Journal of the Society for Information Science*, 41(6), 391-407.

Singular Value Decomposition



Rank-2 Truncated Matrix

	c1	c2	c3	c4	c5	m1	m2	m3	m4
human	0.16	0.40	0.38	0.47	0.18	-0.05	-0.12	-0.16	-0.09
interface	0.14	0.37	0.33	0.40	0.16	-0.03	-0.07	-0.10	-0.04
computer	0.15	0.51	0.36	0.41	0.24	0.02	0.06	0.09	0.12
user	0.26	0.84	0.61	0.70	0.39	0.03	0.08	0.12	0.19
system	0.45	1.23	1.05	1.27	0.56	-0.07	-0.15	-0.21	-0.05
response	0.16	0.58	0.38	0.42	0.28	0.06	0.13	0.19	0.22
time	0.16	0.58	0.38	0.42	0.28	0.06	0.13	0.19	0.22
EPS	0.22	0.55	0.51	0.63	0.24	-0.07	-0.14	-0.20	-0.11
survey	0.10	0.53	0.23	0.21	0.27	0.14	0.31	0.44	0.42
trees	-0.06	0.23	-0.14	-0.27	0.14	0.24	0.55	0.77	0.66
graph	-0.06	0.34	-0.15	-0.30	0.20	0.31	0.69	0.98	0.85
minors	-0.04	0.25	-0.10	-0.21	0.15	0.22	0.50	0.71	0.62

+.94

-.83

Rank-2 Truncated Matrix

	<i>c1</i>	<i>c2</i>	<i>c3</i>	<i>c4</i>	<i>c5</i>	<i>m1</i>	<i>m2</i>	<i>m3</i>
<i>c2</i>	-0.19							
<i>c3</i>	0.00	0.00						
<i>c4</i>	0.00	0.00	0.47					
<i>c5</i>	-0.33	0.58	0.00	-0.31				
<i>m1</i>	-0.17	-0.30	-0.21	-0.16	-0.17			
<i>m2</i>	-0.26	-0.45	-0.32	-0.24	-0.26	0.67		
<i>m3</i>	-0.33	-0.58	-0.41	-0.31	-0.33	0.52	0.77	
<i>m4</i>	-0.33	-0.19	-0.41	-0.31	-0.33	-0.17	0.26	0.56

0.02	
-0.30	0.44

<i>c2</i>	0.91							
<i>c3</i>	1.00	0.91						
<i>c4</i>	1.00	0.88	1.00					
<i>c5</i>	0.85	0.99	0.85	0.81				
<i>m1</i>	-0.85	-0.56	-0.85	-0.88	-0.45			
<i>m2</i>	-0.85	-0.56	-0.85	-0.88	-0.44	1.00		
<i>m3</i>	-0.85	-0.56	-0.85	-0.88	-0.44	1.00	1.00	
<i>m4</i>	-0.81	-0.50	-0.81	-0.84	-0.37	1.00	1.00	1.00

0.92	
-0.72	1.00

Web Search as a huge IR system

- ⌘ A Web crawler (robot) crawls the Web to collect all the pages.
- ⌘ Servers establish a huge inverted indexing database and other indexing databases
- ⌘ At query (search) time, search engines conduct different types of vector query matching.

Inverted index

- ∞ The inverted index of a document collection is basically a data structure that
 - attaches each distinctive term with a list of all documents that contains the term.
- ∞ Thus, in retrieval, it takes constant time to
 - find the documents that contains a query term.
 - multiple query terms are also easy handle as we will see soon.

An example

id_1 : Web mining is useful.

1 2 3 4

id_2 : Usage mining applications.

1 2 3

id_3 : Web structure mining studies the Web hyperlink structure.

1 2 3 4 5 6 7 8

Applications: id_2

Hyperlink: id_3

Mining: id_1, id_2, id_3

Structure: id_3

Studies: id_3

Usage: id_2

Useful: id_1

Web: id_1, id_3

(A)

Applications: $\langle id_2, 1, [3] \rangle$

Hyperlink: $\langle id_3, 1, [7] \rangle$

Mining: $\langle id_1, 1, [2] \rangle, \langle id_2, 1, [2] \rangle, \langle id_3, 1, [3] \rangle$

Structure: $\langle id_3, 2, [2, 8] \rangle$

Studies: $\langle id_3, 1, [4] \rangle$

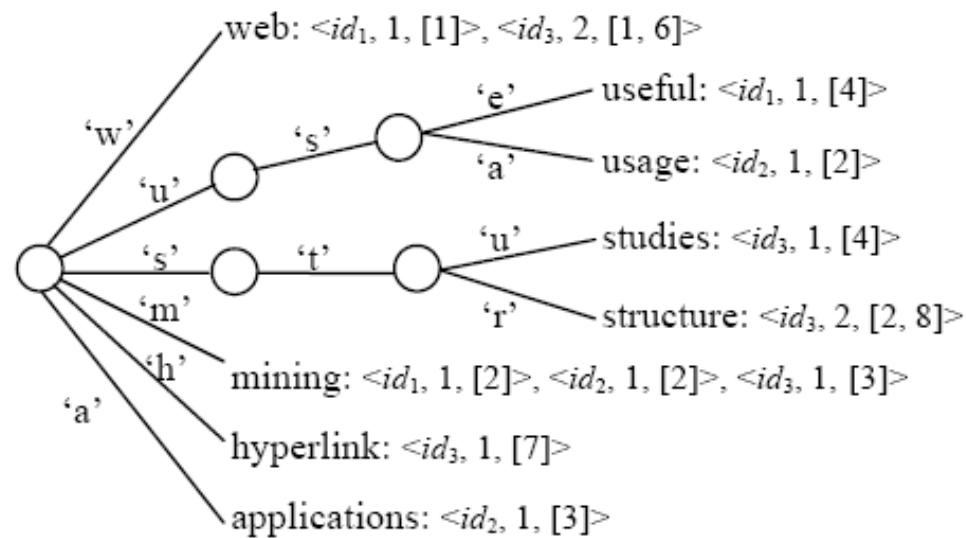
Usage: $\langle id_2, 1, [1] \rangle$

Useful: $\langle id_1, 1, [4] \rangle$

Web: $\langle id_1, 1, [1] \rangle, \langle id_3, 2, [1, 6] \rangle$

(B)

Index construction



Search using inverted index

Given a query q , search has the following steps:

- ∞ Step 1 (**vocabulary search**): find each term/word in q in the inverted index.
- ∞ Step 2 (**results merging**): Merge results to find documents that contain all or some of the words/terms in q .
- ∞ Step 3 (**Rank score computation**): To rank the resulting documents/pages, using,
 - content-based ranking
 - link-based ranking

Different search engines

- ∞ The real differences among different search engines are
 - their index weighting schemes
 - Including location of terms, e.g., title, body, emphasized words, etc.
 - their query processing methods (e.g., query classification, expansion, etc)
 - **their ranking algorithms**
 - Few of these are published by any of the search engine companies. They are tightly guarded secrets.

Page Rank

