

7-Feb-2024

# Information Retrieval and Web Search

## Assignment - 2

Utkay Bhaskar Valapadashu  
11696364

|Ans :-

a) Brutus AND NOT Caesar

Here, Rather than gathering documents present in both postings lists, the approach is to gather those appearing in the first list but not in the second.

The time complexity for this query would be same  $O(\alpha \times \beta)$ .

b) Brutus OR NOT Caesar

Assuming the requirement is to provide a comprehensive list of all documents meeting the query criteria, the time complexity is  $O(N)$ ,

where N represents the total number of documents in the collection. This is because the length of the result list can potentially reach  $N$ , unrestricted by the length of the postings lists.

2Ans:-

(i) (tangerine OR trees) AND (marmalade OR skies)  
AND (Kaleidoscope OR eyes)

Estimate the size of each OR operation by  
adding the frequencies of its disjuncts:

$$1 \geq (\text{tangerine OR trees}) = 316812 + 46653 = 363465$$

$$2 \geq (\text{marmalade OR skies}) = 107913 + 271658 = 379571$$

$$3 \geq (\text{Kaleidoscope OR eyes}) = 87009 + 213312 = 300321$$

To process the query in increasing order of the  
size of each disjunctive term.  $\Rightarrow$ :

The recommended query processing order is:

(Kaleidoscope OR eyes) AND (tangerine OR trees)  
AND (marmalade OR skies)

$\textcircled{3} \rightarrow \textcircled{1} \rightarrow \textcircled{2}$

(ii) tangerine AND (NOT marmalade) AND (NOT trees)

Estimate the size of each term taking  
into account the NOT operation

$$1 \geq \text{tangerine} = 46653$$

2  $\geq$  NOT marmalade = Total no of documents -

$$\text{frequency of "marmalade"} = 500000 - 107913$$

$$= 392087$$

$$\begin{aligned}
 3. > \text{NOT toes} &= \text{Total no of documents} - \text{Frequency of "toes"} \\
 &= 500000 - 316812 \\
 &= 183188
 \end{aligned}$$

To process the query in increasing order of the size of each term. So, the recommended query processing order is:

(NOT toes) AND (NOT marmalade) AND  
 (tangerine)

③ → ② → ①

Ans:-

Let's consider an example to understand:

Suppose we have three terms with the following postings list sizes:

- Term 1: Size = 100
- Term 2: Size = 105
- Term 3: Size = 110

Now, let's say the intersections b/w these posting lists as follows:

1.> intersection of 81 and 82 : 100 steps

2.> intersection of result with 83 : 100 steps (

assuming no intersection with 83)

$$\text{Total steps} = 100 + 105 + 100 + 110 = 315 \text{ steps}$$

However, if we process the terms in a different order (81, 83, 82) as you suggested, we would perform

the following step:

1. Intersection of 81 and 83 : 0 steps (

assuming no intersection with 83)

2. No further intersection needed with 82 ( since intersection with 83 is empty )

$$\text{Total steps} : 100 + 110 + 0 + 0 = 210 \text{ steps.}$$

So processing postings lists in order of size is not always optimal, as demonstrated by this example.

4 Ans :-

postings merge algorithm for  $\text{xe OR Y}$  query.

$\text{UNION}(\text{xe}, \text{y})$

1.  $\text{answer} \leftarrow ()$
2. while  $\text{xe} \neq \text{NIL}$  and  $\text{y} \neq \text{NIL}$
3. do if  $\text{docID}(\text{xe}) = \text{docID}(\text{y})$
4. ~~then ADD(answer, docID(xe))~~
5.  $\hookrightarrow$  then ~~then~~  $\text{ADD}(\text{answer}, \text{docID}(\text{xe}))$   
 $\text{x} \leftarrow \text{next}(\text{xe})$
6.  $\text{y} \leftarrow \text{next}(\text{y})$
7. else if  $\text{docID}(\text{xe}) < \text{docID}(\text{y})$
8.  $\text{xe} \leftarrow \text{next}(\text{xe})$
9. else  $\text{ADD}(\text{answer}, \text{docID}(\text{y}))$
10.  $\text{y} \leftarrow \text{next}(\text{y})$
11. return ( $\text{answer}$ )

- - - - - X - - - - -