Text Indexing

Manish Gupta

The Classic Search Model



Query on Unstructured Data

- Which plays of Shakespeare contain the words *Brutus* AND *Caesar* but NOT *Calpurnia*?
- One could grep all of Shakespeare's plays for *Brutus* and *Caesar*, then strip out lines containing *Calpurnia*?
- Why is that not the answer?
 - Slow (for large corpora)
 - <u>NOT</u> Calpurnia is non-trivial
 - Other operations (e.g., find the word *Romans* near *countrymen*) not feasible
 - Ranked retrieval (best documents to return)

Term-Document Incidence Matrices



Incidence Vectors

- So we have a 0/1 vector for each term.
- To answer query: take the vectors for *Brutus, Caesar* and *Calpurnia* (complemented) → bitwise *AND*.

- 110100 AND 1101111 AND 101111 = **100100**

	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

Bigger Collections

- Consider N = 1 million documents, each with about 1000 words.
- Avg 6 bytes/word including spaces/punctuation
 - 6GB of data in the documents.
- Say there are *M* = 500K *distinct* terms among these.

Can't Build the Matrix

- 500K x 1M matrix has half-a-trillion 0's and 1's.
- But it has no more than one billion 1's.
 - matrix is extremely sparse.
- What's a better representation?
 - We only record the 1 positions.



Inverted Index

- For each term *t*, we must store a list of all documents that contain *t*.
 - Identify each doc by a **docID**, a document serial number
- Can we used fixed-size arrays for this?



Inverted Index

- We need variable-size postings lists
 - On disk, a continuous run of postings is normal and best
 - In memory, can use linked lists or variable length arrays

Sorted by docID



Inverted Index Construction



Initial Stages of Text Processing

- Tokenization
 - Cut character sequence into word tokens
 - Deal with "John's", a state-of-the-art solution
- Normalization
 - Map text and query term to same form
 - You want **U.S.A.** and **USA** to match
- Stemming
 - We may wish different forms of a root to match
 - authorize, authorization
- Stop words
 - We may omit very common words (or not)
 - the, a, to, of

Indexer Steps: Token Sequence

did 1 Sequence of (Modified token, Document ID) pairs. enact 1 julius 1 1 caesar 1 was 1 killed 1 1 the 1 capitol 1 1 brutus 1 killed 1 me Doc 1 Doc 2 2 so 2 let 2 it 2 be 2 with 2 caesar I did enact Julius 2 the So let it be with 2 noble Caesar I was killed 2 brutus Caesar. The noble 2 hath i' the Capitol; 2 Brutus hath told you told 2 vou Brutus killed me. 2 Caesar was ambitious caesar 2 was 2 ambitious



Term

Indexer Steps: Sort

- Sort by terms
 - And then docID



Term	docID	T	erm	docID
I	1	a	mbitious	2
did	1	b	e	2
enact	1	b	rutus	1
julius	1	b	rutus	2
caesar	1	C	apitol	1
1	1	C	aesar	1
was	1	С	aesar	2
killed	1	С	aesar	2
ř	1	d	id	1
the	1	е	nact	1
capitol	1	h	ath	1
brutus	1	1		1
killed	1			1
me	1	i		1
SO	2	it		2
let	2	ju	ılius	1
it	2	k	illed	1
be	2	k	illed	1
with	2	le	et	2
caesar	2	n	ne	1
the	2	n	oble	2
noble	2	S	0	2
brutus	2	th	ne	1
hath	2	tł	ne	2
told	2	to	old	2
you	2	У	ou	2
caesar	2	W	as	1
was	2	w	as	2
ambitious	2	W	vith	2

Indexer Steps: Dictionary & Postings

- Multiple term entries in a single document are merged.
- Split into Dictionary and Postings
- Doc. frequency information is added.





Sec. 1.2

Where do we pay in Storage?



Query Processing: AND

• Consider processing the query:

Brutus AND Caesar

- Locate *Brutus* in the Dictionary;
 - Retrieve its postings.
- Locate *Caesar* in the Dictionary;
 - Retrieve its postings.
- "Merge" the two postings (intersect the document sets):



The Merge

• Walk through the two postings simultaneously, in time linear in the total number of postings entries

If the list lengths are x and y, the merge takes O(x+y) operations.

<u>Crucial</u>: postings sorted by docID.

Boolean Queries: Exact Match

- The Boolean retrieval model is being able to ask a query that is a Boolean expression:
 - Boolean Queries are queries using AND, OR and NOT to join query terms
 - Views each document as a <u>set</u> of words
 - Is precise: document matches condition or not.
 - Perhaps the simplest model to build an IR system on
- Primary commercial retrieval tool for 3 decades.
- Many search systems you still use are Boolean:
 - Email, library catalog, Mac OS X Spotlight

Query Optimization

- What is the best order for query processing?
- Consider a query that is an AND of *n* terms.
- For each of the *n* terms, get its postings, then AND them together.



Query: Brutus AND Calpurnia AND Caesar

Query Optimization Example

- <u>Process in order of increasing freq</u>:
 - start with smallest set, then keep cutting further.



Execute the query as (*Calpurnia* AND *Brutus)* AND *Caesar*.

More General Optimization

- e.g., (madding OR crowd) AND (ignoble OR strife)
- Get doc. freq.'s for all terms.
- Estimate the size of each *OR* by the sum of its doc. freq.'s (conservative).
- Process in increasing order of OR sizes.

Phrase Queries

- We want to be able to answer queries such as "stanford university" as a phrase
- Thus the sentence "I went to university at Stanford" is not a match.
 - The concept of phrase queries has proven easily understood by users; one of the few "advanced search" ideas that works
 - Many more queries are *implicit phrase queries*
- For this, it no longer suffices to store only

<term : docs> entries

A First Attempt: Bi-word Indexes

- Index every consecutive pair of terms in the text as a phrase
- For example the text "Friends, Romans, Countrymen" would generate the biwords
 - friends romans
 - romans countrymen
- Each of these bi-words is now a dictionary term
- Two-word phrase query-processing is now immediate.

Longer Phrase Queries

- Longer phrases can be processed by breaking them down
- stanford university palo alto can be broken into the Boolean query on biwords:

stanford university AND university palo AND palo alto

Without the docs, we cannot verify that the docs matching the above Boolean query do contain the phrase.



Issues for Bi-word Indexes

- False positives, as noted before
- Index blowup due to bigger dictionary
 - Infeasible for more than bi-words, big even for them
- Bi-word indexes are not the standard solution (for all bi-words) but can be part of a compound strategy

Solution 2: Positional Indexes

 In the postings, store, for each *term* the position(s) in which tokens of it appear

<term, number of docs containing term; doc1: position1, position2 ... ; doc2: position1, position2 ... ; etc.>

Positional Index Example

- For phrase queries, we use a merge algorithm recursively at the document level
- But we now need to deal with more than just equality <*be*: 993427;
- *1*: 7, 18, 33, 72, 86, 231; 2: 3, 149;
 - *4*: 17, 191, 291, 430, 434;
 - **5**: 363, 367, ...>



Processing a Phrase Query

- Extract inverted index entries for each distinct term: *to, be, or, not.*
- Merge their *doc:position* lists to enumerate all positions with "*to be or not to be*".
 - **to**:
 - 2:1,17,74,222,551; 4:8,16,190,429,433; 7:13,23,191; ...
 - be:
 - 1:17,19; 4:17,191,291,430,434; 5:14,19,101; ...
- Same general method for proximity searches

Proximity Queries

- LIMIT! /3 STATUTE /3 FEDERAL /2 TORT
 - Again, here, /k means "within k words of".
- Clearly, positional indexes can be used for such queries; bi-word indexes cannot.

Further Reading

- Chapters 1,2,5 of <u>Manning-Raghavan-Schuetze book</u>
 - http://nlp.stanford.edu/IR-book/
- Chapter 3 (Web Search and Information Retrieval) from Mining the Web
 - <u>http://www.cse.iitb.ac.in/soumen/mining-the-web/</u>
- Original publication on SPIMI: Heinz and Zobel (2003)
- F. Scholer, H.E. Williams and J. Zobel. 2002. Compression of Inverted Indexes For Fast Query Evaluation. *Proc. ACM-SIGIR 2002*.
 - Variable byte codes
- V. N. Anh and A. Moffat. 2005. Inverted Index Compression Using Word-Aligned Binary Codes. *Information Retrieval* 8: 151–166.
 - Word aligned codes
- As We May Think -- Vannevar Bush
 - http://www.theatlantic.com/magazine/archive/1945/07/as-we-may-think/303881/