

11-02-2022

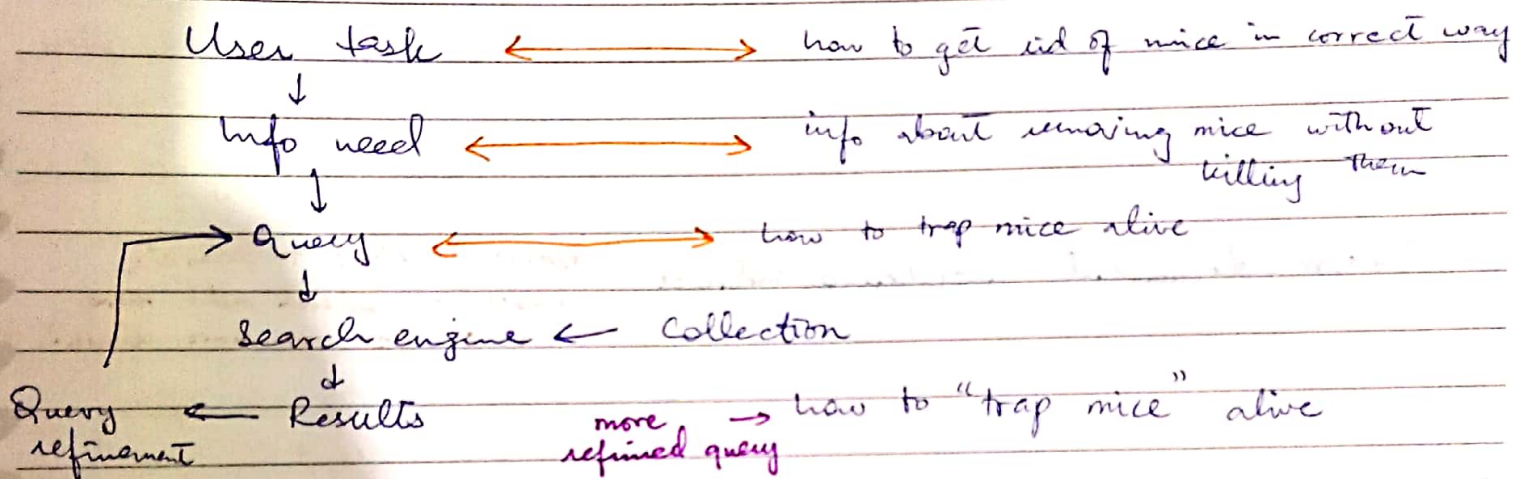
Date

Information Retrieval

- web crawlers

• extraction of essential required info from ~~data~~ unstructured data

• Info Retrieval vs. Relational Database



• Core concepts of IR

- Query Representation
- Document Representation
- Retrieval Model

— Precision: fraction of retrieved docs that are relevant to user's information need

— Recall: Fraction of relevant docs in collection that are retrieved

$$\frac{70}{100}$$

$$\frac{70}{150}$$

Information Retrieval: Challenges

- data is unstructured
- query is unstructured
- computers can't guess person's intent

Date February 14th, 2022

- inferring relevance and intent from data

Unstructured data in 1620: (example)

- slow
- NOT Calpurnia non trivial (include 2, exclude 1)
but queries are not always easy.
- 'find words ~~near~~ Romans near countrymen' not feasible, complex, 'near' is ambiguous.
- ranked retrieval (criteria according to which ranking is done)

Term-document incidence matrices

- makes it easier to search
- easy to understand add more plays, just one more column

1 — contains word

0 — play doesn't contain word

Incidence Vectors:

← not efficient for large data collection
will make extremely sparse matrix

Brutus 1 1 0 1 0 0 AND

Caesar 1 1 0 1 1 1 AND

NOT Calpurnia 1 0 1 1 1 1 = ← after complement of Calpurnia
1 0 0 1 0 0 ← plays containing Brutus and Caesar
but not Calpurnia

Solution: if a word is not present
in a doc. don't add/enter it.
but can't be a matrix now.

Date February 15th, 2022.

Inverted Index

backward index

works with modern IR systems

- each doc has a doc ID

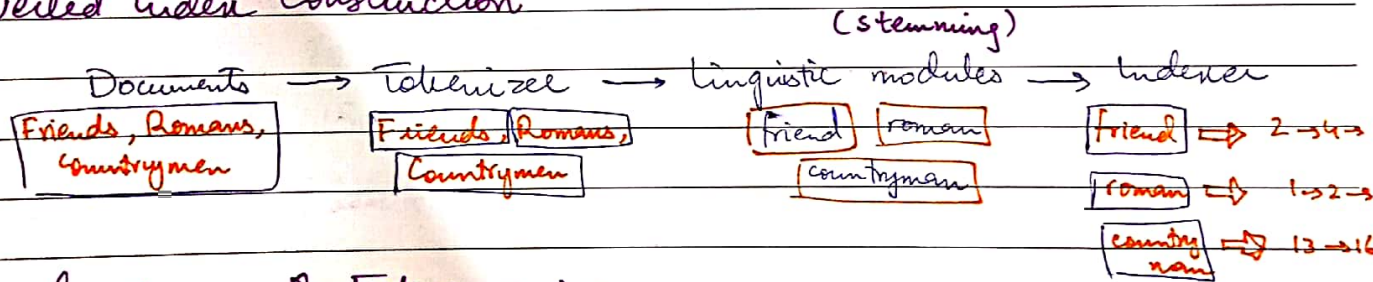
Dictionary $\left[\begin{array}{c} Brutus \\ Caesar \\ Calpurnia \end{array} \right] \begin{array}{c} 1 \\ 2 \\ 4 \end{array} \begin{array}{c} 11 \\ 31 \\ 45 \end{array} \begin{array}{c} 173 \\ 174 \end{array} \begin{array}{c} 179 \end{array}$ ← sorted

- fixed size array can't be used because it can't grow, some positions may be empty and wasted so we use Linked List.

- need variable size postings lists ← sorted by doc ID



Inverted Index Construction



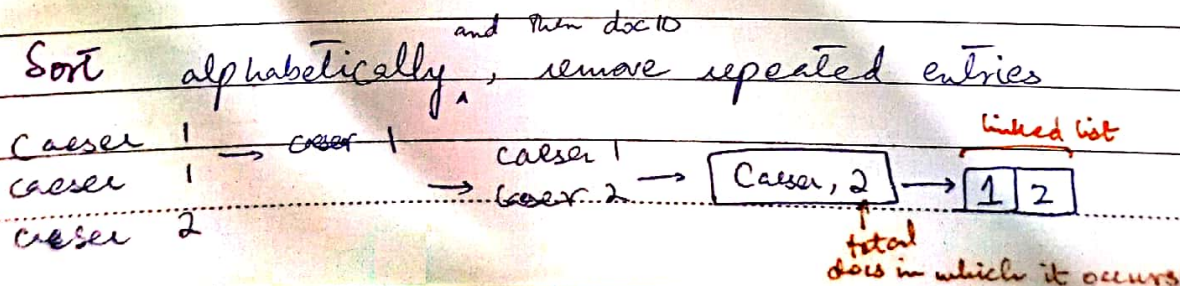
Initial stages of Text processing

- tokenization
- normalisation
- Stemming
- Stop words

① Indexer Steps : Token Sequence

add all words, even repeated ones

② Sort alphabetically, remove repeated entries



Date

- Sorted on Terms firstly so repeated terms come together ← takes linear time for computer to search

Dictionary + Postings

- multiple entries are merged in a single document
- split into dictionary + postings

term	doc. freq.	posting
ambitious	1	2
be	1	2

Issues of Storage: ^{pointers,} term, doc. freq. list in memory
postings in disk/permanent storage
(take it in memory when using)
need be

Query Processing w/ inverted index

AND: Brutus AND Caesar

retrieve postings and merge/intersection

len x	Brutus	<u>2</u>	4	<u>8</u>	16	<u>32</u>	64	128	
len y	Caesar	1	<u>2</u>	3	5	<u>8</u>	13	21	34
		↑							

$O(x+y)$
cuz sorted

set pointers on start, no match? take pointer to next ~~to~~ where number is small.

if not sorted then $x*y$
 $O(xy)$

Algorithm in slides (33 pg)

DIY

take 5 docs, make
inverted index, do
linear search, do
query search

Date February 22nd, 2022.

22-2-2022

Example: Westlaw Boolean Search

- What's statute of limitations in cases involving federal tort claims act?

⇒ LIMIT / 3 STATUTE ACTION / 3 FEDERAL / 2

word should have limit, followed by anything i.e. limited, limits, limitation
 ↓ within 3 words
 OR
 ↓ in same sentence
 TORT / 3 CLAIM

1p
same paragraph

- disjunction OR — conjunction AND (space)

- Brutus and Not Caesar: $O(x+y)$ can be done in linear time
different? add in list, same? pop

- Brutus or not Caesar: DIY.

Brutus \cup (not Caesar)

can be in linear time

- Merging combine 2 lists, OR.

AND: first consider list with min elements (check from freq. list)
 then take intersection with next min list then next.
 AND

Slide 4

* rhyme zone. conf
 Shakespeare
 write 5 search
 features you think
 it could do better

Date 24th February, 2022.

Phrase Queries

- "Stanford university" \leftarrow both words should come in this order like this.
- phrases can't be stored in inverted indexes. $\langle \text{term} : \text{docs} \rangle \times$
maybe Stanford & uni come in a document but inv. index doesn't tell if both occur together or not.

Solution 1: Biword indexes

excluding stop words from phrases	Stanford, 3	1, 3, 5	Stanford uni {1, 5}
	university, 5	1, 5, 11, 16, 23	
	Stanford uni, 1	5	\swarrow don't occur together \searrow occur together
	international cricket		
	cricket council		

Bigrams

Stanford university international
cricket council

- Stanford university palo alto \leftarrow longer queries can't be solved with biword indexes

Stanford uni AND uni palo AND palo alto
{1, 5, 7, 11} {2, 3, 5, 7} {5, 11, 23} \rightarrow {5} can give false positives as no assurance whether these 3 bigrams occur together in doc or not which we want.

- We need to have positions stored with each word.

Issues: ① false pos ② can be part of a compound strategy

Date

Solution 2: Positional Indexes

→ example in slides.

- not for biwords, only for single terms
- documents every occurrence in the doc.
- <term, no. of docs it occurs in>

doc1 ← 1: 7, 18 → position in doc

doc2 ← 2: ...

99 3427: ... >

- proximity queries: positional indexes can be used, biword indexes cannot
- positional index size depends on avg. document size.
- pos. index > 2-4 times than non-pos. index
- pos. index size 35-50% of volume of original text → pos. index is faster than linear search

Combination Schemes

- Positional + Bigram indexes in a hybrid scheme
- Study / Paper showed it takes $\frac{1}{4}$ th of time (faster) but 26% more space.

watch videos

Date

March 1st, 2022.

RCV1 dataset \leftarrow UCI repo

tokens \leftarrow words + stop words
terms \leftarrow only words.

avg. no. bytes per token $<$ avg. no. of bytes per Term

\downarrow
was for, the etc
stop words

\downarrow
only words, short words
'for' 'the' are removed
so average size increases.

- Hardware Stuff

Index Construction

Sort based Index Construction

(TermID, docID) \leftarrow 8 bytes
4B 4B

Term1 \rightarrow 1 school \rightarrow 1
Term2 \rightarrow 2 form \rightarrow 2
:

If $T = 100,000,000$ docs \rightarrow 800,000,000 bytes
/ 1024

Scaling Index Construction

BSBI

/ 1024

= 0.75 GB

$T = 100,000,000 \rightarrow$ too big to load in memory

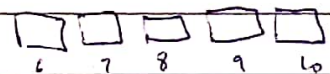
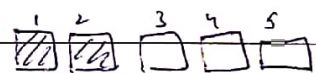
so divide in 10 chunks of 10,000,000 each

① 8 byte (TermID, docID) \rightarrow 100,000,000

② Now ^{have to} sort 100,000,000 by TermID

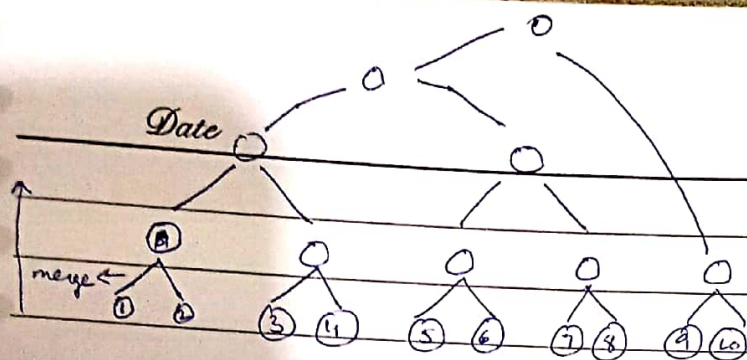
③ divide T and make chunks

④ sort each chunk



Code in slides

now 1, 2 are merged,
after this sort the merged



CON: if ① and ② are 10,000,000 so
200,000,000 needs to be loaded in memory

How to merge sorted runs?

take 5 million chunks each of 10 million, merge them.
(will take 2 weeks)

Multi-way merge

Another way: Take all 10 chunks of 100,000,000 at once
and process a part of all 10 chunks. Take 2 million
from each 10 chunks. Equal to 20 million.

Single Pass in Memory Indexing

- ① generate separate dics for each 10 blocks
no Term - term ID x only Term - DocID
- ② don't sort. accumulate postings in postings list
- ③ inverted index for each block
- ④ separate indexes can be merged into one big index

Code in Slide

Date March 3rd, 2022.

SPIMI : Compression

Store index in ~~DB~~ Main memory. If too large then store in secondary storage. If it's too large then distributed storage. Put it in multiple nodes (systems).

Distributed Indexing:

- one master machine: directs indexing job
important + safe
- master divides indexing tasks + assigns it to idle machines from pool

Parallel Tasks

• (1) Parsers

(2) Writers

after parsing we get:

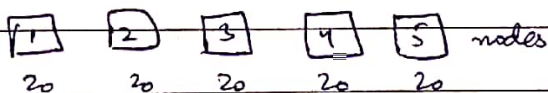
(Term, docID) pairs

sort (term, docID) → postings → inverted index

- split doc. collection, each split is a subset of doc (corresponding to blocks in BSB1/SPIMI).

100 docs

□ master

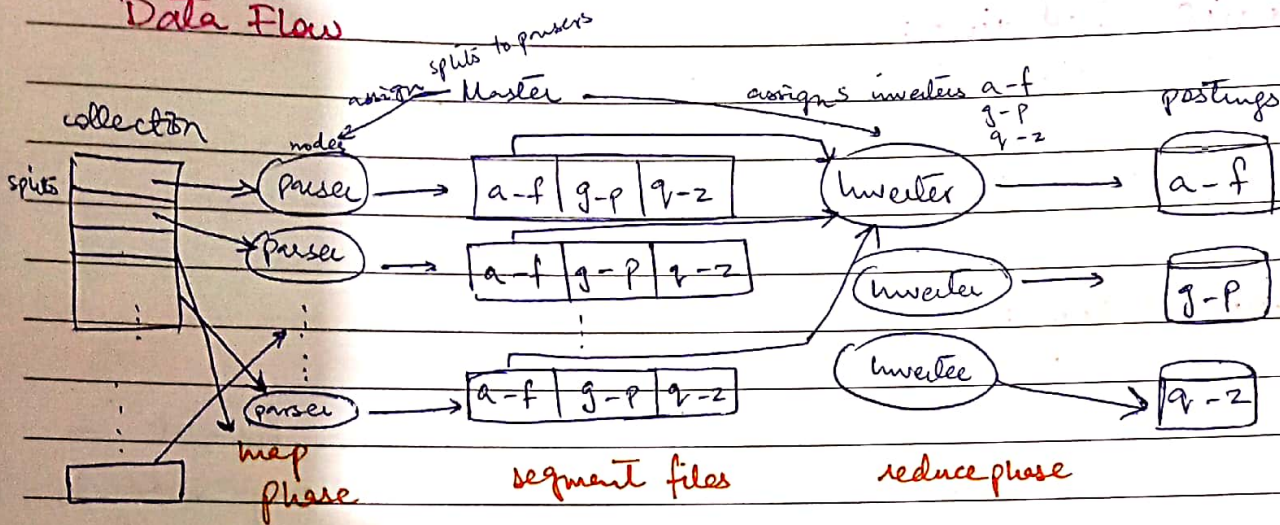


① first parse, each node does it for their own 20

it's possible >1 nodes have same terms so we have to combine all terms afterwards

Date

Data Flow



map phase $\text{Map}(f, [a_1, a_2, \dots, a_n])$ $f(a_1) \rightarrow f(a_2) \dots f(a_n)$

\downarrow parser \downarrow collection

reduce phase $\text{Merge}(a-f_{p_1}, a-f_{p_2}, \dots, a-f_{p_n})$

Parsers

$j=3$ partitions $a-f, g-p, q-z$

MapReduce

application: Hadoop

map: collection \rightarrow list (TermID, docID)

reduce:

example in slides

problem w/ postings based on alphabetic

Terms: issues with

balancing of loads.

if everybody is searching

for Pakistan Then $\boxed{g-p}$

will have too much load

Solution: do partition on

bases of doc. IDs.

Postings

pakistan, 3 1-20

21-40

pakistan, 7 41-60

collection

new load is distributed

Date

Dynamic Indexing:

before this collection was static

March 8th, 2022. (watch video)

Logarithmic Merge:

Z_0 : put terms in it, when reaches n length. Put it in disk with name I_0 . Again build Z_0 (auxiliary index), it merges with I_0 , I_0 gets erased, now I_1 is built and stored in disk.

Algo in slides.

Issues w/ multiple indexes:

incidence matrix

inverted index

- pos. indexing

- boolean queries

Date

March 15th, 2022.

Scoring, Term Weightage, Ranking

• Ranked Retrieval:

Jaccard coefficient:

$$\text{jaccard}(A, B) = |A \cap B| / |A \cup B|$$

ex: doc. size

$$(A, A) = 1$$

affects the score

$$(A, B) = 0 \text{ if } A \cap B = \emptyset$$

- $A + B$ may not have same size
- gives value b/w 0 + 1

- rare terms are more informative

• Query-doc matching scores

- assign score to query/doc pair

- if does not occur in doc: score 0

- if frequency \uparrow , score \uparrow

- in index: replace 1 with freq count of terms in a doc (vector representation)

Bag-of-words Model:

- in vector representation, order of terms occurring is not considered

Term frequency (tf): $tf_{t,d} \rightarrow$ term freq of term t in doc d

Raw term freq:

	doc1	doc2
school:	1	10

← This doesn't mean doc2 is 10x more relevant than doc1.
it only means doc2 is more relevant

- relevancy doesn't increase proportionally with tf.

tf → Term freq
idf → doc freq

Date

used it so we get weights
in smaller range so ease in calculations.

Log-frequency weighting:

$$W_{t,d} = \begin{cases} 1 + \log_{10} tf_{t,d} & \text{if } tf_{t,d} > 0 \\ 0 & \text{otherwise} \end{cases}$$

- now put weights in index

→ calculate score for only common words
b/w query + doc.

$$\text{Score} = \sum_{t \in q \cap d} (1 + \log_{10} tf_{t,d})$$

- now put scores in vector index style

- tf is giving us score

(idf) tf-idf
idf weight: → inverted style

- df_t → doc. frequency of t : no. of docs that contain t
→ total no. of docs

$$idf_t = \log_{10} \left(\frac{N}{df_t} \right)$$

→ how many docs contain the term

$$\left[\begin{aligned} \frac{100}{5} = 20 \text{ idf} = 2 &\rightarrow \text{giving more weight to more rare doc} \\ \frac{100}{30} = 3.3 \text{ idf} = 0.3 \end{aligned} \right]$$

- capricious person

↓
tf ↓

idf ↑

is rare

↓
tf ↑

idf ↓

is common

$$W_{t,d} = 1 + \log_{10} (tf_{t,d}) \times \log_{10} (N / df_t)$$

↑ if freq increases
of a term

tf

idf

↑ if rarity increases
of a term