

Information Retrieval (IR)

Esther Pacitti

Partage de Données a Grande Echelle

Master 2 Informatique

Scenario



Information Need

- Example of an information need in the context of the world wide web:
- Find all documents (information!) about universities in India that
 - (1) offer schools of computer science and management degrees and
 - (2) offers data mining lectures. The information (the document!) should include full curriculum, fees, student campus, e-mail and other contact details.
- Formal representation of an information need = Query

Informal Definition

- Representation, storage, organisation and access of information (information items, information objects, documents). Not tuples.
- Find relevant (useful) information



Web Search Engine

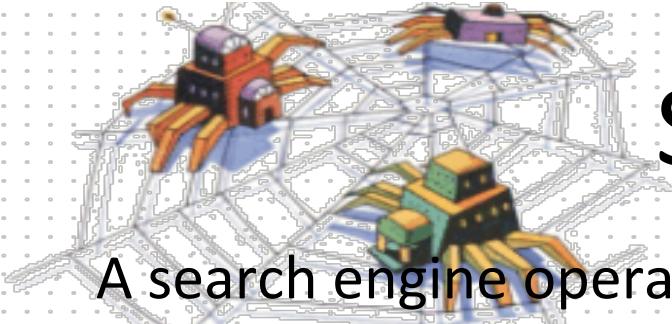
Internet search engines are special sites on the Web that are designed to help people find information stored on other sites. There are differences in the ways

various search engines work, but they all perform three

basic tasks:

1. They search the Internet or select pieces from the Internet based on important keywords.
2. They keep an index of the words they find and where they find them.
3. They allow users to look for words or combinations of words found in that index.

**WORLD
WIDE
WEB**



Search Engine

A search engine operates, in the following order

Crawling

Follow links to find information

Indexing

Record what words appear where

Ranking

What information is a good match to a user query? What information is inherently good?

Displaying

Find a good format for the information

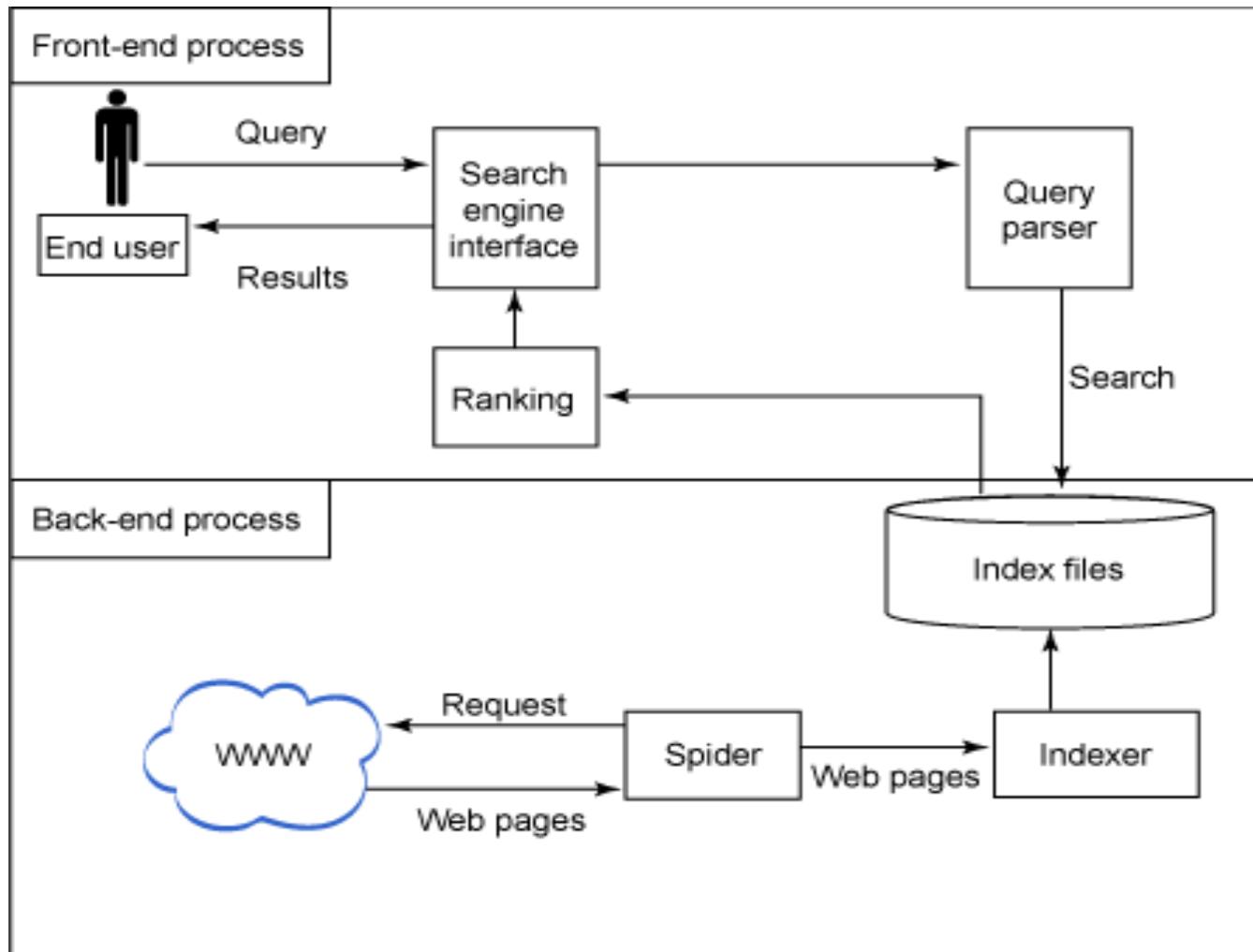
Serving

Handle queries, find pages, display results

Definitions

1. **Spiders:** To find information on the hundreds of millions of Web pages that exist, a search engine employs special software robots, called **spiders**, to build lists of the words found on Web sites.
"Spiders" take a Web page's content and create key search words that enable online users to find pages they're looking for.
2. **Crawling:** When a spider is building its lists, the process is called **Web**. In order to build and maintain a useful list of words, a search engine's spiders have to look at a lot of pages.
3. **Indexing:** For fast accessing of data.
4. **Meta tag:** The contents of each page are then analyzed to determine how it should be **indexed** (for example, words are extracted from the titles, headings, or special fields).

Web Search Engine



Informal Definition

- Goal of an IR system - RECALL
 - Retrieve all relevant documents.
- Goal of an IR system - PRECISION
 - Retrieve the most relevant documents.
- Goal of an IR system:
 - Retrieve as few non-relevant documents as possible.
 - Retrieve relevant documents before non-relevant documents.

Recall and Precision

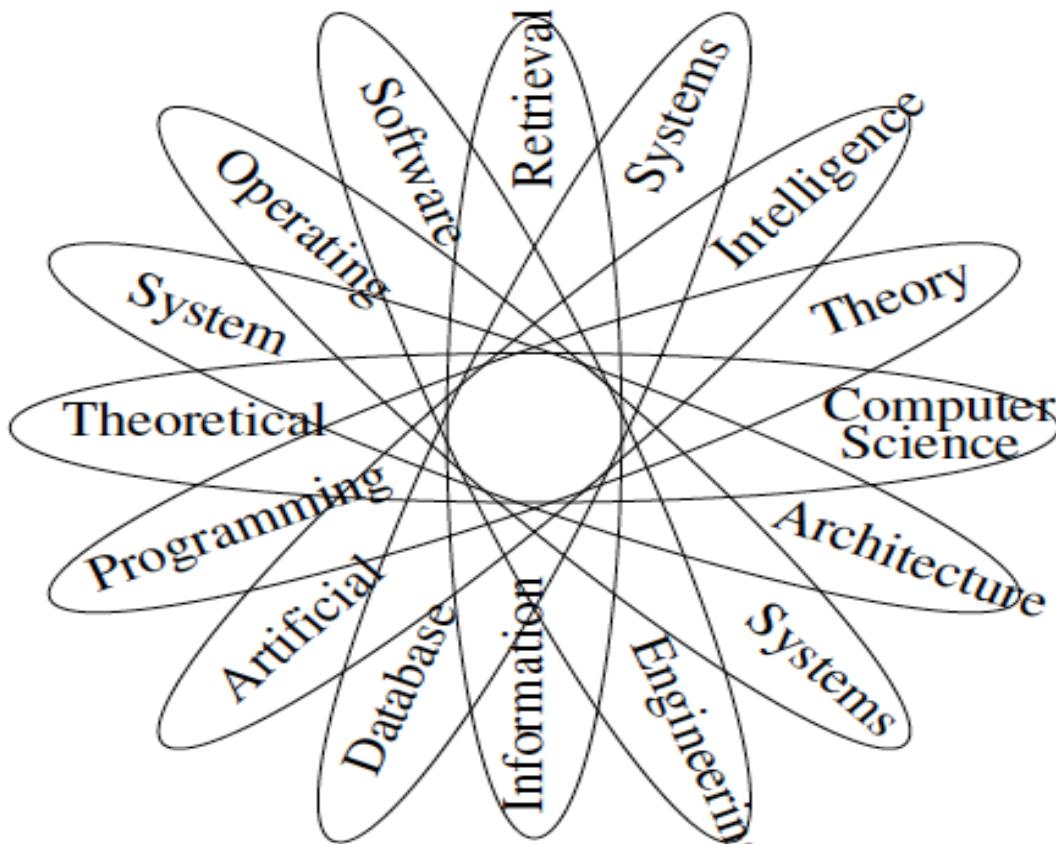
- Precision is the fraction of the documents retrieved that are relevant to the user's information need.

$$\text{precision} = \frac{|\{\text{relevant documents}\} \cap \{\text{retrieved documents}\}|}{|\{\text{retrieved documents}\}|}$$

- Recall is the fraction of the documents that are relevant to the query that are successfully retrieved.

$$\text{recall} = \frac{|\{\text{relevant documents}\} \cap \{\text{retrieved documents}\}|}{|\{\text{relevant documents}\}|}$$

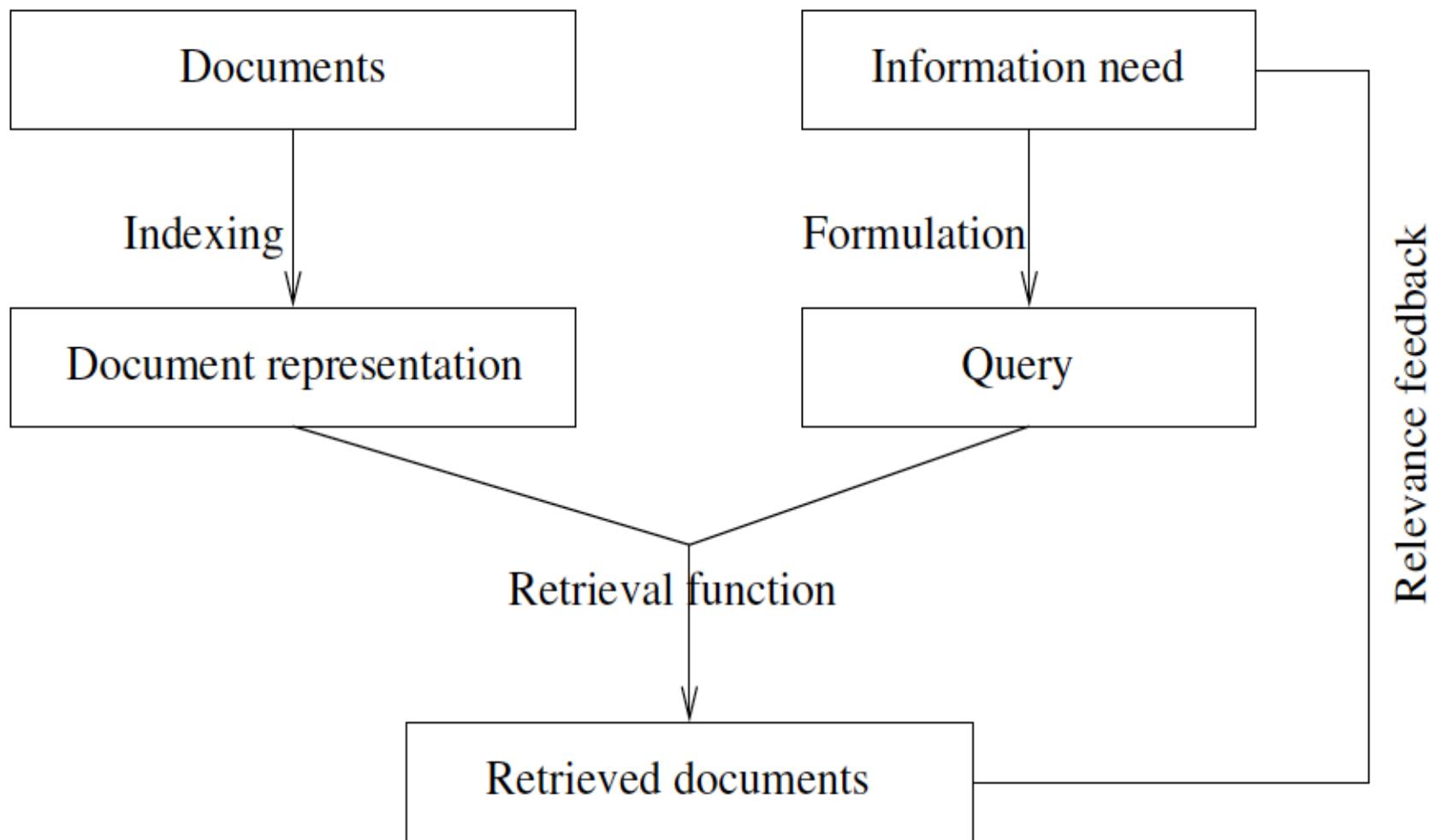
Information Retrieval in Computer Science



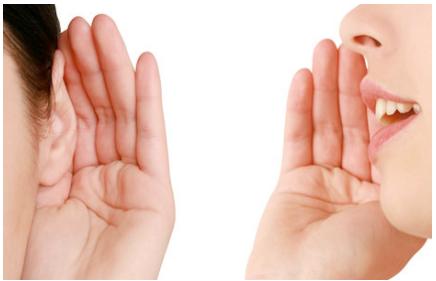
Information Retrieval x Data Retrieval

	Information Retrieval	Data Retrieval
Matching	vague	exact
Model	probabilistic	deterministic
Query language	natural	artificial
Query specification	incomplete	complete
Items wanted	relevant	all (matching)
Error handling	insensitive	sensitive

Model



Information Deluge – BIG DATA



Natural Language



E-mails



Movies



News



Web Sites Pages



Photos



Documents

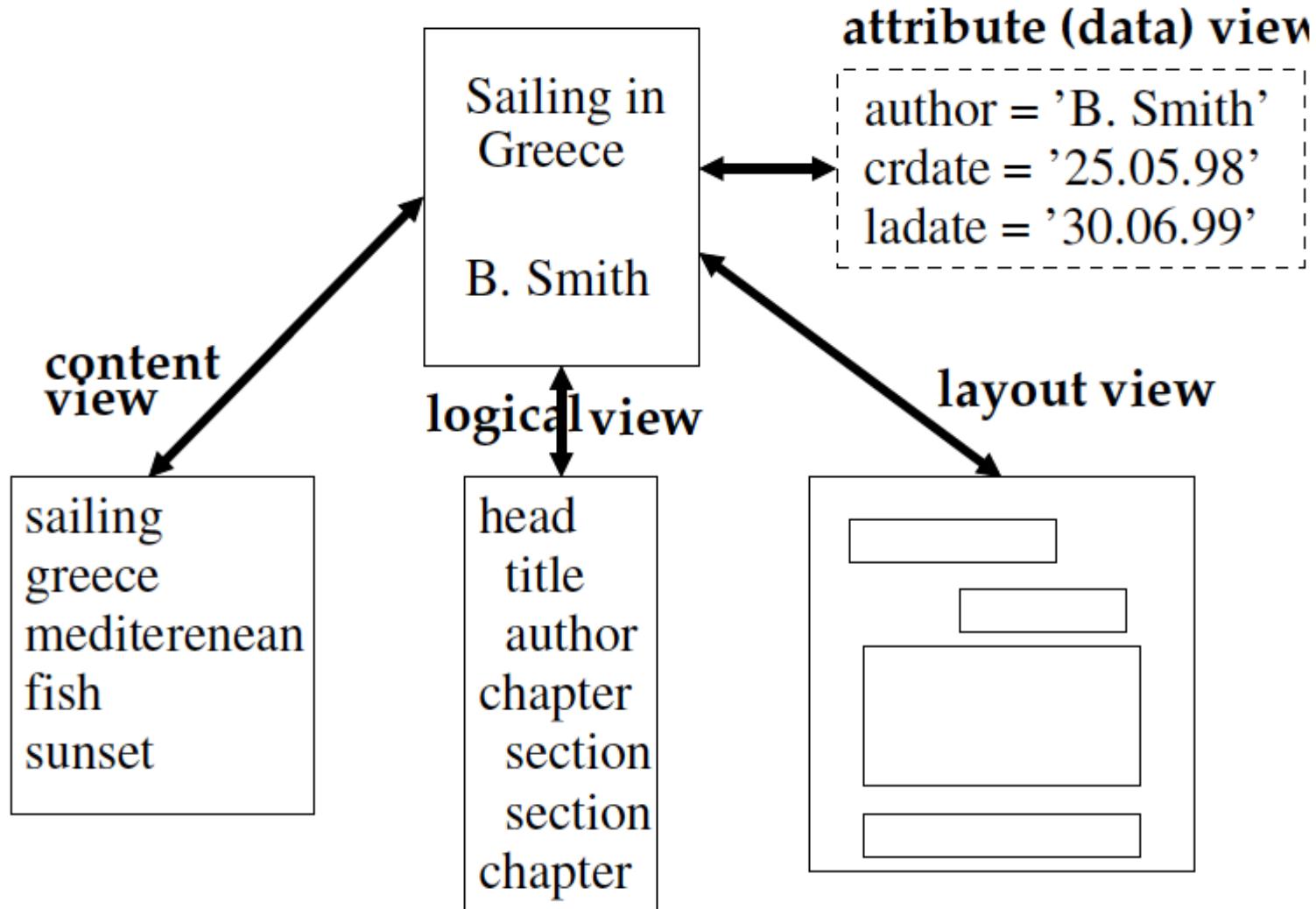
- Unit of retrieval and Free text:
 - text, strings of characters from an alphabet
 - natural language:
 - newspaper articles, journal paper, dictionary definition, e-mail messages
 - size of documents:
 - arbitrary, newspaper article vs journal article vs e-mail
- Sub-document can also be a unit of retrieval (XML element, answer to a question)

History

- Manual IR in libraries: manual indexing; manual categorisation
- 70ies and 80ies: Automatic IR in libraries
- 90ies: IR on the web and in digital libraries
- Success factors: Response time, coverage, interactivity, low (no!) costs, precision-oriented

precision \approx correctness, recall \approx completeness

Document Indexing



Document Indexing – Vector Space Model

- Language used to describe documents and queries
 - Index terms
- **Linear Algebraic** model for representing text documents

Vector Space Model

- Documents and queries are represented as vectors.
- Each dimension corresponds to a separate term. If a term (or keyword) occurs in the document, its value in the vector is non-zero.

$$\begin{aligned}d_j &= (w_{1,j}, w_{2,j}, \dots, w_{t,j}) \\q &= (w_{1,q}, w_{2,q}, \dots, w_{t,q})\end{aligned}$$

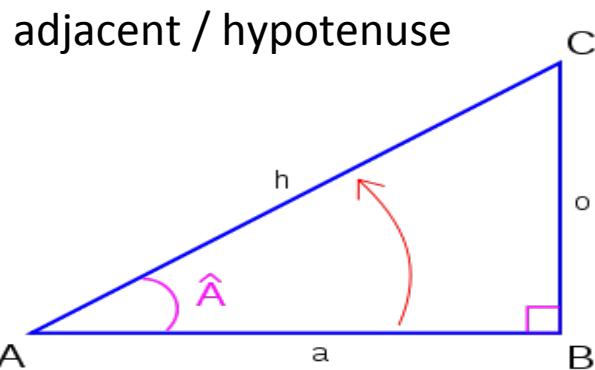
- TF-IDF (Term Frequency Inverse Document Frequency)
 - computes the values of the terms (or keywords) weights taking into account the frequency of the term in a document and in the whole corpus.

Relevance Ranking - Similarity

- Useful to compute de similarity between a query and a document or among documents.
- Compares the deviation of angles between each document vector and the original query vector, where the query is represented as the same kind of vector as the documents.
- It is easier to calculate the cosine of the angle between the vectors, instead of the angle itself.

$$\cos \theta = \frac{\mathbf{d}_2 \cdot \mathbf{q}}{\|\mathbf{d}_2\| \|\mathbf{q}\|}$$

Cosinus



$$\text{similarity} = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^n A_i \times B_i}{\sqrt{\sum_{i=1}^n (A_i)^2} \times \sqrt{\sum_{i=1}^n (B_i)^2}}$$

TF-IDF Hypothesis

Frequency of words is a measurement of a word significance.

... a measurement of the power of a word to discriminate documents by their content ...

Stop Words

- is, a, the, or, and, ...
- not?
- other?
- Stop-word list often defined manually.
- Reduction: between 30 and 50 per cent.

Index Term Weighting

- **Exhaustivity**
 - number of different topics indexed
 - importance of term in a document
 - high exhaustivity: high recall and low precision
- **Specificity**
 - ability of the indexing technique to describe topics precisely
 - number of documents to which a term is assigned in a collection
 - related to the distribution of index terms in collection
 - high specificity: high precision and low recall

Index Term Weighting

- Index term weighting
 - index **term frequency**: occurrence frequency of a term in document
 - **document frequency**: number of documents in which a term occurs

TF-IDF

$$\text{weight}(t, d) = \text{tf}(t, d) \times \text{idf}(t)$$

N	number of documents in collection
$n(t)$	number of documents in which term t occurs
$\text{idf}(t)$	inverse document frequency of term t
$\text{occ}(t, d)$	occurrence of term t in document d
t_{\max}	term in document d with highest occurrence
$\text{tf}(t, d)$	term frequency of t in document d

$$\text{tf}(t, d) := \frac{\text{occ}(t, d)}{\text{occ}(t_{\max}, d)}$$

$$\text{idf}(t) := \log \frac{N}{n(t)}$$

Example

The screenshot shows a Netscape browser window with the following details:

- Menu Bar:** File, Edit, View, Go, Communicator, Help.
- Toolbar:** Back, Forward, Reload, Home, Search, Netscape, Print, Security, Shop, Stop.
- Address Bar:** Bookmarks, Location: <http://www.ourcivilisation.com/sizes/chap7.htm>, What's Related.
- Page Content:**
 - A Study Of Social Decline by Philip Atkinson (1909)**
 - Underline Signs Of Social Deterioration**
 - Technology Making It Worse**
 - "When the life of people is immoral, and their relations are not based on love, but on egoism, then all technical improvements, the increase of man's power over nature, steam, electricity, the telegraph, every machine, gunpowder, and dynamite, produce the impression of dangerous toys placed in the hands of children."**
 - Leo Tolstoy (1828 - 1910)**
 - Making Things Worse Not Better**
 - Technology is the artificial enhancement of human power. It should make us stronger and smarter, however our demented community is discovering that it now has the opposite effect. Nuclear power has terrified and paralysed its creators, while the improved cleverness and flexibility of our machines have caused social chaos and economic stagnation.**
 - Australia - A Nuclear Free Zone (1990)**
 - All over this country are signs announcing the existence of nuclear free zones, erected by councils to announce the unpopularity of nuclear technology. Our nation has its nuclear power generating stations, or nuclear weapons, despite our growing need for energy and the inefficiency of our industry. Such concerns have been ignored by the electorate in favour of conventional (old) technologies. Any government that tried to reverse this situation would be deposed by a wave of public protests from a worried electorate.**
 - The End Of The Need To Work**
 - Benefits Of Technology Threatened By Our Attitude To Employment**
 - My job for fifteen years (1975 - 1990) had been to write computer programs to make people redundant. My efforts were not unique, throughout the western world an army of programmes have been working night and day to get rid of as many jobs as possible. Each job discarded meant improved productivity, and reduced costs. Because of our work, businesses throughout the world have become much more efficient and able to supply better goods and services, at a cheaper price. Nevertheless it would seem we have wasted our time. Industry and commerce can't utilise our improvements because there is no demand. There is no demand because people have no money. Nobody has any money because so many people are out of work.**
 - Luddite Riot**
 - The possibility of the loss of employment was first realised during the onset of the machine age. The invention and application of the steam engine heralded the industrial revolution. It dramatically extended the power and ability of the community. No longer was human strength and endurance the limiting factor in achievements. Machines could be constructed to work harder faster cheaper and more reliably than any team of people, however the initial implementation of machines caused mass unemployment and those who were left behind reacted.**

At the bottom of the browser window, the status bar displays "Document Done".

Nuclear 7	Computer 9
Poverty 5	Unemployment 1
Luddites 3	Machines 19
People 25	And 49

$$\begin{aligned} \text{Weight(machine)} &= \\ &19/25 \times \log(100/50) \\ &= 0.76 \times 0.3013 = 0.228988 \end{aligned}$$

$$\begin{aligned} \text{Weight(luddite)} &= \\ &3/25 \times \log(100/2) \\ &= 0.12 \times 1.69897 = 0.2038764 \end{aligned}$$

$$\begin{aligned} \text{Weight(poverty)} &= 5/25 \times \log(100/2) = 0.2 \\ &\times 1.69897 = 0.339794 \end{aligned}$$

Inverted Index

- In computer science, an inverted index is an index data structure storing a mapping from content, such as words or numbers, to its locations in a database file. The purpose of an inverted index is to allow **fast full text searches**, at a cost of increased processing when a document is added to the database (Wikipedia)

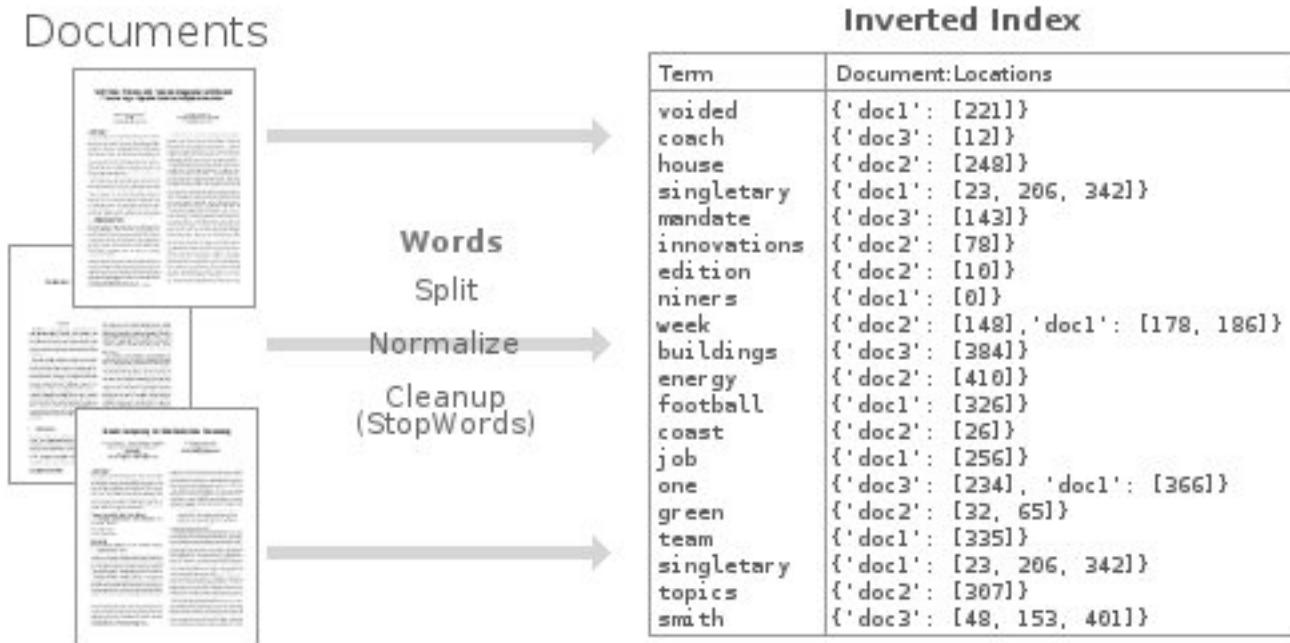
Documents:

D0 = "it is what it is"
D1 = "what is it"
D2 = "it is a banana"

Inverted Index:

"a": {2}
"banana": {2}
"is": {0, 1, 2}
"it": {0, 1, 2}
"what": {0, 1}

Inverted Index



Queries

- Simple queries: key words
- Best Match (query, documents)
 - Compare the terms in a document and query
 - Compute "similarity » (cosinus) between each document in the collection and the query based on the terms they have in common
 - Sorting the document in order of decreasing similarity with the query
 - The outputs are a ranked list and displayed to the user - the top ones are more relevant as judged by the system

Query Processing – Top-k processing

Top-k query processing
=

Finding k documents that have the highest overall grade wrt to the score:

$\text{COS}(Q, D_i)$, considering the whole corpus of documents.

In other words inspect the inverted index using for instant a Fagin's Algorithm

Top-k processing

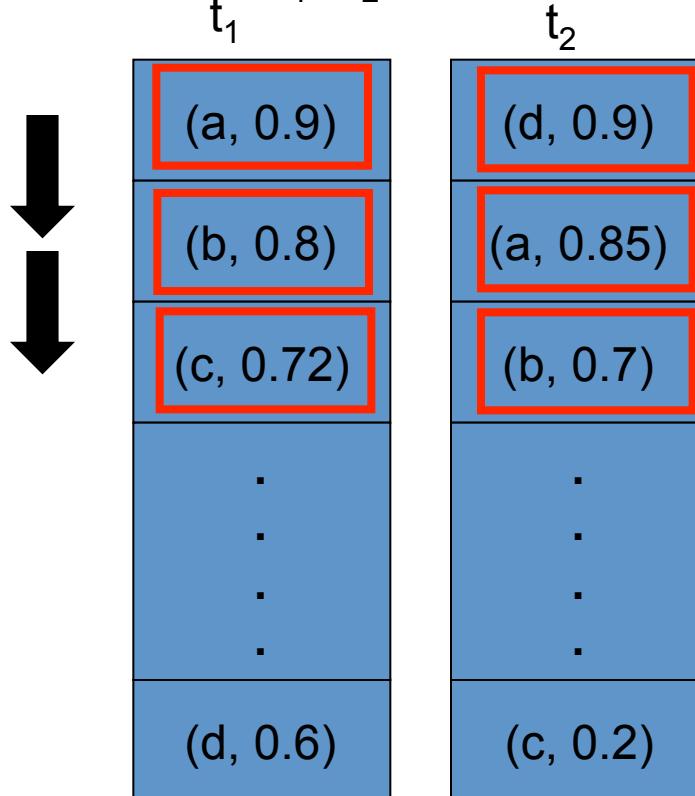
Input: inverted index, query vector, similarity measure (cosinus)

Output: top-k documents

Example – Fagin’s Algorithm

STEP 1

- $k = 3$ and $Q = t_1, t_2$
- Read $tf-idf$ weights (w_i) from every term in the sorted inverted lists (in parallel)
- Stop when k documents have been seen **in common from all lists**
- Ex: Return the the $k=3$ most relevant documents (a, b, c, etc) for the query and $Q = t_1, t_2$



Docs	w_1	w_2	$\text{Cos}(Q, d_i)$
a	0.9	0.85	
d		0.9	
b	0.8	0.7	
c	0.72		

Inverted Indexes are ordered in decreasing order (monotonicity)

Example – Fagin’s Algorithm

STEP 2

- Random access to find missing term weights.

t_1	t_2
(a, 0.9)	
(b, 0.8)	
(c, 0.72)	
.	
.	
.	
.	
(d, 0.6)	
	(d, 0.9)
	(a, 0.85)
	(b, 0.7)
	.
	.
	.
	(c, 0.2)

Docs	w_1	w_2	$\text{Cos}(Q, d_i)$
a	0.9	0.85	
d	0.6	0.9	
b	0.8	0.7	
c	0.72	0.2	

Example – Fagin’s Algorithm

STEP 3

- Compute the Cosinus score of the seen documents.
- Return the k highest scored documents.

t_1	t_2
(a, 0.9)	(d, 0.9)
(b, 0.8)	(a, 0.85)
(c, 0.72)	(b, 0.7)
.	.
.	.
.	.
.	.
(d, 0.6)	(c, 0.2)

Docs	w_1	w_2	$\text{Cos}(Q, d_i)$
a	0.9	0.85	0.85
d	0.6	0.9	0.6
b	0.8	0.7	0.7
c	0.72	0.2	0.2

Advantages

- Simple model based on linear algebra
- Term weights not binary
- Allows computing a continuous degree of similarity between queries and documents
- Allows ranking documents according to their possible relevance
- Allows partial matching