Deconstructing Wikipedia vs. Project Gutenberg





Introduction to Information Retrieval CS 150 Donald J. Patterson

Content adapted from Hinrich Schütze http://www.informationretrieval.org

OVERVIEW

- Boolean Retrieval
- Weighted Boolean Retrieval
- Zone Indices
- Term Frequency Metrics
- The full vector space model



FROM THE BOTTOM

- "Grep"
 - Querying without an index or a crawl
 - Whenever you want to find something you look through the entire document for it.
 - Example:
 - You have the collected works of Shakespeare on disk
 - You want to know which play contains the words
 - "Brutus AND Caesar"

FROM THE BOTTOM

- "Grep"
 - "Brutus AND Caesar" is the query.
 - This is a boolean query. Why?
 - What other operators could be used?
 - The grep solution:
 - Read all the files and all the text and output the intersection of the files



FROM THE BOTTOM

- "Grep"
 - Slow for large corpora
 - Calculating "NOT" requires exhaustive scanning
 - Some operations not feasible
 - Query: "Romans NEAR Countrymen"
 - Doesn't support ranked retrieval
- Moving beyond grep is the motivation for the inverted index.



OVERVIEW

- Boolean Retrieval
- Weighted Boolean Retrieval
- Zone Indices
- Term Frequency Metrics
- The full vector space model

BOOLEAN RETRIEVAL

Our inverted index is a 2-D array or Matrix

A Column For Each Document

	Anthony	Julius	The	Hamlet	Othello	Macbeth
	and	Caesar	Tempest			
	Cleopatra					
Anthony	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

...

- Boolean Query
 - Queries are boolean expressions
 - Search returns all documents which satisfy the expression
 - Does Google use the Boolean model?



BOOLEAN RETRIEVAL

Boolean Query

. . .

- Straightforward application of inverted index
- where cells of inverted index are (0,1)
 - indicating presence or absence of a term

Document

		Anthony and Cleonatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
,	Anthony		1	0	0	0	1
Term	Anthony	1	1	0	U	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

BOOLEAN RETRIEVAL

• Boolean Query

. . .

- 0/1 vector for each term
- "Brutus AND Caesar AND NOT Calpurnia =
- Perform bitwise Boolean operation on each row:
 - 110100 AND 110111 AND !(010000) = 100100 Document

		Anthony and	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
		Cleopatra		-			
Term	Anthony	ī	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

- Boolean Query
 - A big corpus means a sparse matrix
 - A sparse matrix motivates the introduction of the posting
 - Much less space to store
 - Only recording the "1" positions



- Boolean Query
 - Query processing on postings
 - Brutus AND Caesar
 - Locate the postings for Brutus
 - Locate the postings for Caesar
 - Merge the postings



- Boolean Query
 - Merging -> walk through the two posting simultaneously
 - postings sorted by doc ID



- Boolean Query
 - Merging -> walk through the two posting simultaneously
 - postings sorted by doc ID



- Boolean Query
 - Merging -> walk through the two posting simultaneously
 - postings sorted by doc ID



- Boolean Query
 - Merging -> walk through the two posting simultaneously
 - postings sorted by doc ID



- Boolean Query
 - Merging -> walk through the two posting simultaneously
 - postings sorted by doc ID



- Boolean Query
 - Merging -> walk through the two posting simultaneously
 - postings sorted by doc ID



- Boolean Query
 - Merging -> walk through the two posting simultaneously
 - postings sorted by doc ID



- Boolean Query
 - Merging -> walk through the two posting simultaneously
 - postings sorted by doc ID



- Boolean Query
 - Merging -> walk through the two posting simultaneously
 - postings sorted by doc ID



- Boolean Query
 - Merging -> walk through the two posting simultaneously
 - postings sorted by doc ID





BOOLEAN RETRIEVAL

4

5

6

7

8

9

Boolean Query

- An algorithm based on postings
- Linear in the size of the postings INTERSECT (p_1, p_2)
 - 1 $answer \leftarrow <>$
 - 2 while $p_1 \neq nil$ and $p_2 \neq nil$
 - 3 do if $docID(p_1) = docID(p_2)$

then $ADD(answer, docID(p_1))$

- $p_1 \leftarrow next(p_1)$
- $p_2 \leftarrow next(p_2)$
 - else if $docID(p_1) < docID(p_2)$
 - then $p_1 \leftarrow next(p_1)$ else $p_2 \leftarrow next(p_2)$
- 10 return answer

BOOLEAN RETRIEVAL

- Boolean Query
 - Is the algorithmic complexity better than scanning?
 - Where would you put more complex formulae?

INTERSECT (p_1, p_2) 1 $answer \leftarrow <>$ 2while $p_1 \neq nil$ and $p_2 \neq nil$ do if $docID(p_1) = docID(p_2)$ 3 then $ADD(answer, docID(p_1))$ 4 $p_1 \leftarrow next(p_1)$ 5 $p_2 \leftarrow next(p_2)$ 6 7 else if $docID(p_1) < docID(p_2)$ 8 then $p_1 \leftarrow next(p_1)$ else $p_2 \leftarrow next(p_2)$ 9

10 return answer

- Boolean Queries
 - Exact match
 - Views each document as a "bag of words"
 - Precise: a document matches or it doesn't
 - Primary commercial retrieval tool for 3 decades
 - Professional searchers (e.g., lawyers) still like
 Boolean queries
 - Why?
 - No question about what you are getting

BUILDING UP OUR QUERY TECHNOLOGY

- Linear on-demand retrieval (aka grep)
- 0/1 Vector-Based Boolean Queries
- Posting-Based Boolean Queries

BUILDING UP OUR QUERY TECHNOLOGY

- Linear on-demand retrieval (aka grep)
- 0/1 Vector-Based Boolean Queries
- Posting-Based Boolean Queries

- Deconstruct what is happening here:
 - http://www.rhymezone.com/shakespeare/



QUERYING BOOLEAN MODEL VS. RANKED RETRIEVAL METHODS

- Only game for 30 years
- uses precise queries
- user decides relevance
- stayed current with proximity queries
- precise controlled queries
- transparent queries
- controlled queries

- Appeared with www
- uses "free-text" queries
- system decides relevance
- works with enormous

corpora

• "no guarantees" in queries

QUERYING - BOOLEAN SEARCH EXAMPLE

WESTLAW CASE STUDY

- Largest commercial (paying subscribers) legal search service (started in 1975, ranking added in 1992)
- Tens of terabytes of data, 700,000 users
- Majority of users still use boolean queries (default in 2005)
 - Example:
 - What is the statute of limitations in cases involving federal tort claims act?
 - LIMIT! /3 STATUTE ACT /S FEDERAL /2 TORT /3 CLAIM
 - /3 = within 3 words. /S same sentence

QUERYING - BOOLEAN SEARCH EXAMPLE

WESTLAW CASE STUDY

- Example:
 - Requirements for disabled people to be able to access a workplace
 - disabl! /p access! /s work-site work-place employment /3 place
 - space is a disjunction not a conjunction
 - long precise queries, proximity operators,
 incrementally developed, not like web search
 - preferred by professionals, but not necessarily better

BUILDING UP OUR QUERY TECHNOLOGY

- "Matching" search
 - Linear on-demand retrieval (aka grep)
 - 0/1 Vector-Based Boolean Queries
 - Posting-Based Boolean Queries
- Ranked search
 - Parametric Search

RANKED SEARCH

- Rather than saying
 - (query, document) matches or not (0,1)
 - ("Capulet", "Romeo and Juliet") = 1
- Now we are going to assign rankings
 - (query, document) in {0,1}
 - ("capulet", "Romeo and Juliet") = 0.7

METADATA

- "structured additional information about a document."
 - The author of a document
 - The creation date of a document
 - The title of a document
 - The location where a document was created
- author, creation date, title, location are fields
- searching for "William Shakespeare" in a doc differs from searching for "William Shakespeare" in the author of a doc

- supports searching on meta-data explicitly
- a parametric search interface allows a mix of full-text query and meta-data queries
- Example:
 - <u>http://carsquare.com/</u>

- Result is a large table
- Columns are fields
- Searching for "Lamborghini" only applied to make/model field



- Example:
 - <u>http://www.ocregister.com/realestate/</u>





- Example:
 - http://www.ocregister.com/realestate/
 - 92614: 218 results



PARAMETRIC SI

- Example:
 - http://www.oci
 - 92614: 218 res





View Details

\$999,800 5 Bedrooms 3 Baths 2,801 Sqft Single Family Residence

\$929,000

3 Baths

2,601 Sqft

Residence

\$839,000

2,341 Sqft

Residence

Single Family

3 Baths

Single Family

3 Salerno

Irvine, CA 92614

largest sorrento model in a private cul de sac location in One of the most desirable westpark neighborhood across the park/school grounds. brand new interior...

Save View #1

21 Decente 4 Bedrooms

Irvine, CA 92614

beautiful curb appeal! quiet interior location, cathedral ceilings. convenient main floor bed w/full bath. custom paint. separate laundry room. new roll ...

Save View #2

24 Toscany 4 Bedrooms

Irvine, CA 92614

largest model in the impeters promenade plan 234 home with a recent major kitchen & living area designer upgrades.custom maple/cherry wood kitchen cabinets, lapis...

Save View #3

- In these examples we select field values
 - Values could be hierarchical
 - USA -> California -> SB County -> Montecito
- It is a paradigm for navigating through a corpus
 - e.g, "Aerospace companies in Brazil" can be found by combining "Geography" and "Industry"
- Approach:
 - Filter for relevant documents
 - Run text searches on subset

- Index support for parametric search
 - Must be able to support queries of the form:
 - Find pdf documents that contain "Westmont"
 - Field selection and text query
- Field selection approach
 - Use inverted index of field values
 - (field value, docID)
 - organized by field name
 - Using same compression and sorting techniques

BUILDING UP OUR QUERY TECHNOLOGY

- "Matching" search
 - Linear on-demand retrieval (aka grep)
 - 0/1 Vector-Based Boolean Queries
 - Posting-Based Boolean Queries
- Ranked search
 - Parametric Search
 - Zones



ZONES

- A zone is an extension of a field
- A zone is an identified region of a document
 - e.g., title, abstract, bibliography
 - Generally identified by mark-up in a document
 - <title>Romeo and Juliet</title>
- Contents of zone are free text
 - Not a finite vocabulary
- Indices required for each zone to enable queries like:
 - (instant in TITLE) AND (oatmeal in BODY)
- Doesn't cover "all papers whose authors cite themselves"
 - Why?

PARAMETRIC/ZONE SEARCH

- Now, we crawl the corpus
- We parse the document keeping track of terms, fields and docIDs
- Instead of building just a (term, docID) pair
- We build (term, field, docID) triples
- These can then be combined into postings like this:



PARAMETRIC/ZONE SEARCH

- So are we just creating a database?
 - Not really.
 - Databases have more functionality
 - Transactions
 - Recovery
 - Our index can be recreated. Not so with database.
 - Text is never stored outside of indices
- We are focusing on optimized indices for text-oriented queries not a full SQL engine

BUILDING UP OUR QUERY TECHNOLOGY

- "Matching" search
 - Linear on-demand retrieval (aka grep)
 - 0/1 Vector-Based Boolean Queries
 - Posting-Based Boolean Queries
- Ranked search
 - Parametric Search
 - Zones
 - Scoring



