

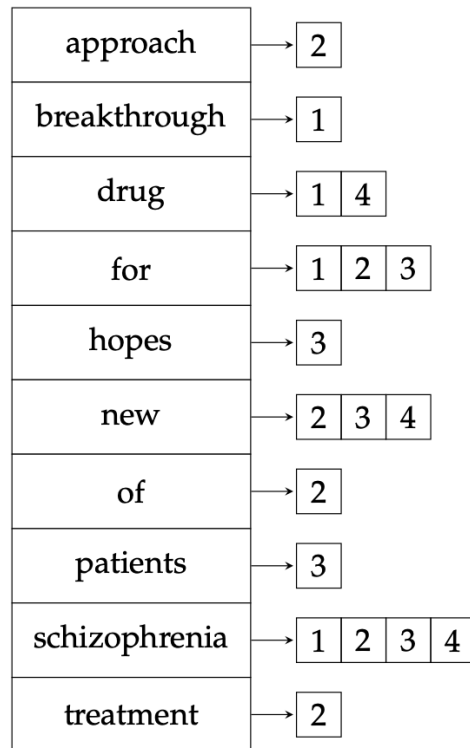
Assignment 1 Solution

Problem 1

1. Draw the term-document incidence matrix for this document collection.

	Doc1	Doc2	Doc3	Doc4
approach	0	1	0	0
breakthrough	1	0	0	0
drug	1	0	0	1
for	1	1	1	0
hopes	0	0	1	0
new	0	1	1	1
of	0	1	0	0
patients	0	0	1	0
schizophrenia	1	1	1	1
treatment	0	1	0	0

2. Draw the inverted index representation for this collection.



3. What are the returned results for these queries:

(a) schizophrenia AND drug

Doc1, Doc4

(b) for AND NOT (drug OR approach)

(drug OR approach): **Doc1, Doc2, Doc4.**

for AND NOT (drug OR approach):

Doc3

Problem 2

1. Write out a postings merge algorithm, in the style of Figure 1.6 in IIR, for an x OR y query.

```
1: procedure UNION( $p_1, p_2$ )
2:    $answer \leftarrow \langle \rangle$ 
3:   while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$  do
4:     if  $docID(p_1) = docID(p_2)$  then
5:       ADD( $answer, docID(p_1)$ )
6:        $p_1 \leftarrow next(p_1)$ 
7:        $p_2 \leftarrow next(p_2)$ 
8:     else if  $docID(p_1) < docID(p_2)$  then
9:       ADD( $answer, docID(p_1)$ )
10:       $p_1 \leftarrow next(p_1)$ 
11:    else
12:      ADD( $answer, docID(p_2)$ )
13:       $p_2 \leftarrow next(p_2)$ 
14:   while  $p_1 \neq \text{NIL}$  do
15:     ADD( $answer, docID(p_1)$ )
16:      $p_1 \leftarrow next(p_1)$ 
17:   while  $p_2 \neq \text{NIL}$  do
18:     ADD( $answer, docID(p_2)$ )
19:      $p_2 \leftarrow next(p_2)$ 
```

2. Write out a postings merge algorithm, in the style of Figure 1.6 in IIR, for an x AND NOT y query.

```
1: procedure MINUS( $p_1, p_2$ )
2:    $answer \leftarrow \langle \rangle$ 
3:   while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$  do
4:     if  $docID(p_1) = docID(p_2)$  then
5:        $p_1 \leftarrow next(p_1)$ 
6:        $p_2 \leftarrow next(p_2)$ 
7:     else if  $docID(p_1) < docID(p_2)$  then
8:       ADD( $answer, docID(p_1)$ )
9:        $p_1 \leftarrow next(p_1)$ 
10:    else
11:       $p_2 \leftarrow next(p_2)$ 
12:   while  $p_1 \neq \text{NIL}$  do
13:     ADD( $answer, docID(p_1)$ )
14:      $p_1 \leftarrow next(p_1)$ 
```

Problem 3

The estimated (worst-case) lengths for the results of the OR operations are:

- (tangerine OR trees): $87009 + 316812 = 403821$
- (marmalade OR skies): $107913 + 271658 = 379571$
- (kaleidoscope OR eyes): $46653 + 213312 = 259965$

These operations must be done first (in any order), and then the AND operations should be applied on the results in ascending order of size. Therefore, the recommended order is:

```
(1) ← tangerine OR trees
(2) ← marmalade OR skies
(3) ← kaleidoscope OR eyes
(4) ← (3) AND (2)
return (4) AND (1)
```

Problem 4

Suppose that the list of all documents is long, with length D .

The naive evaluation would first compute NOT y , and then x OR (NOT y):

- 1: $result \leftarrow \text{MINUS}(\text{documents}, y)$
- 2: **return** UNION($x, result$)

This evaluation would require one scan of the document list to compute NOT y . Most of the time, the postings list for y is much shorter than the list of documents, so the result of NOT y is most of the document list. Computing the union of x with NOT y requires a complete scan of the result of NOT y , so there is a second scan of length D . In total, the list of all documents is scanned twice. This makes the query expensive.

To be more efficient, first apply DeMorgan's Law to the query:

$$x \text{ OR NOT } y = \text{NOT} (\text{NOT} (x \text{ OR NOT } y)) = \text{NOT} (\text{NOT } x \text{ AND } y) = \text{NOT} (y \text{ AND NOT } x)$$

Then the algorithm is:

- 1: $result \leftarrow \text{MINUS}(y, x)$
- 2: **return** MINUS($\text{documents}, result$)

Computing y AND NOT x is linear in the lengths of postings for x and y , which are usually much shorter than D . Then computing NOT the result requires just one scan of the documents list. So in total, this algorithm requires only one scan of the list of all documents, which is more efficient.