

Ans 2 (a)

Sahil Tindel  
18CS10048

Ans 2 (a) Advantage: Lookup is faster

Disadvantage: Cannot handle wild card queries easily.

(b) There can only be one original vocabulary term from which it was derived because of the \$ terminal character.

For eg. book & okbo would have given ~~same~~ same permuterm vocabulary if we don't use \$.

(c) \$red

(d) For trigram, we would search for terms:

\$re

red

& take intersection.

~~The term reduction~~

"The term 'recovered' would match the boolean query but not permuterm query.

$$(e) \text{ Jaccard Coeff} = \frac{|A \cap B|}{|A \cup B|}$$

For a bigram index is

bar :- ba, ar

barbaric :- ba, ar, rb, ba, ar, ri, ic

Removing repetition,

$$(bar) \cap (barbaric) \Rightarrow ba, ar$$

$$(bar \cup barbaric) \Rightarrow ba, ar, rb, ri, ic$$

$$\therefore \text{Jaccard coeff} = \frac{2}{5} = \underline{0.4}$$

A3

(a) Reading time for collection of data from disk:

$$= (\text{no. of tokens})^* (\text{bytes per token})^* (\text{time for transfer each byte}) + \text{disk seek time}$$

where disk seek time

$$= (4 \times 10^7)(8)(2 \times 10^{-8}) + 5 \times 10^{-3}$$

$$= 6.4 + 0.005 = 6.4058$$



(b) Total searching time = (no. of blocks)<sup>\*</sup> / (N log N)<sup>\*</sup> (time for low level operation)

$$= (20) (N \log N) (10^{-8})$$

$$= 2 \times 10^{-7} \times N \log N$$

where  $N$  = no. of documents  
(not given in ques)

(c) Total time for writing sorted blocks to disk

$$= (\text{no. of blocks})^* (\text{time to write single block})$$

where

time to write single block = (No. of terms per block + no. of positing in each block)<sup>\*</sup> (size of term id)<sup>\*</sup>  
(time to transfer 1 byte data).

$$= (20) \left( \frac{250000}{20} + \frac{4 \times 10^7}{20} \right) (4) (2 \times 10^{-8})$$

$$= (25 + 4000) (4) (2)$$

$$= 32200 \text{ seconds.}$$

(d) Assuming buffer for each block to contain  $\frac{1}{10^{th}}$  of the block content.

then, required time

$$= \text{disk seek time} + \text{time for reading part of sorted posix list from disk} + \text{time for writing final index.}$$

$$\text{now, disk seek time} = (\text{No of blocks})^* (\text{time to read full block given we have } \frac{1}{10^{th}} \text{ of block})^* (\text{avg seek time})^* 2$$

&

$$\begin{aligned} \text{Time for loading part of sorted posix list} &= \text{time for writing final index} \\ &= (\text{no of terms} + \text{total posting}) \\ &\quad * (\text{size of term id}) * (\text{time to transfer single byte}) - \end{aligned}$$

$$\begin{aligned} \therefore \text{reqd time} &= 20^* \left( \frac{4 \times 10^7 \times 8 \times 2 \times 10^{-8} \times 10}{20} \right) \\ &\quad + (5 \times 10^{-3})^* 2 + 20^* \\ &\quad + (250000 + \end{aligned}$$



Ans 4

		Judge 1	
		R	N
Judge 2	R	2	5
	N	5	3

$$(9) P(a) = P(\text{Agreement}) = \frac{2+3}{15} = \frac{1}{3}$$

$$P(NR) = \frac{5+5+3+3}{30} = \frac{16}{30} = \frac{8}{15}$$

$$P(R) = \frac{5+5+2+2}{30} = \frac{14}{30} = \frac{7}{15}$$

$$P(e) = P(\text{Agree Chance}) = P(NR)^2 + P(R)^2$$
$$= \frac{113}{225}$$

$$K = \frac{P(a) - P(e)}{1 - P(e)} = \frac{0.33 - 0.502}{1 - 0.502}$$

$$K = -0.3453$$

(b) ~~if both~~ Retrieved = 4, 5, 6, 7, 3, 8  
Relevant IDs: 3, 4  
Precision =  $\frac{\text{Relevant}}{\text{Retrieved}}$

$$= \frac{2}{6} = \frac{1}{3}$$

Recall =  $\frac{\text{Retrieved}}{\text{Relevant}}$

$$= \frac{2}{2} = 1$$

(c) Relevant IDs: 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14.

Precision =  $\frac{6}{6} = 1$

Recall =  $\frac{6}{12} = \frac{1}{2}$



(d) w

Bonus :: We can show results for different queries to users & ~~crowd~~ crowd source their opinions about the relevance of the result.

We will show them various results & use the  $k$ -~~rank~~ rank of that ~~product~~ product in the result by our IR system.

In other words, ~~Then~~ we will get the relevance score of a product & then use it to find score of our IR system.