

Information Retrieval (Search)

IR

Artificial Intelligence → IR

Information Retrieval

- **Search**
- Using a computer to find relevant pieces of information
- **Text search**
- Idea popularized in the article *As We May Think* by Vannevar Bush in 1945

Where (or for what) do you do text search?

- World Wide Web
 - Using, e.g., Google, Yahoo
- Library catalog
- Personal (desktop) search
 - Email, files
- Within a document
 - Search-n-replace a word
- Specific domain/database
 - Medline (free)
 - Westlaw (for a fee)

Terminology

- **Query**
 - What you tell the computer to look for
- **Document**
 - What you are hoping to find
 - A webpage that contains the info you're after
 - A specific file on your computer
 - A specific email in your mail box

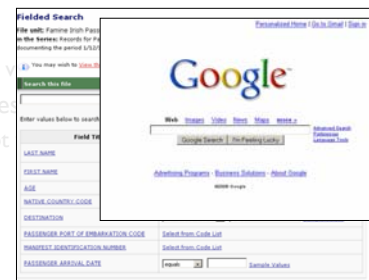
Type of search

- Flat text
 - Query: robot vision
- Quoted phrases
 - Query: "robot vision"



Type of search

- Flat text
 - Query: robot vision
- Quoted phrases
 - Query: "robot vision"
- Fielded search

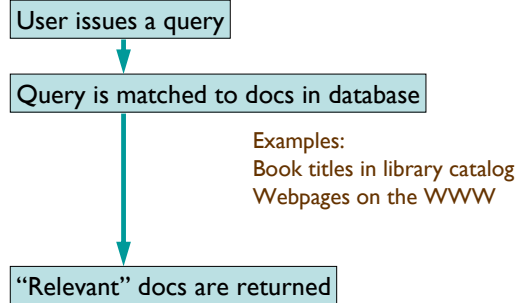


Type of search

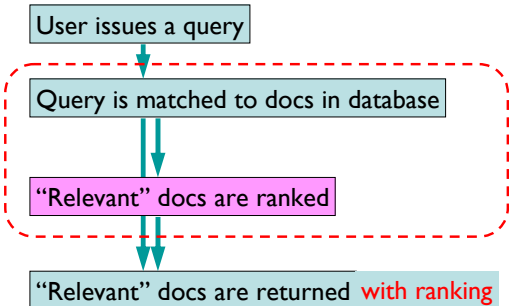
- Flat text
 - Query: robot vision
- Quoted phrases
 - Query: "robot vision"
- Fielded search
- Boolean operators
 - Query: flu and swine not human

The screenshot shows a search interface with three sections. Each section has a 'Category' dropdown, a 'Field Name' dropdown, an 'Operator' dropdown, and a 'Keyword(s)' text box. The first section has 'environmental conditions' as the category and 'exposure medium' as the field name. The second section has 'materials and substances' as the category and 'material or substance name' as the field name. The third section has 'statistical' as the category and 'comparison rate' as the field name. Boolean operators (AND, OR, NOT) are also visible between the sections.

The process



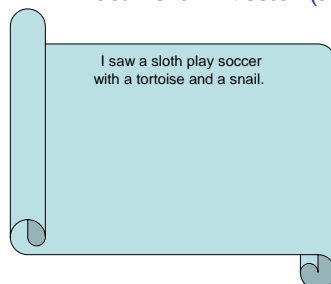
The process



Finding and comparing documents

- The **vector space model** is one method that performs a ranked search
- Represent a document as a **vector**, i.e., a list of individual words
 - Represent the query as a vector
 - Compare the two vectors mathematically

Document → Vector (simple version)



a | and | I | play | saw | sloth | snail | soccer | tortoise | with

Compare document with query

Document: a | and | I | play | saw | sloth | snail | soccer | tortoise | with

Query: shell | tortoise

1 match

Compare document with query

Document 1:	a	and	I	play	saw	sloth	snail	soccer	tortoise	with	1 match
Document 2:	birds	blue	fly	in	sky	the					0 match
Document 3:	blue	field	found	in	jewelry	shell	soccer	tortoise			2 matches
Query:	shell	tortoise									

Ranked search result:
Document 3
Document 1

Vector space model

- Vectors are very, very long
 - We say it is a “high-dimensional” problem
 - # dimensions = size of vocabulary
- Very computationally intensive
- Any other problems?



Variation: term weighting

Some words are more discriminating than others. E.g., “the” appears in just about every document

- **Term frequency (TF)**
 - E.g., The more times “Potter” is in the doc, the more likely the doc is about him
- **Inverse document frequency (IDF)**
 - The more documents there are containing a certain word, the less likely that word is important

Use term frequency to improve search

Document 1:	a	and	I	play	saw	sloth	snail	soccer	tortoise	with	
	3	1	1	1	1	1	1	1	1	1	Score: 1
Document 2:	birds	blue	fly	in	sky	the					
	1	2	1	1	1	1					Score: 0
Document 3:	blue	field	found	in	jewelry	shell	soccer	tortoise			
	2	1	1	1	1	1	1	2			Score: 3
Query:	shell	tortoise									
	1	1									

Ranked search result:
Document 3
Document 1

Two different notations to encode the same vector information

a	and	I	play	saw	sloth	snail	soccer	tortoise	with
3	1	1	1	1	1	1	1	1	1

{3:a 1:and 1:I 1:play 1:saw 1:sloth 1:snail 1:soccer 1:tortoise 1:with}

Preparing documents for vector space model

- **Stemming**
 - Potter’s = Potters = Potter
- **Stop-words**
 - Ignore words like “the”, “of”, ...
- **Use statistical properties of text**
 - E.g, Data from Jamie Callan’s Characteristics of Text, 1997 (Sample of 19 million words)

Commonest fifty words

<i>f</i>	<i>f</i>	<i>f</i>
the 1,130,021	from 96,900	or 54,958
of 547,311	he 94,585	about 53,713
to 516,635	million 93,515	market 52,110
a 464,736	year 90,104	they 51,359
in 390,819	its 86,774	this 50,933
and 387,703	be 85,588	would 50,828
that 204,351	was 83,398	you 49,281
for 199,340	company 83,070	which 48,273
is 152,483	an 76,974	bank 47,940
said 148,302	has 74,405	stock 47,401
it 134,323	are 74,097	trade 47,310
on 121,173	have 73,132	his 47,116
by 118,863	but 71,887	more 46,244
as 109,135	will 71,494	who 42,142
at 101,779	say 66,807	one 41,635
mr 101,679	new 64,456	their 40,910
with 101,210	share 63,925	

Finding documents

- Brute-force approach?
 - Look through every single document every time you have a query
- Efficient way?
 - Make an index

Evaluating IR methods

- Precision
 - How many of the returned documents are relevant?
- Recall
 - How many of the relevant documents are returned?
 - Cannot be the sole criterion in evaluation
- Fall-out
 - How many of the non-relevant documents are returned?
- Can combine these criteria

Web Search

Artificial Intelligence → Information Retrieval → Web Search

What's special about web search?

- Hyperlinks
- Size—scalability issues
- Dynamic content
- Untrained users
- Economic model (advertising)

“Crawling” the web

- Following the links to determine the link structure
- What are some issue and considerations?
 - Broken links, timeouts, ... cause failures
 - Update frequency
 - Coverage, duplicate detection
 - Legal issues (owners don't want their pages indexed)
 - Advertising links
 - Types of content
 - ...

Web search through link analysis

- Find relevant webpages by analyzing the **link structure**, not by the content
- Most famous algorithm is PageRank
- There are other kinds of link analysis
 - E.g., citation analysis—count the number of references to individual research papers (CiteSeer)

PageRank

- Important part of Google's success (although most search engines use something like PageRank nowadays)
- Rank pages not just by how **relevant** they are, but also by how **important** they are
- Estimate importance by considering a link as a vote
 - The more pages link to you, the more important you are

The PageRank idea

- Many pages link to my page
- → there are many ways to get to my page
- → the probability of getting to my page is high
- → I am important

Start from a random page

Repeat:

Click on a random link → go to that page

Do a large number of such simulations.
Where do you end up after a large number of clicks?

For each page, how many visitors end up there?
→ Give the ranks by importance of all the pages

Google can combine this with:
TF
IDF
voodoo
...

Web search is big business! Advertising

- **The advertiser**
 - Buy words (e.g., “digital camera”)
 - Then if my search has those words, I'll see their ad
- **The webmaster**
 - I want to put ads on my site (revenue)
 - I give space on my site to a search engine company and they fill it with relevant ads
- **The user**
 - Sees sponsored results