

INFORMATION RETRIEVAL

- Problem Formulation
- Issues in IR
- Term-Document Incidence Matrix
- Inverted Index
- Ranked Retrieval

IMDAD ULLAH KHAN

Information Retrieval (IR)

Information Retrieval (IR) is finding material (usually documents) of an unstructured nature (usually text) that satisfies an information need from within large collections (usually stored on computers).

Applications:

- Web search, E-mail search, File search in your laptop, corporate knowledge bases, legal information retrieval
- Search and communication are most popular uses of the computer
- Search is one of the the most fundamental computational problem

The Information Retrieval Problem

- Given a collection of documents
Assume it is a static collection
- **Goal:** Retrieve documents with information that is *“relevant”* to the user’s need and helps the user complete a task

- How to “effectively” describe information resources (documents containing information)?
 - Feature extraction and Representation
 - Organization and Storage
 - Indexing
- How to find the “appropriate” information resources?
 - Accessing
 - Filtering
 - Retrieving
 - Ranking
 - Presenting

Information Retrieval: Keywords

- Automated IR systems designed in 60's to search news articles, scientific papers, patents, legal abstracts for **keyword queries**
- Keywords are short and inexpressive
- Problem of **synonymy** (different words with same meaning)
 - “soccer” would not retrieve documents containing the word “football”
- Problem of **polysemy** (same words with different meanings)
 - “Lie” could mean to lay down and to tell something untruthful
 - “fair” could be an event or could refer to skin color, or as in “just”
- Early IR system had two distinct features
 - 1 **Retrieval done by experts** (reference librarians, patents attorneys)
 - 2 Search space (knowledge base) were collection of **expert-written** docs

Now everyone is an author and everyone is a searcher

Information Retrieval: Models

We briefly discuss two main information retrieval models

- Boolean Retrieval
- Ranked Retrieval

The goal is to set the stage for **web information retrieval**, where there are many other issues

Boolean Retrieval

Boolean Retrieval

- Boolean retrieval model pose any query in the form of a Boolean expression of terms
- It combine sentences/documents with operators AND, OR, and NOT
- It views each document as a set of words/terms
- It is precise: document matches a condition or not
- It was the primary commercial retrieval tool for 3 decades
- Many search systems still in use are Boolean
 - Email, library catalog, Mac OS X Spotlight
- Main techniques used in Boolean retrieval are
 - Term Document Incidence Matrix
 - Inverted Index

Term-Document Incidence Matrix

- Store the collection of documents in Boolean form ▷ set of words
- Term-document incidence matrix, M
- A row for each term and a column for each document

$$M(t, d) = \begin{cases} 1 & \text{if doc } d \text{ contains term } t \\ 0 & \text{otherwise} \end{cases}$$

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	1	1	0	0	0	1
Brutus	1	1	0	1	0	0
Caesar	1	1	0	1	1	1
Calpurnia	0	1	0	0	0	0
Cleopatra	1	0	0	0	0	0
mercy	1	0	1	1	1	1
worser	1	0	1	1	1	0
...						

Term-Document Incidence Matrix

Retrieve Shakespeare's novel(s) that have words

Brutus AND Caesar AND NOT Calpurnia

- Brute force: linear scan all docs for each query term
- Bitwise AND of vectors (rows) for Brutus, Caesar and NOT Calpurnia

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	1	1	0	0	0	1
Brutus	1	1	0	1	0	0
Caesar	1	1	0	1	1	1
¬Calpurnia	1	0	1	1	1	1
Cleopatra	1	0	0	0	0	0
mercy	1	0	1	1	1	1
worser	1	0	1	1	1	0
AND	1	0	0	1	0	0

- Returns two documents, “Antony and Cleopatra” and “Hamlet”

Term-Document Incidence Matrix: Limitations

- Consider $N = 10^6$ documents, each with about 1000 terms (tokens)
- 10^9 tokens at avg 6 bytes per token \rightarrow 6GB
- Assume there are $|\Sigma| = 500,000$ distinct terms in the collection
- Size of incidence matrix is then $500,000 \times 10^6$ (half a trillion bits)
- Generally, the term-document matrix is very sparse (contains no more than a billion 1's)

The Inverted Index

- Σ : set of all terms (also called lexicon, vocabulary)
- A **postings list** for each term in Σ
 - A list of document IDs in which the term occurs

Brutus → 1 → 2 → 4 → 11 → 31 → 45 → 173 → 174

Caesar → 1 → 2 → 4 → 5 → 6 → 16 → 57 → 132 → 179

Calpurnia → 2 → 31 → 54 → 101

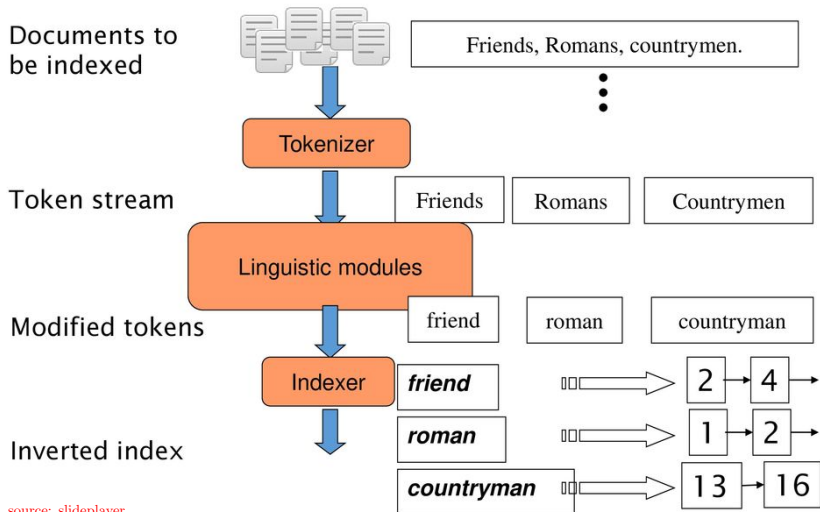
- Sort each posting and the list of posting for optimal processing
- Save length of lists at header

Brutus 8 → 1 → 2 → 4 → 11 → 31 → 45 → 173 → 174

Caesar 9 → 1 → 2 → 4 → 5 → 6 → 16 → 57 → 132 → 179

Calpurnia 4 → 2 → 31 → 54 → 101

The Inverted Index - Steps



source: slideplayer

The Inverted Index - Token sequence

- Sequence of (Modified token, Document ID) pairs

Doc 1

I did enact Julius
Caesar I was killed
i' the Capitol;
Brutus killed me.

Doc 2

So let it be with
Caesar. The noble
Brutus hath told you
Caesar was ambitious

Term	docID
I	1
did	1
enact	1
julius	1
caesar	1
I	1
was	1
killed	1
i'	1
the	1
capitol	1
brutus	1
killed	1
me	1
so	2
let	2
it	2
be	2
with	2
caesar	2
the	2
noble	2
brutus	2
hath	2
told	2
you	2
Caesar	2
was	2
ambitious	2

The Inverted Index - Sort

- Sort by terms
- then by doc-ID

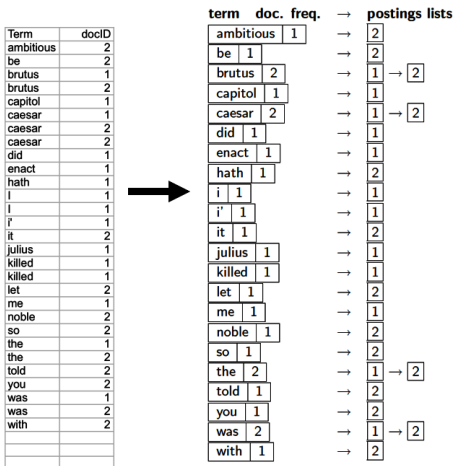
Term	docID
I	1
did	1
enact	1
julius	1
caesar	1
I	1
was	1
killed	1
i'	1
the	1
capitol	1
brutus	1
killed	1
me	1
so	2
let	2
it	2
be	2
with	2
caesar	2
the	2
noble	2
brutus	2
hath	2
told	2
you	2
caesar	2
was	2
ambitious	2



Term	docID
ambitious	2
be	2
brutus	1
brutus	2
capitol	1
caesar	1
caesar	2
caesar	2
did	1
enact	1
hath	1
I	1
I	1
i'	1
it	2
julius	1
killed	1
killed	1
let	2
me	1
noble	2
so	2
the	1
the	2
told	2
you	2
was	1
was	2
with	2

The Inverted Index - Dictionary And Postings

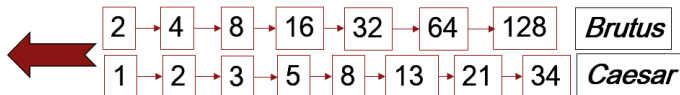
- List document entries for each term
- Split into Dictionary and Postings
- Add number of docs information



The Inverted Index - Query processing

Query: Brutus AND Caesar

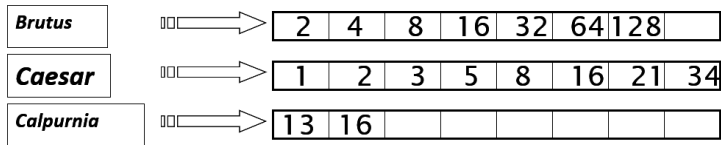
- Locate Brutus in the dictionary and retrieve its postings
- Locate Caesar in the dictionary and retrieve its postings
- “Merge” the two postings (to get intersection)
 - Walk through the two postings simultaneously, in time linear in the total number of postings entries
 - Let the list lengths be x and y , merge takes $O(x + y)$ operations (since postings are sorted by docID)



The Inverted Index - Query Optimization

Query: Brutus AND Calpurnia AND Caesar

- What is the best sequence of binary intersections to get $X \cap Y \cap Z$
- Process in order of increasing posting lengths
 - start with smallest pair of sets, then keep cutting further
 - this is why we kept document lengths in dictionary



Execute the query as (Calpurnia AND Brutus) AND Caesar

The Inverted Index - More general merges

Query: Brutus AND NOT Caesar

- Complement posting would be very long
- Can we still compute intersection in $O(x + y)$?

Query: (Brutus OR Caesar) AND NOT (Antony OR Cleopatra)

- Can we evaluate any set algebraic expression in “linear” time?

Query: “Punjab University” as a phrase

- Documents containing “University of Punjab” are not a match
- **Biword Indexes**: Indexes bigrams in the text as phrases
 - **Issues with Biword Indexes**: Index blowup due to bigger dictionary
- **Positional Indexes**: also cater for context
 - In the postings, for each term store the position(s) where it appears
<term, number of docs containing term;
doc1: position1, position2 ... ;
doc2: position1, position2 ... ;
etc >
 - **Issues with Positional Indexes**:
 - A positional index is 2 – 4 times as large as a non-positional index
 - Positional index size 35 – 50% of volume of original text

Limitations of Boolean Retrieval

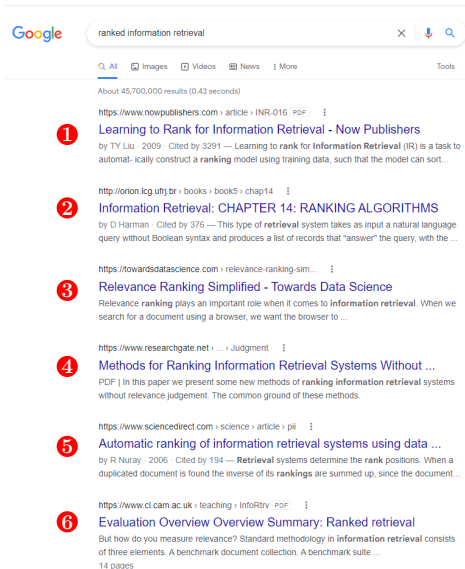
- Exact matching may retrieve too few (~ 0) or too many documents (~ 1000)
- Difficult to come up with a manageable query
 - AND gives too few, OR gives too many
- The burden is on the user to formulate a good Boolean query
- All terms are equally weighted
- Output is not ordered in useful fashion
- Complex query syntax (too much work for non-experts)

Ranked Retrieval

Ranked Retrieval

- Ranks documents by **relevance** to the query
- Return top k documents by relevance
- Allow for **free text queries**: Rather than considering query language of operators and expressions, it consider words of human language
- The size of the result is not an issue ($k \sim 10, 20, 100$)

Ranked Retrieval



Google ranked information retrieval

Q All Images Videos News More Tools

About 45,700,000 results (0.43 seconds)

- <https://www.nowpublishers.com/article/INR-016> PDF

Learning to Rank for Information Retrieval - Now Publishers

by TY Liu · 2009 · Cited by 3291 — Learning to rank for Information Retrieval (IR) is a task to automatically construct a ranking model using training data, such that the model can sort...
- <http://orion.lcg.ufrj.br/books/book5/chap14>

Information Retrieval: CHAPTER 14: RANKING ALGORITHMS

by D Harman · Cited by 376 — This type of retrieval system takes as input a natural language query without Boolean syntax and produces a list of records that "answer" the query, with the ...
- <https://towardsdatascience.com/relevance-ranking-sim...>

Relevance Ranking Simplified - Towards Data Science

Relevance ranking plays an important role when it comes to information retrieval. When we search for a document using a browser, we want the browser to ...
- <https://www.researchgate.net/.../Judgment>

Methods for Ranking Information Retrieval Systems Without ...

PDF | In this paper we present some new methods of ranking information retrieval systems without relevance judgement. The common ground of these methods.
- <https://www.sciencedirect.com/science/article/pii...>

Automatic ranking of information retrieval systems using data ...

by R Nuray · 2006 · Cited by 194 — Retrieval systems determine the rank positions. When a duplicated document is found the inverse of its rankings are summed up, since the document...
- <https://www.cl.cam.ac.uk/teaching/InfoRtriv/> PDF

Evaluation Overview Overview Summary: Ranked retrieval

But how do you measure relevance? Standard methodology in information retrieval consists of three elements. A benchmark document collection. A benchmark suite ...
14 pages

Scoring as Basis of Ranked Retrieval

- How to rank-order documents in collection with respect to a query?
- Assign a relevance score - say in $[0,1]$ - to each document
- e.g query with one term :
 - If the term does not occur in the document, score is 0
 - The more frequent the query term in document; the higher the score

Scoring Methods for Ranked Retrieval

Different measures to calculate relevance/similarity in ranked retrieval

- Exact Match Models
 - Unranked Boolean retrieval model
 - Ranked Boolean retrieval model (based on set/bag of words model)
- Best Match Model (Vector Space Model)
 - TF-IDF (and Cosine Similarity)

VECTOR SPACE MODEL

- Best-match models 'compute' the degree to which a document is relevant to a query
- It is expressed as $sim(q, d)$ (similarity between query and document)
- How to compute the similarity between q and d ??
- **Vector space model** represent both documents and queries by real vectors
- Similarity is evaluated in the vector-space
- e.g. cosine similarity between the vector representation of query and document

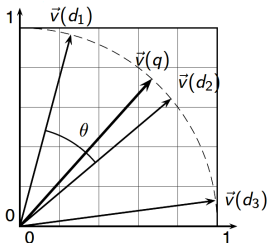
SUMMARY – VECTOR SPACE MODEL RANKING

- Represent the query by the TF-IDF vector
- Represent each document by the TF-IDF vector
- Compute cosine similarity between the query and document vectors
- Rank documents with respect to the query by cosine similarity
- Return top k to the user

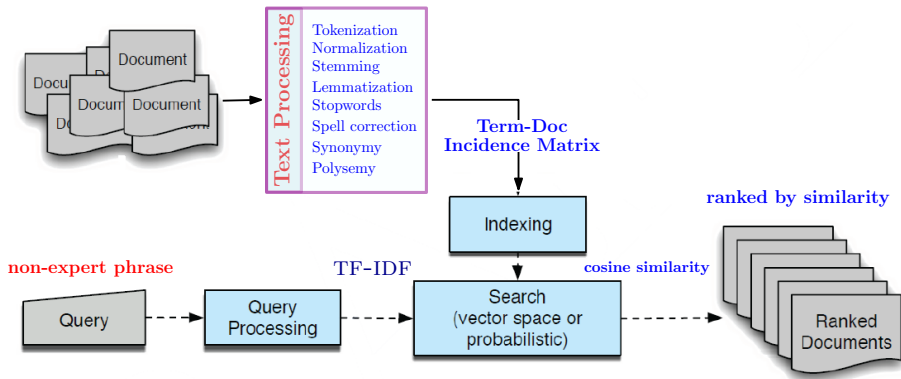
VECTOR SPACE MODEL

$$\cos(\mathbf{q}, \mathbf{d}) = \frac{\mathbf{q} \cdot \mathbf{d}}{\|\mathbf{q}\| \|\mathbf{d}\|} = \frac{\sum_{i=1}^{|\mathcal{V}|} q_i d_i}{\sqrt{\sum_{i=1}^{|\mathcal{V}|} q_i^2 \sum_{i=1}^{|\mathcal{V}|} d_i^2}}$$

- q_i is the TF-IDF weight of term i in the the query
- d_i is the TF-IDF weight of the term i in the document



Information Retrieval



VECTOR SPACE MODEL -EXAMPLE

- **Query:** “best car insurance”
- **Document:** “car insurance auto insurance”

word	query					document				product
	tf-raw	tf-wght	df	idf	weight	tf-raw	tf-wght	weight	n'lized	
auto	0	0	5000	2.3	0	1	1	1	0.52	0
best	1	1	50000	1.3	1.3	0	0	0	0	0
car	1	1	10000	2.0	2.0	1	1	1	0.52	1.04
insurance	1	1	1000	3.0	3.0	2	1.3	1.3	0.68	2.04

Key to columns: tf-raw: raw (unweighted) term frequency, tf-wght: logarithmically weighted term frequency, df: document frequency, idf: inverse document frequency, weight: the final weight of the term in the query or document, n'lized: document weights after cosine normalization, **product: the product of final query weight and final document weight**