

Outlines

- We will discuss two more topics.
 - Boolean retrieval
 - Posting list

Web Information Retrieval

Textbook by
 Christopher D. Manning, Prabhakar Raghavan,
 and Hinrich Schütze
Notes Revised by X. Meng for SEU
 May 2014

2

Tokenizing And Preprocessing

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:



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

Posting Lists

3

4

4

Generate Posting

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

5

5

Sort Postings

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

6

6


Creating Postings Lists, Determine Document Frequency

term	docID			
ambitious	2			
be	2			
brutus	1			
brutus	2			
capitol	1			
canisur	1			
canisur	2			
canisur	2			
did	1			
emact	1			
emact	1			
hath	1			
i	1			
i	1			
i	1			
i	1			
it	2			
julius	1			
killed	1			
killed	1			
let	2			
let	2			
me	1			
mobile	2			
so	2			
the	1			
the	2			
told	2			
you	2			
was	1			
was	2			
with	2			


7

Split the Result into Dictionary and Postings File

BRUTUS	→	1	2	4	11	31	45	173	174	
CAESAR	→	1	2	4	5	6	16	57	132	...
CALPURNIA	→	2	31	54	101					
⋮										



dictionary

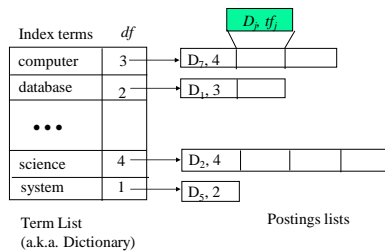


postings

8

8

Inverted Index



In-Class Work

- Creating an inverted index for the following documents
 - “The quick brown fox jumps over the lazy dog”
 - “The quick fox run”
 - “The lazy dog sleep”
- Assume:
 - “the” and “over” are stopwords that are not indexed
 - various forms of verbs are not indexed separately, only the original form (stem) is indexed

11

Creating an Inverted Index

- Create an empty index term list I;
- For each document, D, in the document set V
 - For each (non-zero) token, T, in D:
 - If T is not already in I
 - Insert T into I;
 - Find the location for T in I;
 - If (T, D) is in the posting list for T
 - increase its term frequency for T;
 - Else
 - Create (T, D);
 - Add it to the posting list for T;

Boolean Retrieval

Processing Boolean queries

Query optimization

12

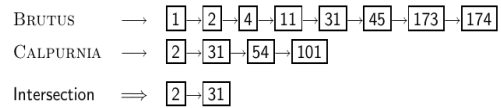
Simple Conjunctive Query (two terms)

- Consider the query: BRUTUS AND CALPURNIA
- To find all matching documents using inverted index:
 - Locate BRUTUS in the dictionary (term list)
 - Retrieve its postings list from the postings file
 - Locate CALPURNIA in the dictionary
 - Retrieve its postings list from the postings file
 - Intersect the two postings lists
 - Return intersection to user

13

13

Intersecting Two Posting Lists



- The complexity is linear in the length of the postings lists.
- Note: This only works if postings lists are sorted.

14

14

Algorithm of Intersecting Two Posting Lists

Assume posting lists are sorted by docID

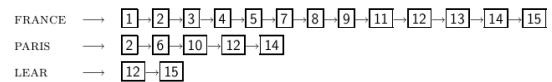
```

INTERSECT( $p_1, p_2$ )
1   $answer \leftarrow \langle \rangle$ 
2  while  $p_1 \neq NIL$  and  $p_2 \neq NIL$ 
3  do if  $docID(p_1) = docID(p_2)$ 
4      then  $ADD(answer, docID(p_1))$ 
5           $p_1 \leftarrow next(p_1)$ 
6           $p_2 \leftarrow next(p_2)$ 
7  else if  $docID(p_1) < docID(p_2)$ 
8      then  $p_1 \leftarrow next(p_1)$ 
9      else  $p_2 \leftarrow next(p_2)$ 
10 return  $answer$ 
  
```

15

15

Query Processing: In-Class Work



Compute hit list for ((paris AND NOT france) OR lear)

16

16

Boolean Queries

- The Boolean retrieval model can answer any query that is a Boolean expression.
 - Boolean queries are queries that use AND, OR and NOT to join query terms.
 - Views each document as a set of terms.
 - Precise: either document matches condition or not, nothing in between.
- Primary commercial retrieval tool for three decades
- Many professional searchers (e.g., lawyers) still like Boolean queries.
 - You know exactly what you are getting.
- Many search systems you use are also Boolean: spotlight, email, intranet etc.

17

17

Outline

- Processing Boolean queries
- Query optimization

18

Query Optimization

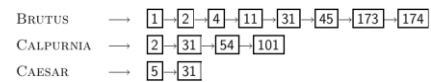
- Consider a query that is an *and* of n terms, $n > 2$
- For each of the terms, get its postings list, then *and* them together
- Example query: BRUTUS AND CALPURNIA AND CAESAR
- What is the best order for processing this query? That is, should we process BRUTUS first? CALPURNIA first? Or CAESAR first?

19

19

Query Optimization

- Example query: BRUTUS AND CALPURNIA AND CAESAR
- Simple and effective optimization: Process in order of increasing frequency
- Start with the shortest postings list, then keep cutting further
- In this example, first CAESAR, then CALPURNIA, then BRUTUS



20

20

Optimized Intersection Algorithm for Conjunctive Queries

```

INTERSECT( $\langle t_1, \dots, t_n \rangle$ )
1   $terms \leftarrow \text{SORTBYINCREASINGFREQUENCY}(\langle t_1, \dots, t_n \rangle)$ 
2   $result \leftarrow \text{postings}(\text{first}(terms))$ 
3   $terms \leftarrow \text{rest}(terms)$ 
4  while  $terms \neq \text{NIL}$  and  $result \neq \text{NIL}$ 
5  do  $result \leftarrow \text{INTERSECT}(result, \text{postings}(\text{first}(terms)))$ 
6      $terms \leftarrow \text{rest}(terms)$ 
7  return  $result$ 

```

21

21

More General Optimization

- Example query: (MADDING OR CROWD) and (IGNOBLE OR STRIFE)
- Get frequencies for all terms
- Estimate the size of each *or* by the sum of its frequencies (conservative)
- Process in increasing order of *or* sizes

22

22