

COEN 169

Retrieval Models

Yi Fang

Department of Computer Engineering

Santa Clara University

The Search Task

- Given a query and a corpus, find relevant items

query: a textual description of the user's information need

corpus: a repository of documents

relevance: satisfaction of the user's information need

What is a Retrieval Model?

A formal method that predicts the degree of relevance of a document to a query

Relevance

- Many factors affect whether a document satisfies a particular user's information need
- Topicality, freshness, diversity, novelty, reading level, authority, etc.
- **Topical relevance**: the document is on the same topic as the query
- For now, we will only try to predict topical relevance

Overview of Retrieval Models

- Boolean
 - unranked Boolean retrieval
 - ranked Boolean retrieval
- Vector space models
 - SMART, Lucene
- Probabilistic models
 - Statistical language models
 - Lemur
 - Two Poisson model
 - Okapi
- Citation/Link analysis models
 - PageRank
 - Google
 - Hub & authorities

Boolean Retrieval

Assumption: the user can represent their information need using boolean constraints: **AND**, **OR**, and **AND NOT**

- lincoln
- president **AND** lincoln
- president **AND** (lincoln **OR** abraham)
- president **AND** (lincoln **OR** abraham) **AND NOT** car
- president **AND** (lincoln **OR** abraham) **AND NOT** (car **OR** automobile)
- Parentheses can specify the order of operations
 - A **OR** (B **AND** C) does not equal (A **OR** B) **AND** C

Boolean Retrieval

- If the query is too specific, precision may be high, but recall will probably be low
- If the query is too broad, recall may be high, but precision will probably be low
- Extreme cases:
 - a query that retrieves a single relevant document will have perfect precision, but low recall (unless only that one document is relevant)
 - a query that retrieves the entire collection will have perfect recall, but low precision (unless the whole collection is relevant)

Boolean Retrieval

- Easy for the system (no ambiguity in the query)
 - ✓ the burden is on the user to formulate the right query
- The user gets transparency and control
 - ✓ lots of results → the query is too broad
 - ✓ no results → the query is too narrow
- Common strategy for finding the right balance is:
 - ✓ if the query is too broad, add AND or AND NOT constraints
 - ✓ if the query is too narrow, add OR constraints

The Binary Full-text Representation

Document-term matrix

	<i>a</i>	<i>aardvark</i>	<i>abacus</i>	<i>abba</i>	<i>able</i>	...	<i>zoom</i>
<i>doc_1</i>	1	0	0	0	0	...	1
<i>doc_2</i>	0	0	0	0	1	...	1
::	::	::	::	::	::	...	0
<i>doc_m</i>	0	0	1	1	0	...	0

- 1 = the word appears in the document at least once
- 0 = the word does not appear in the document
- Does not represent word frequency, order, or position information

Processing a Boolean Query

Query: *Jack AND Jill*

- doc_1 Jack and Jill went up the hill
- doc_2 To fetch a pail of water.
- doc_3 Jack fell down and broke his crown,
- doc_4 And Jill came tumbling after.
- doc_5 Up Jack got, and home did trot,
- doc_6 As fast as he could caper,
- doc_7 To old Dame Dob, who patched his nob
- doc_8 With vinegar and brown paper.

Processing a Boolean Query

Query: *Jack* **AND** *Jill*

<i>docid</i>	<i>text</i>	<i>Jack</i>	<i>Jill</i>
<i>doc_1</i>	Jack and Jill went up the hill	1	1
<i>doc_2</i>	To fetch a pail of water.	0	0
<i>doc_3</i>	Jack fell down and broke his crown,	1	0
<i>doc_4</i>	And Jill came tumbling after.	0	1
<i>doc_5</i>	Up Jack got, and home did trot,	1	0
<i>doc_6</i>	As fast as he could caper,	0	0
<i>doc_7</i>	To old Dame Dob, who patched his nob	0	0
<i>doc_8</i>	With vinegar and brown paper.	0	0

Processing a Boolean Query

Query: *Jack* **AND** *Jill*

	<i>Jack</i>	<i>Jill</i>	<i>Jack</i> AND <i>Jill</i>
<i>doc_1</i>	1	1	1
<i>doc_2</i>	0	0	0
<i>doc_3</i>	1	0	0
<i>doc_4</i>	0	1	0
<i>doc_5</i>	1	0	0
<i>doc_6</i>	0	0	0
<i>doc_7</i>	0	0	0
<i>doc_8</i>	0	0	0

Processing a Boolean Query

Query: *Jack* **OR** *Jill*

	<i>Jack</i>	<i>Jill</i>	<i>Jack</i> OR <i>Jill</i>
<i>doc_1</i>	1	1	1
<i>doc_2</i>	0	0	0
<i>doc_3</i>	1	0	1
<i>doc_4</i>	0	1	1
<i>doc_5</i>	1	0	1
<i>doc_6</i>	0	0	0
<i>doc_7</i>	0	0	0
<i>doc_8</i>	0	0	0

Processing a Boolean Query

Query: *Jack* **AND** (*up* **OR** *down*)

	<i>up</i>	<i>down</i>	<i>up</i> OR <i>down</i>	<i>Jack</i>	<i>Jack</i> AND (<i>up</i> OR <i>down</i>)
<i>doc_1</i>	1	0	1	1	1
<i>doc_2</i>	0	0	0	0	0
<i>doc_3</i>	0	1	1	1	1
<i>doc_4</i>	0	0	0	0	0
<i>doc_5</i>	1	0	1	1	1
<i>doc_6</i>	0	0	0	0	0
<i>doc_7</i>	0	0	0	0	0
<i>doc_8</i>	0	0	0	0	0

Processing a Boolean Query

- Query: *Jack* **AND NOT** *Jill*

	<i>Jack</i>	<i>Jill</i>	NOT <i>Jill</i>	<i>Jack</i> AND NOT <i>Jill</i>
<i>doc_1</i>	1	1	0	0
<i>doc_2</i>	0	0	1	0
<i>doc_3</i>	1	0	1	1
<i>doc_4</i>	0	1	0	0
<i>doc_5</i>	1	0	1	1
<i>doc_6</i>	0	0	1	0
<i>doc_7</i>	0	0	1	0
<i>doc_8</i>	0	0	1	0

The Binary Full-text Representation

	<i>a</i>	<i>aardvark</i>	<i>abacus</i>	<i>abba</i>	<i>able</i>	...	<i>zoom</i>
<i>doc_1</i>	1	0	0	0	0	...	1
<i>doc_2</i>	0	0	0	0	1	...	1
::	::	::	::	::	::	...	0
<i>doc_m</i>	0	0	1	1	0	...	0

- Based on Zipf's law, this representation is not efficient
- There are lots of zeros!

Sparse Representation of an Inverted List

<i>a</i>	<i>aardvark</i>	<i>abacus</i>	<i>abba</i>	<i>able</i>	...	<i>zoom</i>
<i>df=3421</i>	<i>df=22</i>	<i>df=19</i>	<i>df=2</i>	<i>df=44</i>		<i>df=1</i>
1	33	2	33	66		54
33	56	10	150	134		
45	86	15		176		
::	::	::		::		
1022	1011	231		432		

- Variable-length inverted lists
 - represent only the 1's
 - each document has a unique identifier (docid)
 - *df* = number of documents in which the term appears at least once
- Why do we store the *df* in the index?

Unranked Boolean

- Retrieve the set of documents that match the boolean
- query (an “exact-match” retrieval model)
- Returns results in no particular order (ordered by date?)
- This is problematic with large collections
 - requires complex queries to reduce the result set to a manageable size
- Can we do better?

Ranked Boolean

<i>University</i>	<i>North</i>	<i>Carolina</i>	<i>UNC</i>
<i>df=6</i>	<i>df=4</i>	<i>df=3</i>	<i>df=5</i>
1, 4	1, 4	1, 4	1, 4
10, 1	10, 5	10, 5	10, 1
15, 2	16, 1	16, 1	16, 4
16, 1	68, 1		33, 2
33, 5			56, 10
67, 7			

- *docid* = document identifier
- *tf* = term frequency (# of times the term appears in the document)

Ranked Boolean

- Compute ranking scores based on term frequency
- Score computation:
 - A AND B: take the **minimum** frequency associated with expression A and expression B as the ranking score
 - A OR B: take the **sum** of frequencies associated with expression A and expression B

Ranked Boolean

- Query: *University* AND *North* AND *Carolina* OR *UNC*

<i>University</i>	<i>North</i>	<i>Carolina</i>	<i>Result_1</i>
<i>df=6</i>	<i>df=4</i>	<i>df=3</i>	<i>count=3</i>
1, 4	1, 4	1, 4	1, 4
10, 1	10, 5	10, 5	10, 1
15, 2	16, 1	16, 1	16, 1
16, 1	68, 1		
33, 5			
68, 7			

- AND: *min*
- OR: *sum*

Ranked Boolean

- Query: $(University \text{ AND } North \text{ AND } Carolina) \text{ OR } UNC$

<i>University</i>	<i>North</i>	<i>Carolina</i>	<i>UNC</i>	<i>Query</i>
<i>df=6</i>	<i>df=4</i>	<i>df=3</i>	<i>df=5</i>	<i>count=5</i>
1, 4	1, 4	1, 4	1, 4	1, 8
10, 1	10, 5	10, 5	10, 1	10, 2
15, 2	16, 1	16, 1	16, 4	16, 5
16, 1	68, 1		33, 2	33, 2
33, 5			56, 10	56, 10
68, 7				

Ranked Boolean

- Advantages:
 - same as unranked boolean: efficient, predictable, easy to understand, works well when the user knows what to look for
 - the user may be able to find relevant documents quicker and may not need to examine the entire result set
- Disadvantages:
 - same as unranked boolean: works well when the user knows what to look for

Best-Match Retrieval Models

- So far, we've discussed 'exact-match' models
- Now, we start discussing 'best-match' models
- Best-match models predict the degree to which a document is relevant to a query
- Ideally, this is would be expressed as **RELEVANT**(q,d)
- In practice, it is expressed as **SIMILAR**(q,d)
- How might you compute the similarity between q and d?