



CalPolyPomona

Retrieval Models I – Boolean Retrieval

Lecture 6

CS 4250 – Web Search and Recommender Systems

1

Retrieval Models

- Provide a mathematical framework for defining the search process
 - includes explanation of assumptions
 - basis of many ranking algorithms
- Progress in retrieval models has corresponded with improvements in effectiveness
- Sometimes trial and error
- Different theories about relevance

CS 4250 – Web Search and Recommender Systems

2

Relevance

- Complex concept that has been studied for some time
 - Many factors to consider
 - People often disagree when making relevance judgments
- Retrieval models make various assumptions about relevance to simplify problem
 - e.g., *topical* vs. *user* relevance
 - e.g., *binary* vs. *multi-valued* relevance

CS 4250 – Web Search and Recommender Systems

3

Retrieval Model Overview

- Original models
 - Boolean retrieval
 - Vector Space model
- Probabilistic Models
 - BM25
 - Language models
- Combining evidence
 - Inference networks
 - Learning to Rank

CS 4250 – Web Search and Recommender Systems

4

Boolean Retrieval

CS 4250 – Web Search and Recommender Systems

5

Boolean Retrieval

- Two possible outcomes for query processing
 - TRUE and FALSE
 - “exact-match” retrieval
- Query usually specified using Boolean operators
 - AND, OR, NOT
 - proximity operators also used

CS 4250 – Web Search and Recommender Systems

6

Searching by Numbers

- Sequence of queries driven by number of retrieved documents
 - e.g. “lincoln” search of news articles
 - president **AND** lincoln
 - president **AND** lincoln **AND NOT** (automobile **OR** car)
 - president **AND** lincoln **AND** biography **AND** life **AND** birthplace **AND** gettysburg **AND NOT** (automobile **OR** car)
 - president **AND** lincoln **AND** (biography **OR** life **OR** birthplace **OR** gettysburg) **AND NOT** (automobile **OR** car)
- Parentheses specify the order of operations
 - A **OR** (B **AND** C) does not equal (A **OR** B) **AND** C
- If query too broad, user adds **AND** constraints
- If query too narrow, user adds **OR** constraints

CS 4250 – Web Search and Recommender Systems

7

Examples

president AND lincoln would retrieve:

Ford Motor Company today announced that Darryl Hazel will succeed Brian Kelley as president of Lincoln Mercury.

What if we change it to:

president AND lincoln AND NOT (automobile OR car)

Would it retrieve

Lincoln's body departs Washington in a nine-car funeral train.

CS 4250 – Web Search and Recommender Systems

8

Boolean Retrieval

- Advantages
 - Results are predictable, relatively easy to explain
 - Many different features can be incorporated
 - Efficient processing since many documents can be eliminated from search
- Disadvantages
 - Effectiveness depends entirely on user
 - Simple queries usually don't work well (why?)
 - Complex queries are difficult (why?)

CS 4250 – Web Search and Recommender Systems

9

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

CS 4250 – Web Search and Recommender Systems

10

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.

CS 4250 – Web Search and Recommender Systems

11

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

CS 4250 – Web Search and Recommender Systems

12

Processing a Boolean Query

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

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

CS 4250 – Web Search and Recommender Systems

13

Processing a Boolean Query

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

	<i>Jack</i>	<i>Jill</i>	<i>Jack AND 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

CS 4250 – Web Search and Recommender Systems

14

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

CS 4250 – Web Search and Recommender Systems

15

Processing a Boolean Query

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

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

CS 4250 – Web Search and Recommender Systems

16

Processing a Boolean Query

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

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

CS 4250 – Web Search and Recommender Systems

17

Processing a Boolean Query

• Query: Jack **AND NOT** Jill

	Jack	Jill	Jack AND NOT Jill
doc_1	1	1	
doc_2	0	0	
doc_3	1	0	
doc_4	0	1	
doc_5	1	0	
doc_6	0	0	
doc_7	0	0	
doc_8	0	0	

CS 4250 – Web Search and Recommender Systems

18

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		
<i>doc_2</i>	0	0		
<i>doc_3</i>	1	0		
<i>doc_4</i>	0	1		
<i>doc_5</i>	1	0		
<i>doc_6</i>	0	0		
<i>doc_7</i>	0	0		
<i>doc_8</i>	0	0		

CS 4250 – Web Search and Recommender Systems

19

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

CS 4250 – Web Search and Recommender Systems

20

Fixed-length inverted list (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

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

CS 4250 – Web Search and Recommender Systems

21

Variable length inverted list

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

- Represent only the 1's
- Much more efficient

CS 4250 – Web Search and Recommender Systems

22

So far: 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? document ID? title?)
- This is problematic with large collections
 - requires complex queries to reduce the result set to a manageable size
- Can we do better?

CS 4250 – Web Search and Recommender Systems

23

Ranked Boolean

University	North	Carolina	UNC
df=6	df=4	df=3	df=5
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)*

CS 4250 – Web Search and Recommender Systems

24

Ranked Boolean

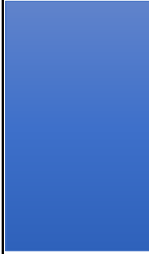
- Compute ranking scores based on term frequency
- Score computation:
 - A AND B: take the **minimum** frequency associated with expression A and expression B as the ranking score
 - A OR B: take the **sum** of frequencies associated with expression A and expression B

CS 4250 – Web Search and Recommender Systems

25

Ranked Boolean

- Query: (University AND North AND Carolina) OR UNC

University	North	Carolina	Result_I
df=6	df=4	df=3	count=3
1, 4	1, 4	1, 4	
10, 1	10, 5	10, 5	
15, 2	16, 1	16, 1	
16, 1	68, 1		
33, 5			
68, 7			

- AND -> **min**
- OR -> **sum**

CS 4250 – Web Search and Recommender Systems

26

Ranked Boolean

- Query: (University AND North AND Carolina) OR UNC

University	North	Carolina	Result_I
df=6	df=4	df=3	count=3
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*

CS 4250 – Web Search and Recommender Systems

27

Ranked Boolean

- Query: (University AND North AND Carolina) OR UNC

Result_I	UNC	Query
count=3	df=5	count=
1, 4	1, 4	
10, 1	10, 1	
16, 1	16, 4	
	33, 2	
	56, 10	

- AND -> *min*
- OR -> *sum*

CS 4250 – Web Search and Recommender Systems

28

Ranked Boolean

- Query: (University AND North AND Carolina) OR UNC

Result_I	UNC	Query
count=3	df=5	count=5
1, 4	1, 4	1, 8
10, 1	10, 1	10, 2
16, 1	16, 4	16, 5
	33, 2	33, 2
	56, 10	56, 10

- AND -> *min*
- OR -> *sum*

CS 4250 – Web Search and Recommender Systems

29

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: only works well when the user knows what to look for

CS 4250 – Web Search and Recommender Systems

30