

Information Retrivel: Assignment 1

Shifei Chen

Exercise 1.1

new → 1, 4
home → 1, 2, 3, 4
sales → 1, 2, 3, 4
top → 1
forecasts → 1
rise → 2, 4
in → 2, 3
july → 2, 3, 4
increase → 3

Exercise 1.2

a

Terms	Doc 1	Doc 2	Doc 3	Doc 3
breakthrough	1	0	0	0
drug	1	1	0	0
for	1	0	1	1
schizophrenia	1	1	1	1
new	0	1	1	1
approach	0	0	1	0
treatment	0	0	1	0
of	0	0	1	0
hopes	0	0	0	1
patients	0	0	0	1

b

breakthrough → 1
drug → 1, 2
for → 1, 3, 4
schizophrenia → 1, 2, 3, 4
new → 2, 3, 4
approach → 3
treatment → 3
of → 3
hopes → 4
patients → 4

Exercise 1.3

a

$\text{schizophrenia AND drug} = 1111 \text{ AND } 1100 = 1100$
So this query will return Doc 1 and Doc 2

b

for $\text{AND NOT}(\text{drug OR approach})$
 $= 1011 \text{ AND NOT } (1100 \text{ OR } 0010)$
 $= 1011 \text{ AND NOT } 1110$
 $= 1011 \text{ AND } 0001$
 $= 0001$
So this query will return Doc 4

Use the term-document incidence matrix in Figure 1.1 to return the documents related to the query “(Brutus OR Caesar) AND NOT(Antony OR Cleopatra)”

$(\text{Brutus OR Caesar}) \text{ AND NOT}(\text{Antony OR Cleopatra})$
 $= (110100 \text{ OR } 110111) \text{ AND NOT } (110001 \text{ OR } 100000)$
 $= 110111 \text{ AND NOT } 110001$
 $= 110111 \text{ AND } 001110$
 $= 000110$
So this query will return Hamlet and Othello

Can you find a general way to process arbitrary Boolean queries as the query in previous exercise?

I think the general process for boolean queries should be like

1. Build the term-document incidence matrix for this particular document collection
2. Replace all of the terms in the query with their binary representations
3. Calculate the binary result by executing binary operations, such as AND, OR or NOT
4. Look up in the term-document incidence matrix to figure out which document is in the final binary result

It might also be a good idea to convert the query into a postfix expression with the Shunting Yard Algorithm¹. For example the query in the previous question “(Brutus OR Caesar) AND NOT(Antony OR Cleopatra)” would be turned into “Brutus Caesar OR Antony Cleopatra OR NOT AND”.

¹<http://www.oxfordmathcenter.com/drupal7/node/628>