#### Multimedia Information Extraction and Retrieval

Introduction

#### Ralf Moeller Hamburg Univ. of Technology

#### Literature

- <u>Christopher D. Manning, Prabhakar Raghavan and</u> <u>Hinrich Schütze</u>, *Introduction to Information* <u>Retrieval</u>, *Cambridge University Press. 2008*.
- http://nlp.stanford.edu/IR-book/informationretrieval-book.html



#### Unstructured (text) vs. structured (database) data in 1996



#### Unstructured (text) vs. structured (database) data in 2006



## **Unstructured data in 1650**

- 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*?
  - 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)
    - Later lectures

#### **Term-document incidence**

	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
Brutus . Calpurr	AND <b>Caesar</b> bu nia	t NOT		1 if pla word, (	y cont ) othe	ains rwise

#### **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 110111 AND 101111 = 100100.

#### **Answers to query**

#### Antony and Cleopatra, Act III, Scene ii

Agrippa [Aside to DOMITIUS ENOBARBUS]: Why, Enobarbus,

When Antony found Julius *Caesar* dead, He cried almost to roaring; and he wept When at Philippi he found *Brutus* slain.

#### • Hamlet, Act III, Scene ii

Lord Polonius: I did enact Julius **Caesar** I was killed i' the Capitol; **Brutus** killed me.

## **Bigger corpora**

- Consider N = 1M documents, each with about 1K terms.
- Avg 6 bytes/term incl spaces/punctuation
  - 6GB of data in the documents.
- Say there are m = 500K <u>distinct</u> 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.



• We only record the 1 positions.





- For each term *T*, we must store a list of all documents that contain *T*.
- Do we use an array or a list for this?



What happens if the word **Caesar** is added to document 14?

#### **Inverted** index

- Linked lists generally preferred to arrays
  - Dynamic space allocation
  - Insertion of terms into documents easy
  - Space overhead of pointers



Posting



#### **Indexer** steps



#### **Indexer steps**



Term	Doc#	Term	Doc#
1	1	amhitious	2
did	1	he	2
enact	1	hrutus	1
iulius	1	hrutus	2
caesar	1	canitol	1
l	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		1
killed	1		1
me	1	i'	1
S0	2	it	2
let	2	iulius	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

#### **Indexer steps**

- Multiple term entries in a single document are merged.
- Frequency information is added.

Why frequency? Will discuss later.

Term	Doc #	Т	erm
ambitious	2	a	ambiti
be	2	b	e
brutus	1	b	rutus
brutus	2	b	rutus
capitol	1	C	apitol
caesar	1	C	aesai
caesar	2	C	aesai
caesar	2	C	lid
did	1	е	enact
enact	1	h	nath
hath	1		
	1	i'	J
	1	it	t
i'	1	i	ulius
it	2	k	illed
julius	1	1	et
killed	1	n	ne
killed	1	n	noble
let	2		0
me	1	+	he he
noble	2	+	ho
SO	2	t.	
the	1		
the	2	y M	vas
told	2	V	vas
you	2	V	vas
was	1	V	VILII
was	2		
with	2		
		L	

Term	Doc #	Term freq
ambitious	2	1
be	2	1
brutus	1	1
brutus	2	1
capitol	1	1
caesar	1	1
caesar	2	2
did	1	1
enact	1	1
hath	2	1
Ι	1	2
i'	1	1
it	2	1
julius	1	1
killed	1	2
let	2	1
me	1	1
noble	2	1
<b>S</b> 0	2	1
the	1	1
the	2	1
told	2	1
you	2	1
was	1	1
was	2	1
with	2	1

# • The result is split into a *Dictionary* file and a *Postings* file.

Term	Doc #	Freq
ambitious	2	1
be	2	1
brutus	1	1
brutus	2	1
capitol	1	1
caesar	1	1
caesar	2	2
did	1	1
enact	1	1
hath	2	1
I	1	2
i'	1	1
it	2	1
julius	1	1
killed	1	2
let	2	1
me	1	1
noble	2	1
<b>S</b> 0	2	1
the	1	1
the	2	1
told	2	1
you	2	1
was	1	1
was	2	1
with	2	1

Term	N docs	Coll free
ambitious	1	1
be	1	1
brutus	2	2
capitol	1	1
caesar	2	3
did	1	1
enact	1	1
hath	1	1
1	1	2
i'	1	1
it	1	1
julius	1	1
killed	1	2
let	1	1
me	1	1
noble	1	1
SO	1	1
the	2	2
told	1	1
you	1	1
was	2	2
with	1	1



#### • Where do we pay in storage?



# The index we just built

- How do we process a query?
  - 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:





• 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.
- Primary commercial retrieval tool for 3 decades.
- Professional searchers (e.g., lawyers) still like Boolean queries:
  - You know exactly what you're getting.

#### **Example: WestLaw**

http://www.westlaw.com/

- Largest commercial (paying subscribers) legal search service (started 1975; ranking added 1992)
- Tens of terabytes of data; 700,000 users
- Majority of users *still* use boolean queries
- Example query:
  - What is the statute of limitations in cases involving the federal tort claims act?
  - LIMIT! /3 STATUTE ACTION /S FEDERAL /2 TORT /3 CLAIM
- /3 = within 3 words, /S = in same sentence

#### **Example: WestLaw**

http://www.westlaw.com/

- Another example query:
  - Requirements for disabled people to be able to access a workplace
  - disabl! /p access! /s work-site work-place (employment /3 place
- Note that SPACE is disjunction, not conjunction!
- Long, precise queries; proximity operators; incrementally developed; not like web search
- Professional searchers often like Boolean search:
  - Precision, transparency and control
- But that doesn't mean they actually work better....

#### Boolean queries: More general merges

 Exercise: Adapt the merge for the queries: Brutus AND NOT Caesar Brutus OR NOT Caesar

Can we still run through the merge in time O (x+y) or what can we achieve?

# Merging

What about an arbitrary Boolean formula? (Brutus OR Caesar) AND NOT (Antony OR Cleopatra)

- Can we always merge in "linear" time?
  - Linear in what?
- Can we do better?

# **Query optimization**

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



Query: Brutus AND Calpurnia AND Caesar



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

## **More general optimization**

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

#### Exercise

 Recommend a query processing order for

(tangerine OR trees) AND (marmalade OR skies) AND (kaleidoscope OR eyes)

Term	Freq
eyes	213312
kaleidoscope	<b>87009</b>
marmalade	107913
skies	<b>271658</b>
tangerine	<b>46653</b>
trees	316812

## **Query processing exercises**

- If the query is *friends* AND romans AND (NOT countrymen), how could we use the freq of countrymen?
- Exercise: Extend the merge to an arbitrary Boolean query. Can we always guarantee execution in time linear in the total postings size?
- Hint: Begin with the case of a Boolean *formula* query: in this, each query term appears only once in the query.

#### What's ahead in IR? Beyond term search

- What about phrases?
  - Stanford University
- Proximity: Find *Gates NEAR Microsoft*.
  - Need index to capture position information in docs. More later.
- Zones in documents: Find documents with (*author = Ullman*) AND (text contains *automata*).

#### **Evidence** accumulation

#### • 1 vs. 0 occurrence of a search term

- 2 vs. 1 occurrence
- 3 vs. 2 occurrences, etc.
- Usually more seems better
- Need term frequency information in docs

## **Ranking search results**

- Boolean queries give inclusion or exclusion of docs.
- Often we want to rank/group results
  - Need to measure proximity from query to each doc.
  - Need to decide whether docs presented to user are singletons, or a group of docs covering various aspects of the query.

#### IR vs. databases: Structured vs unstructured data

 Structured data tends to refer to information in "tables"

Employee	Manager	Salary
Smith	Jones	50000
Chang	Smith	60000
lvy	Smith	50000

*Typically allows numerical range and exact match* (for text) queries, e.g., *Salary < 60000 AND Manager = Smith.* 

#### **Unstructured data**

- Typically refers to free text
- Allows
  - Keyword queries including operators
  - More sophisticated "concept" queries e.g.,
    - find all web pages dealing with *drug abuse*
- Classic model for searching text documents

#### Semi-structured data

- In fact almost no data is "unstructured"
- E.g., this slide has distinctly identified zones such as the *Title* and *Bullets*
- Facilitates "semi-structured" search such as
  - *Title* contains <u>data</u> AND *Bullets* contain <u>search</u>

... to say nothing of linguistic structure

# More sophisticated semi-structured search

- Title is about <u>Object Oriented Programming</u> AND Author something like <u>stro\*rup</u>
- where \* is the wild-card operator
- Issues:
  - how do you process "about"?
  - how do you rank results?
- The focus of XML search.

# **Clustering and classification**

- Given a set of docs, group them into clusters based on their contents.
- Given a set of topics, plus a new doc *D*, decide which topic(s) *D* belongs to.

# The web and its challenges

- Unusual and diverse documents
- Unusual and diverse users, queries, information needs
- Beyond terms, exploit ideas from social networks
  - link analysis, clickstreams ...
- How do search engines work? And how can we make them better?

# More sophisticated *information* retrieval

- Cross-language information retrieval
- Question answering
- Summarization
- Text mining