CS105 Introduction to Information Retrieval

Lecture: Yang Mu UMass Boston

> Slides are modified from: http://www.stanford.edu/class/cs276/

Information Retrieval

- 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).
 - These days we frequently think first of web search, but there are many other cases:
 - E-mail search
 - Searching your laptop
 - Corporate knowledge bases
 - Legal information retrieval

Unstructured (text) vs. structured (database) data in the mid-nineties



Unstructured (text) vs. structured (database) data today



Basic assumptions of Information Retrieval

- Collection: A set of documents
 - Assume it is a static collection for the moment
- Goal: Retrieve documents with information that is relevant to the user's information need and helps the user complete a task

The classic search model



How good are the retrieved docs?

- *Precision* : Fraction of retrieved docs that are relevant to the user's information need
- *Recall* : Fraction of relevant docs in collection that are retrieved

Term-document incidence

Unstructured data in 1620

- Which plays of Shakespeare contain the words *Brutus AND Caesar* but *NOT Calpurnia*?
- One could get 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



Bigger collections

- 500K (distinct words) x 1M (documents) 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

The key data structure underlying modern IR

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



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

• Sequence of (Modified token, Document ID) pairs.

Doc 1

Doc 2

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

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



Indexer steps: Sort

- Sort by terms
 - And then docID



_		
lerm	docID	
1	1	
did	1	
enact	1	
julius	1	
caesar	1	
1	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	

Torm	doolD
ombitiouro	0001D
ampluous	2
be	Z
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

Indexer steps: Dictionary & Postings

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



Will discuss later.

		term doc. freq.	\rightarrow	postings list
Term	docID	ambitious 1	\rightarrow	2
ambitious	2	be 1	\rightarrow	2
be	2			
brutus	1	brutus 2	\rightarrow	$ 1 \rightarrow 2 $
brutus	2	capitol 1	\rightarrow	1
capitol	1		1	
caesar	1	caesar 2	\rightarrow	$1 \rightarrow 2$
caesar	2	did 1	\rightarrow	1
caesar	2			-
did	1	enact I	\rightarrow	1
enact	1	hath 1	\rightarrow	2
hath	1			1
<u> </u>	1		\rightarrow	1
<u> </u>	1	i' 1	\rightarrow	1
i'	1	+ 1		
it	2		\rightarrow	2
julius	1	julius 1	\rightarrow	1
killed	1			1
killed	1	killed 1	-	1
let	2	let 1	\rightarrow	2
me	1	me 1	\rightarrow	1
noble	2			
SO	2	noble 1	\rightarrow	2
the	1	so 1	\rightarrow	2
the	2			
told	2	the 2	\rightarrow	$1 \rightarrow 2$
you	2	told 1	\rightarrow	2
was	1			
was	2	you 1	\rightarrow	2
with 2	was 2	\rightarrow	$1 \rightarrow 2$	
	with 1	\rightarrow	2	

Query processing with an inverted index

The index we just built

How do we process a query? Our focus
 Later - what kinds of queries can we process?

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):



21

IR System Components

- Text processing forms index words (tokens).
- Indexing constructs an *inverted index* of word to document pointers.
- Searching retrieves documents that contain a given query token from the inverted index.
- Ranking scores all retrieved documents according to a relevance metric.