

→ Define Information Retrieval?

Information Retrieval (IR) deals with representation, storage and organization of unstructured data. Information retrieval is the process of searching within a document collection for a particular need (a query).

→ IR can be defined as a Software Program that deals with the organization, storage, retrieval and evaluation of information from document repositories, particularly textual information.

→ Traditional examples of information-retrieval systems are online library catalogs and online document management systems such as those that store newspaper articles.

→ The data in such systems are organized as a collection of documents. A newspaper article and a catalog entry (in a library catalog) are examples of documents.

→ It uses a particular keyword (or) set of keywords in document for retrieval.  
↳ user search for a keyword

→ Eg The keywords "Stock" and "Scandal",  
locate articles about Stock-market  
The keyword "database system" may be used  
locate books on database systems.

→ Keyword-based information retrieval can  
be used not only for retrieving textual data  
also for retrieving other types of data,  
as video and audio data.

For eg a video movie may have associated  
with its keywords associated with it such  
as its title, director, actors and genre.

→ While an image or video clip may have  
tags, which are keywords describing the  
image (or) video clip associated with it.

Differences b/w Database System and  
information retrieval systems

→ Database Systems deal with the update  
and with the associated transactional requirements  
of concurrency control and durability  
whereas in information retrieval (systems)  
it is viewed as less important.

→ Database systems deal with structured information organized with relatively complex data models (such as relational model or object-oriented data models). Whereas in information retrieval the information is stored in unstructured documents.

### Keyword Search

In fulltext retrieval, all the words in each document are considered to be keywords.

→ we use the word term to refer to the words in a document.

→ Information retrieval systems typically allow query expressions formed using keywords and the logical connectives and, or and not.

→ Ranking of documents on the basis of estimated relevance to a query is critical.

Relevance ranking is based on factors such as

Term frequency (TF) :- Frequency of occurrence of query keyword in document.

## Inverse document frequency :- (IDF)

How many documents the query keyword is in. Fewer give importance to keyword.

## Hypolinks to documents

More links to a documents.

## Relevance Ranking using Terms

TF-IDF (Term Frequency / Inverse Document frequency) Ranking

Let  $n(d)$  = number of terms in the document

$n(d, t)$  = number of occurrences of term  $t$  in the document  $d$ .

### Relevance

$$TF(d, t) = \log \left( 1 + \frac{n(d, t)}{n(d)} \right)$$

If two terms "database", "silberschatz" are there

The document containing "silberschatz" but not "database" should be ranked higher than document containing term "database" but not "silberschatz"

To fix this problem, weight are assigned to term using the inverse document frequency (IDF) defined as

$$IDF(t) = \frac{1}{n(t)}$$

$n(t)$  denotes no of documents.

The relevance of a document  $d$  to a set of terms  $Q$  is then defined as

$$\pi(d, Q) = \sum_{t \in Q} TF(d, t) * IDF(t)$$

- Most systems add to the above model
- words, that occur in title, author list, section headings etc, are given greater importance
- words whose first occurrence is late in the document are given lower importance
- very common words such as "a", "an", "the", "in", "it" etc are eliminated called stop words.

→ Proximity; If keywords in query occur close together in the document, the document

has higher importance than if they occur apart.

→ Documents are returned in decreasing order of relevance score. usually only few documents are returned, not all.

### → Similarity Based Retrieval

Similarity based retrieval retrieve documents similar to given document.

→ Similarity may be defined on the basis

of common words

E.g. find k terms in A with highest

$TF(\text{d}, t)/n(t)$  and use these terms to find relevance of other documents.

→ The model of documents as points and vectors in an n-dimensional space is called the vector space model.

→ The resultant set of documents is likely to be what the user intended to find. This is called

- relevance feedback.
- Relevance feedback can also be used to help users find relevant documents from a large set of documents matching the given query keywords.
  - The cosine of the angle between the vectors of two documents is used as a measure of their similarity.

Relevance using Hyperlinks :- Popularity Ranking

The basic idea of Popularity ranking, also called Prestige ranking is to find Pages that are popular and to grant them higher than other pages that contain the specified keywords.

Eg: The term "google" may occur in vast numbers of pages, but the page google.com is most popular among the pages that contain the term "google". The page google.com should therefore be ranked as the most relevant answer to a query consisting of the term "google".

use number of hyperlinks to a site  
a measure of the popularity or prestige  
of the site.

Hub and authority based ranking

→ A hub is a page that stores links to many pages (on a topic).

→ An authority is a page that contains actual information on a topic.

→ Each page gets a hub prestige based on prestige of authorities that it points to.

→ Each page gets an authority prestige based on prestige of hubs that points to it.

→ Page Rank :- The web search engine Google introduced PageRank, which is a measure of popularity of a page based on popularity of pages that link to the page.

one drawback of PageRank algorithm is that it assigns a measure of popularity that does not take query keywords into account.

→ Search Engine Spammer refers to the practice of creating webpages or sets of web pages designed to get a high relevance rank for some queries, even though the sites are not actually popular site.

Eg: A Travel site may want to be ranked high for queries with the keyword "travel". It can get high TF-IDF scores by repeating the word "travel" many times in its page.

→ Combining TF-IDF and Popularity Ranking Measures. → By using machine learning techniques.

Synonyms, homonyms and ontologies

"Motorcycle Maintenance" replaced by "Motorcycle repair" → words with same meaning

Homonyms: - object with different meanings

Eg: - Table  
= single word with multiple meaning  
dinnerable  
(or)  
table is a database

Concept based querying :- What concept word in a document represents & similarly to understand what concepts a user looking for, and to return documents that the concepts that user is interested in.

Advantages :- A query in one language retrieve documents in other languages, so long they relate to the same concept.

→ Automated translated mechanisms can be used. If user does not understand the language in which the document is written.

ontologies :- ontologies are hierarchical structures that reflects relationships between concepts.

→ The most common relation is, the is-a relationship.

Eg:- A leopard is-a mammal and a mammal is-a animal.

→ other relationships such as Part-of are also possible.

Eg: An airplane wing is Part-of an airplane

→ There are two basic measures for assessing the quality of text retrieval  
2 Performance Metrics  
Precision — This is the Percentage of retrieved data that are actually relevant to the query (ie, "correct" responses).

It is formally represented as

$$\text{Precision} = \frac{|\{\text{Relevant}\} \cap \{\text{Retrieved}\}|}{|\{\text{Retrieved}\}|}$$

Recall :- This is the Percentage of records that are relevant to the query and were actually retrieved. It is formally represented as

$$\text{recall} = \frac{|\{\text{Relevant}\} \cap \{\text{Retrieved}\}|}{|\{\text{Relevant}\}|}$$

→ F-Score =  $\frac{\text{recall} \times \text{precision}}{(\text{recall} + \text{precision})^2}$

## Indexing of documents

- An efficient index structure is important for efficient processing of queries in information retrieval system.
- Documents that contain a specified keyword can be located efficiently by using an inverted index that maps each keyword to a list  $S_i$  of (identifiers of) documents that contain  $K_i$ .

### Inverted Index

An Inverted Index is a data structure used in information retrieval systems to efficiently retrieve documents or web pages containing specific term or set of terms.

Document<sub>1</sub> → The, quick, brown, fox, jumped, over lazy, dog.

Document<sub>2</sub> → The, Lazy, dog, slept, in, the, sun

Create an index of terms

The → Document<sub>1</sub>, document<sub>2</sub> (1, 1) (1, 7)

quick → Document<sub>1</sub> (1, 2)

brown → Doc<sub>1</sub> (1, 3)

fox → Doc<sub>1</sub> (1, 4)

jumped → Doc<sub>1</sub>, (1, 8)

over → Doc<sub>1</sub> (1, 6)

lazy → Doc<sub>1</sub>, Doc<sub>2</sub> (1, 8) (2, 2)

dog → Doc<sub>1</sub>, Doc<sub>2</sub> (1, 9) (2, 3)

slept → Doc<sub>2</sub> (2, 4)

in → Doc<sub>2</sub> (2, 5)

sun → Doc<sub>2</sub> (3, 7)

Example  
→ The documents d<sub>1</sub>, d<sub>9</sub>, d<sub>21</sub> contain the term "silberschatz", the inverted list for the keyword silberschatz would be "d<sub>1</sub>; d<sub>9</sub>; d<sub>21</sub>".  
→ It also provides the list of locations within the documents where the keyword appears.

Eg:- If "silberschatz" appeared at Position 21 in d<sub>1</sub>, positions 1 and 19 in d<sub>9</sub>, positions 4, 29 and 46 in d<sub>21</sub>, the inverted list with positions would be "d<sub>1</sub> | 21", "d<sub>9</sub> | 1, 19", "d<sub>21</sub> | 4, 29, 46".

→ and operation: Finds documents that contain all of k<sub>1</sub>, k<sub>2</sub> ... k<sub>n</sub>.

Intersection S<sub>1</sub> ∩ S<sub>2</sub> ∩ ... ∩ S<sub>n</sub>

or operation : documents that contain  
one of  $k_1, k_2 \dots k_n$   
union,  $S_1 \cup S_2 \cup \dots \cup S_n$ .

→ not operation: documents that contain a specified keyword e.g.

## Merging of intersection/Union.

# Measuring Retrieval Effectiveness

→ Information-retrieval Systems save space by using index structures that supports only appropriate retrieval, result in

→ false negative (false drop) - some relevant documents may not be retrieved

→ False Positive → Some irrelevant data may be retrieved.

→ For many applications a good index should not permit any false drops, but may permit a few false positives.

## → Recall vs Precision

1. Can increase recall by retrieving many documents (due to a low level of relevance ranking), but may irrelevant documents would be fetched, reducing precision.

→ Measures of retrieval effectiveness  
Recall as a function of number of documents fetched over Precision as a function of recall

Eg Precision of 75% at recall of 50%  
and 60% at a recall of 75%

Problem is which documents are actually relevant and which are not.

## → Crawling and Indexing the web

→ web crawlers are programs that locate and gather information on web.

→ They recursively follow the hyperlinks present in the known documents to find the other documents.

→ Starting from an initial set

→ The fetched documents

- Handled over to an indexing system
- Can be discarded after indexing as a cached copy.

disadvantage:- crawling the entire web

- would take a very large amount of time.
- Search engines typically cover only a part of the web not all of it.
- Take months to perform a single crawl.

→ Crawling is done by multiple processes on multiple machines running in parallel.

- Set of links to be crawled stored in a database.
- New links found in crawled pages are added to this set to be crawled later.

→ Indexing process also runs on multiple machines.

- Creates a new copy of index instead modifying old index.

- old index is used to answer queries.
  - After a crawl is "completed" new index becomes "old" index
- multiple machines used to answer queries.
- indices may be kept in memory.
  - Queries may be routed to different machines for load balancing.

→ Information Retrieval: Beyond Ranking of Pages

1. Diversity of Query Results.

2. Information Extraction Systems. Convert information from textual form to a more structured form.

3. Question Answering :- Information retrieval systems focus on finding documents relevant to a given query. However, the answer to a query may lie in just one part of a document, or in small parts of several documents.

→ Question answering systems attempt to provide direct answers to questions by users.

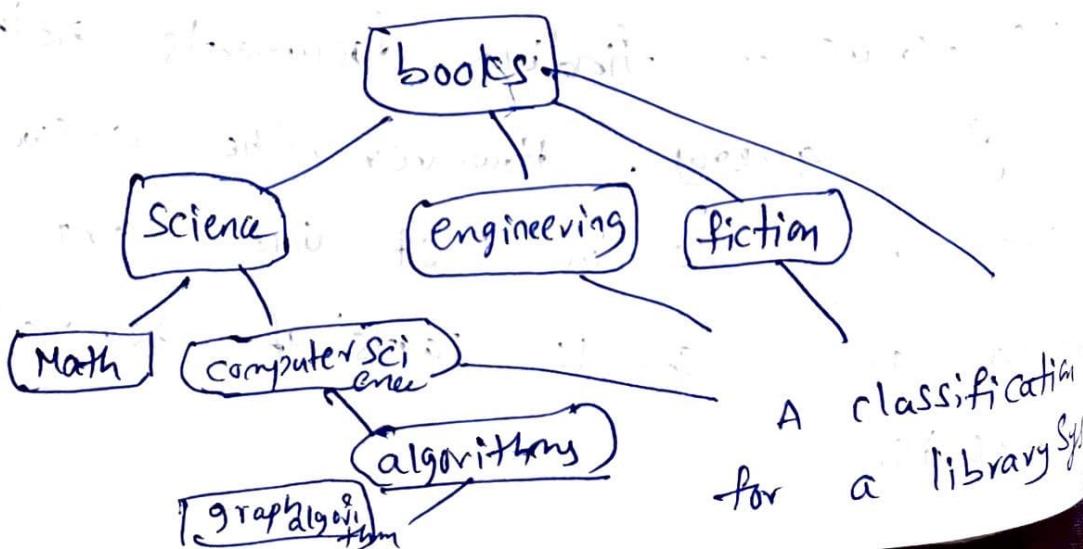
#### 4. Querying structured data:

##### Directories and categories

1. Storing related documents together
  - a. library facilitates browsing.  
Users can see not only requested document but also related ones.

→ Browsing is facilitated by classification system that organizes logically related documents together.

→ To keep related books close together libraries use a classification hierarchy



A classification for a library system

- A directory is simply a classification DAG structure.
  - Each leaf of the directory stores links to documents on the topic represented by the leaf.
  - Internal nodes may also contain links.
  - Documents can reside in multiple places in a hierarchy in an information retrieval system. Since physical location is not important.
  - classification hierarchy is thus Directed Acyclic Graph (DAG)
- A classification DAG for a Library information Retrieval System

