

Ques. a. d1: New York Times d2: New York Post  
d3: Los Angeles Times

Total no. of documents:  $N=3$ , Hence idf values of terms arranged alphabetically are as follows:

angels :-  $\log_2 \frac{3}{1} = 1.584$

los :-  $\log_2 \frac{3}{2} = 0.584$

new :-  $\log_2 \frac{3}{2} = 0.584$

post :-  $\log_2 \frac{3}{1} = 1.584$

times :-  $\log_2 \frac{3}{2} = 0.584$

york :-  $\log_2 \frac{3}{2} = 0.584$

For all documents, we calculate the tf scores as follows:-

	angels	los	new	post	times	york
d1	0	0	1	1	1	1
d2	0	0	0	0	0	1
d3	1	1	0	0	1	0

Now, multiplying tf-idf values we get term vectors for documents

$d_3 := W_{\text{angels}} \rightarrow 1 * 1.584 = 1.584$

$W_{\text{los}} \rightarrow 0 * 1.584 = 0$

$W_{\text{times}} \rightarrow 0 * 0.584 = 0$

$d_2 := W_{\text{new}} \rightarrow 1 * 0.584 = 0.584$

$W_{\text{york}} \rightarrow 1 * 0.584 = 0.584$

$W_{\text{post}} \rightarrow 1 * 1.584 = 1.584$

$d_1 := W_{\text{new}} \rightarrow 1 * 0.584 = 0.584$

$W_{\text{york}} \rightarrow 1 * 0.584 = 0.584$

$W_{\text{times}} \rightarrow 1 * 0.584 = 0.584$

$d_1 = \langle 0.584, 0.584, 0.584 \rangle$

$d_2 = \langle 0.584, 0.584, 1.584 \rangle$

$d_3 = \langle 1.584, 1.584, 0.584 \rangle$

Q1 b. When computing tf-idf values for the query terms we divide the frequency by the maximum frequency (2) and multiply with the idf values.

$$2 \quad 0 \quad 0 \quad (2/2) * 0.584 = 0.584 \quad 0 \quad (1/2) * 0.584 = 0.292$$

We calculate the length of each document  $q$  of the query :-

$$d_1 = \sqrt{(0.584)^2 + (0.584)^2 + (0.584)^2} = 1.011$$

$$d_2 = \sqrt{(0.584)^2 + (1.584)^2 + (0.584)^2} = 1.786$$

$$d_3 = \sqrt{(1.584)^2 + (1.584)^2 + (1.584)^2} = 2.316$$

$$q = \sqrt{(0.584)^2 + (0.292)^2} = 0.652$$

Then the similarity values are:-

$$\cos \text{sim}(d_1, q) = \frac{[0 + 0 + (0.584)^2 + 0 + 0.584 * 0.292 + 0]}{(1.011 * 0.652)} = 0.776$$

$$\cos \text{sim}(d_2, q) = \frac{[0 + 0 + 0.584 * 0.584 + 0 + 0 + 0]}{1.786 * 0.652} = 0.292$$

$$\cos \text{sim}(d_3, q) = \frac{[0 + 0 + 0 + 0 + 0.584 * 0.292 + 0]}{2.316 * 0.652} = 0.112$$

$$+82.0 = +82.0 * 1 \leftarrow \text{unit w}$$

$$+82.0 = +82.0 * 1 \leftarrow \text{unit w}$$

$$+82.0 = +82.0 * 1 \leftarrow \text{unit w}$$

$$\langle +82.0, +82.0, +82.0 \rangle = 1b$$

$$\langle +82.1, +82.0, +82.0 \rangle = 5b$$

$$\langle +82.0, +82.1, +82.1 \rangle = 5b$$

	$g_1$	$g_2$	$d_1$	$d_2$	$j = (d_2) \cap g_1$
aerodinax	0	1	0	0	
banana	1	1	1	1	$g_1 \cap g_2 \cap d_1 \cap d_2$
campus	0	0	0	1	
columbians	0	1	0	0	
Cruz	0	0	1	1	
mascot	0	0	0	1	
mountains	0	1	0	0	
santa	0	0	1	1	
slug	1	1	1	0	

The Reinforcement feedback mechanism :- ( $\alpha, \beta, \gamma$ )

$$\vec{q}_m = \alpha \cdot \vec{q}_o + \beta \frac{\sum_{dj \in ED_f} \vec{d_j}}{|ED_f|} + \gamma \frac{\sum_{dj \in D_{nr}} \vec{d_j}}{|D_{nr}|}$$

$\alpha = \beta = \gamma = 1$  i.e.  $\vec{q}_m = \vec{q}_o + \vec{d}_f + \vec{d}_{nr}$

(b) reward reinforced through a vector

0.0	0.0	0	ab
40.0	3.0	1	bb
80.0	98.0	2	abb
-51.0	25.0	3	bab
51.0	98.0	4	bb
31.0	92.0	5	ab
21.0	84.0	6	bb
21.0	92.0	7	ab
45.0	92.0	8	bb
25.0	92.0	9	ab

it has trouble in learning figures, so not

possible -> {last, ab, bb}

total = {last, last, ab, last, abb, last, last, last}

$$Q3. RET(s1) = \{d_2, d_5, d_{150}, d_{250}, d_{11}, d_{33}, d_{50}, d_{600}, d_{800}, d_{520}\}$$

$$RET(s2) = \{d_{250}, d_{400}, d_{150}, d_{210}, d_{999}, d_3, d_{501}, d_{800}, d_{205}, d_{300}\}$$

a) Precision =  $\frac{\text{No. of relevant documents retrieved}}{\text{No. of documents retrieved}}$

Recall =  $\frac{\text{No. of relevant documents retrieved}}{\text{Total no. of relevant documents (REL)}}$

The set REL has 25 relevant documents in this case

$$REL = \{d_1, d_5, d_6, d_{10}, d_{88}, d_{150}, d_{200}, d_{210}, d_{250}, d_{300}, d_{405}, d_{450}, d_{472}, d_{500}, d_{501}, d_{530}, d_{545}, d_{590}, d_{600}, d_{700}, d_{720}, d_{800}, d_{1000}\}$$

For  $s_1$ , checking for retrieved documents are relevant or not;  
 $\{d_2, d_{11}, d_{33}, d_{50}, d_{520}\}$  = Non-relevant  
 $\{d_5, d_{150}, d_{250}, d_{600}, d_{300}\}$  = relevant

Retrieved documents	Retrieved docs so far	Precision	Call
d <sub>2</sub>	0	0.0	0.0
d <sub>5</sub>	1	0.5	0.04
d <sub>150</sub>	2	0.67	0.08
d <sub>250</sub>	3	0.75	0.12
d <sub>11</sub>	3	0.60	0.12
d <sub>33</sub>	3	0.50	0.12
d <sub>50</sub>	3	0.43	0.12
d <sub>600</sub>	4	0.50	0.16
d <sub>500</sub>	5	0.56	0.20
d <sub>520</sub>	5	0.50	0.20

For  $s_2$ , checking if documents are relevant or not;

$\{d_{999}, d_3, d_{205}\}$  = Non-relevant

$\{d_{250}, d_{400}, d_{150}, d_{210}, d_{201}, d_{800}, d_{300}\}$  = Relevant

Retrieved documents	Relevant Ref. so far	Precision	Recall
d <sub>250</sub>	1	1.0	0.04
d <sub>400</sub>	2	1.0	0.08
d <sub>500</sub>	3	1.0	0.12
d <sub>150</sub>	4	1.0	0.16
d <sub>210</sub>	4	0.80	0.16
d <sub>999</sub>	4	0.67	0.16
d <sub>3</sub>	5	0.71	0.20
d <sub>501</sub>	6	0.75	0.24
d <sub>800</sub>	6	0.67	0.24
d <sub>205</sub>	7	0.70	0.28
d <sub>300</sub>	7	0.70	0.28
d <sub>10</sub>	7	0.70	0.28

b) F1 Measure :-

$$F1 = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

Computing the F1 measure for each document retrieved for both the systems :-

For system s<sub>1</sub>, F1 is as follows :-

For s<sub>1</sub>,

Retrieved documents	Precision	Recall	F1
d <sub>2</sub>	0.0	0.0	0.0
d <sub>5</sub>	0.5	0.04	0.07
d <sub>150</sub>	0.67	0.08	0.14
d <sub>250</sub>	0.75	0.12	0.21
d <sub>50</sub>	0.60	0.12	0.20
d <sub>33</sub>	0.50	0.12	0.19
d <sub>50</sub>	0.43	0.12	0.19
d <sub>500</sub>	0.50	0.16	0.24
d <sub>500</sub>	0.56	0.20	0.29
d <sub>520</sub>	0.50	0.20	0.29

For S2:-

Retrieved Documents	Precision	Recall	F1
d <sub>250</sub>	1.0	0.04	0.08
d <sub>400</sub>	1.0	0.08	0.15
d <sub>50</sub>	1.0	0.12	0.21
d <sub>210</sub>	1.0	0.16	0.23
d <sub>999</sub>	0.80	0.16	0.27
d <sub>3</sub>	0.67	0.16	0.26
d <sub>501</sub>	0.71	0.20	0.31
d <sub>800</sub>	0.75	0.24	0.36
d <sub>205</sub>	0.67	0.24	0.36
d <sub>300</sub>	0.70	0.28	0.40

c) To determine which system is better, we can analyze the F1 measures graph. Based on F1 scores, S2 system generally performs better at combining precision & recall, especially as more documents are retrieved.

Q4. a) i) cosine similarity :- Between query Q and a document D<sub>i</sub> is calculated as,

$$\text{Cosine Similarity} = \frac{Q \cdot D_i}{|Q| * |D_i|}$$

$$\text{Query Vector } Q = [2, 1, 1, 0, 2, 0, 3, 0]$$

$$|Q| = \sqrt{2^2 + 1^2 + 1^2 + 2^2 + 3^2} = \sqrt{19}$$

$$\text{For Each document ; } D_{OC1} = [0, 3, 4, 0, 0, 2, 4, 0]$$

$$D_{OC2} = [5, 5, 0, 0, 4, 0, 4, 3] \quad D_{OC3} = [3, 0, 4, 3, 4, 0, 0, 5]$$

$$D_{OC4} = [0, 7, 0, 3, 2, 0, 4, 3] \quad D_{OC5} = [0, 1, 0, 0, 0, 5, 4, 2]$$

$$D_{OC6} = [2, 0, 2, 0, 0, 4, 0, 1] \quad D_{OC7} = [3, 5, 3, 4, 0, 0, 4, 2]$$

$$D_{OC8} = [0, 3, 0, 0, 0, 4, 4, 2] \quad D_{OC9} = [0, 0, 3, 3, 3, 0, 0, 1]$$

$$D_{OC10} = [0, 5, 0, 0, 0, 1, 4, 2]$$

For DOC1 :-  $\Phi \cdot D1 \Rightarrow 0+3+4+0+0+0+12+0$   
 $\Rightarrow 19$

$$\text{cosine similarity} = \frac{19}{\sqrt{19} \times \sqrt{55}} = \frac{19}{29.14} = 0.6746$$

$$\text{For DOC2 :- cosine similarity} = \frac{35}{\sqrt{19} \times \sqrt{91}} = \frac{35}{41.45} = 0.844$$

b) Binary term weights with Dot Product :-

$$DOC1 = [0, 1, 1, 0, 0, 1, 1, 0] \quad |D1| = 2$$

$$DOC2 = [1, 1, 0, 0, 1, 0, 1, 1] \quad |D2| = \sqrt{5}$$

$$DOC3 = [1, 0, 1, 1, 1, 0, 0, 1] \quad |D3| = \sqrt{5}$$

$$DOC4 = [0, 1, 0, 1, 1, 0, 1, 0] \quad |D4| = \sqrt{5}$$

$$DOC5 = [0, 1, 0, 0, 0, 1, 1, 1] \quad |D5| = 2$$

$$DOC6 = [1, 0, 1, 0, 0, 1, 0, 1] \quad |D6| = 2$$

$$DOC7 = [1, 1, 1, 1, 0, 0, 1, 1] \quad |D7| = \sqrt{6}$$

$$DOC8 = [0, 1, 0, 0, 0, 1, 1, 1] \quad |D8| = 2$$

$$DOC9 = [0, 0, 1, 1, 1, 0, 0, 1] \quad |D9| = 2$$

$$DOC10 = [0, 1, 0, 0, 0, 1, 1, 1] \quad |D10| = 2$$

$$\vec{\Phi} = [2, 1, 1, 0, 2, 0, 3, 0]$$

Dot Product of binary term weights arranged in decreasing order are as follows:-

$$\begin{aligned} & DOC2: 4 \leftarrow DOC7: 4 \leftarrow DOC1: 3 \leftarrow DOC3: 3 \leftarrow DOC4: 3 \\ & \leftarrow DOC5: 2 \leftarrow DOC5: 2 \leftarrow DOC6: 2 \leftarrow DOC8: 2 \leftarrow \\ & DOC9: 2 \leftarrow DOC10: 2 \end{aligned}$$

$$\Phi_0 = \text{avg vec}$$

Relevant Documents vector sum, computing average

$$\Phi_0 = [1, 1, 0, 0, 0, 0, 0, 0, 0, 0]$$

Vector for relevant documents :-

$$\text{DOC4} = [0, 1, 0, 1, 0, 0, 1, 0, 0, 0]$$

$$\text{DOC5} = [1, 1, 0, 0, 1, 0, 0, 0, 0, 0]$$

$$\text{DOC6} = [1, 1, 0, 0, 0, 1, 0, 0, 0]$$

$$\text{DOC7} = [1, 1, 0, 0, 0, 0, 0, 0, 0, 1]$$

$$\text{SUM}_R = [3, 4, 0, 1, 1, 0, 2, 0, 0, 1]$$

$$\text{Avg}_R = \frac{1}{4} [3, 4, 0, 1, 1, 0, 2, 0, 0, 1]$$

$$\text{Avg}_R = [0.75, 1.0, 0.25, 0.25, 0, 0.5, 0, 0, 0.25]$$

Non-relevant documents vector sum, computing average vector for non-relevant documents :-

$$\text{DOC1} = [1, 1, 1, 0, 1, 1, 0, 0, 0, 0]$$

$$\text{DOC2} = [1, 1, 0, 0, 1, 0, 0, 0, 0, 0]$$

$$\text{DOC3} = [1, 0, 1, 0, 0, 0, 0, 0, 0, 0]$$

$$\text{Sum}_{NR} = [3, 2, 2, 0, 2, 0, 0, 0, 0, 0]$$

$$\text{Avg}_{NR} = \frac{1}{3} [3, 2, 2, 0, 2, 0, 0, 0, 0, 0]$$

$$\text{Avg}_{NR} = [1, 0.66, 0.66, 0, 0.66, 0, 0, 0, 0, 0]$$

substituting everything in Rocchio's formula :-

$$Q_{opt} = 1.0 * Q_0 + 1.0 * \text{Avg}_R - 0.5 * \text{Avg}_{NR}$$

$$Q_0 = [1, 1, 0, 0, 0, 0, 0, 0, 0, 0] + [0.75, 1, 0, 0.25, 0.25, 0, 0.5, 0, 0, 0] \\ + [0.5, 0.33, 0.33, 0, 0.33, 0, 0, 0, 0, 0]$$

$$\therefore Q_0 = [1.25, 1.67, -0.33, 0.25, -0.08, 0, 0.5, 0, 0, 0.25]$$

$$[1.25, 1.67, -0.33, 0.25, -0.08, 0, 0.5, 0, 0, 0.25] = Q_{opt}$$

Rocchio's formula :-  
Given query + relevance document - non-relevance document

$$\text{For } \text{DOC}3 : \text{Cosine similarity} = \frac{18}{4.35 \times 8.66} = 0.477$$

$$\text{For } \text{DOC}4 = 0.5657$$

$$\text{For } \text{DOC}5 = 0.4397$$

$$\text{For } \text{DOC}6 = 0.2753$$

$$\text{For } \text{DOC}7 = 0.1711$$

$$\text{For } \text{DOC}8 = 0.5130$$

$$\text{For } \text{DOC}9 = 0.3902$$

$$\text{For } \text{DOC}10 = 0.4994$$

For cosine similarity, decreasing order is,

$\text{DOC}2, \text{DOC}7, \text{DOC}1, \text{DOC}4, \text{DOC}8, \text{DOC}10, \text{DOC}3, \text{DOC}5, \text{DOC}9, \text{DOC}6$

ii) Dice's coefficient :-  $\frac{2 \times (q \cdot Di)}{T\phi + 1Di}$ ,

The Dice's coefficient between query & each document is as follows :-

$$\text{DOC}1 = 1.7273, \quad \text{DOC}2 = 2.3333, \quad \text{DOC}3 = 1.2857$$

$$\text{DOC}4 = 1.6429, \quad \text{DOC}5 = 1.2381, \quad \text{DOC}6 = 0.6667$$

$$\text{DOC}7 = 1.7333, \quad \text{DOC}8 = 1.3626, \quad \text{DOC}9 = 0.9474$$

$$\text{DOC}10 = 1.4167$$

In decreasing order of dice's coefficient,

DOC2

DOC7

DOC1

DOC4

DOC10

DOC8

DOC3

DOC5

DOC9

DOC6

Q5. Rocchio's formula :-

$$q_1 = \alpha q_0 + \beta \frac{1}{|D_R|} \sum_{d \in D_R} D_d - \gamma \frac{1}{|D_N|} \sum_{d \in D_N} D_d$$

$$q_{opt} = \alpha q_0 + \beta \frac{1}{|D_R|} \sum_{j \in D_R} D_j - \gamma \frac{1}{|D_N|} \sum_{j \in D_N} D_j$$

$$\beta = 1.0, \gamma = 0.5$$

Relevant Docs = {DOC4, DOC5, DOC6, DOC7}

Non-Relevant Docs = {DOC1, DOC2, DOC3}

Terms :- dog, race, greyhound, track, kitting, iditarod, sled, husky, malamute, alaska.

	Doc1	Doc2	Doc3	Doc4	Doc5	Doc6	Doc7
dog	1	1	1	1	1	1	1
race	1	1	0	1	1	1	1
greyhound	1	1	0	0	0	0	0
track	1	0	0	0	0	0	0
kitting	1	1	0	0	0	1	0
iditarod	0	0	0	1	0	0	1
sled	0	0	0	1	1	1	0
husky	0	0	0	0	1	0	0
malamute	0	0	0	0	1	0	0
alaska	0	0	0	0	0	1	1

The significance increase in weights of terms like 'dog' & 'race' after applying Rocchio's relevance feedback is because they appear frequently in the relevant documents which were positively marked by the user.

The decrease in term weights typically occurs for terms that are more common in non-relevant documents. This reflects the idea that terms frequently appearing in non-relevant documents are less likely to be useful for improving the search results.

## Q6. Ranked Document list :-

- 1) DOC1
  - 2) DOC2 (Relevant)
  - 3) DOC3 (Relevant)
  - 4) DOC4 (Relevant)
  - 5) DOC5
  - 6) DOC6
  - 7) DOC7
  - 8) DOC8 (Relevant)
  - 9) DOC9
  - 10) DOC10
- Total Relevant Docs = 5

Calculation precision for each rank position, then use that to find interpolated precision at specified recall

Position 1 :-  $0 \text{ relevant} / 1 \text{ retrieved} = 0.0$

Position 2 :-  $1 \text{ relevant} / 2 \text{ retrieved} = 0.5$

Position 3 :-  $2 \text{ relevant} / 3 \text{ retrieved} = 0.67$

Position 4 :-  $3 \text{ relevant} / 4 \text{ retrieved} = 0.75$

Position 5 :-  $3 \text{ relevant} / 5 \text{ retrieved} = 0.6$

Position 6 :-  $3 \text{ relevant} / 6 \text{ retrieved} = 0.5$

Position 7 :-  $3 \text{ relevant} / 7 \text{ retrieved} = 0.43$

Position 8 :-  $4 \text{ relevant} / 8 \text{ retrieved} = 0.5$

Position 9 :-  $4 \text{ relevant} / 9 \text{ retrieved} = 0.44$

0.0, 0.1 = minimum & maximum (0.0 to 0.5)

0.2 = Maximum & Minimum (0.0 to 0.5)

0.3 = Middle & Minimum through 0.5

0.4 = Middle & Through 0.5

6. Number of documents retrieved = 10  
 Number of relevant documents = 5  
 Number of relevant documents retrieved = 4 (positions 2, 3, 5, 8)

Rank Position	Precision	Recall
1	$0/1 = 0$	$0/5 = 0$
2	$1/2 = 0.5$	$1/5 = 0.2$
3	$2/3 = 0.67$	$2/5 = 0.4$
4	$3/4 = 0.75$	$3/5 = 0.6$
5	$3/5 = 0.60$	$3/5 = 0.6$
6	$3/6 = 0.5$	$3/5 = 0.6$
7	$3/7 = 0.43$	$3/5 = 0.6$
8	$4/8 = 0.5$	$4/5 = 0.8$
9	$4/9 = 0.44$	$4/5 = 0.8$
10	$4/10 = 0.4$	$4/5 = 0.8$

Interpreted Precision values:

for each recall level we take the maximum precision observed recall level greater than or equal to that recall level

Recall Level

Interpolated Precision

0.75

0.7

0.75

0.2

0.75

0.3

0.75

0.4

0.75

0.5

0.75

0.6

0.75

0.7

0.5

0.8

0.5

0

0.9

0

1.0

0

Rank	Relevant retrieved	Total retrieved	Precision	Recall
1	0	1	0.0	0.0
2	1	2	0.5	0.2
3	2	3	0.67	0.4
4	3	4	0.75	0.6
5	3	5	0.6	0.6
6	3	6	0.5	0.6
7	3	7	0.43	0.6
8	4	8	0.5	0.8
9	4	9	0.44	0.8
10	4	10	0.4	0.8

\* Interpolated Precision for specified recall

Recall level

Interpolated Precision

$$0.0 : 0.0 \quad \text{Precision} = 0.0$$

$$0.1 : 0.5 \quad \text{Precision} = 0.5$$

$$0.2 : 0.5 \quad \text{Precision} = 0.5$$

$$0.3 : 0.67 \quad \text{Precision} = 0.67$$

$$0.4 : 0.67 \quad \text{Precision} = 0.67$$

$$0.5 : 0.75 \quad \text{Precision} = 0.75$$

$$0.6 : 0.75 \quad \text{Precision} = 0.75$$

$$0.7 : 0.75 \quad \text{Precision} = 0.75$$

$$0.8 : 0.5 \quad \text{Precision} = 0.5$$

$$0.9 \quad \text{Precision} = 0.5$$

$$1.0 \quad \text{Precision} = 0.4$$

Q7] a. Total no. of docs in collection = 10,000

No. of relevant documents = 150

No. of documents retrieved by the system = 250

No. of relevant documents retrieved = 125

$$\text{Precision} = \frac{\text{Relevant documents retrieved}}{\text{Total documents retrieved}} = \frac{125}{500} = 0.5 \text{ or } 50\%$$

$$\text{Recall} = \frac{\text{Relevant documents retrieved}}{\text{Total relevant documents}} = \frac{125}{250} = \frac{5}{6} \approx 83.3\%$$

b) Precision  $\Rightarrow$  50% precision means that half of the documents retrieved by the system are actually relevant. In this scenario, it means that for every 2 patients the system identifies as needing a retest, 1 patient is not actually diabetic or didn't have a mood test on that malfunction day.

Recall  $\Rightarrow$  83.3% recall indicates that the system was able to identify 83.3% of all diabetic patients who need to repeat their test. This means the system missed about 16.7% of the diabetic patients who should be contacted for a retest.

The system's recall is quite high, meaning it retrieves most of the relevant patients. The precision is lower, suggesting that the hospital should contact many patients unnecessarily. In medical scenarios high recall is critical since missing a relevant patient could lead to serious health consequences. However, improving precision would reduce unnecessary communication with patients who don't need retests.

c) If the system is tuned for 100% recall, it would attempt to retrieve all relevant documents, this would result in a significant drop in precision. The system would retrieve many more irrelevant documents, meaning the hospital would contact large numbers of non-diabetic or unaffected patients.

- d) In this scenario, recall is more important than precision. That is because :-
- 1) Missing even a single diabetic patient, who needs a repeat test could lead to health risks. It is better to ensure that as many relevant patients as possible are contacted unnecessarily (low precision)
  - 2) A high recall ensures that most, if not all of the diabetic patients who had a mood test on the malfunction day are identified.

Value of weighing factor  $\alpha$  in F-score which is the weighted harmonic mean of precision and recall are as follows:-

$$F_{\beta} = \frac{(1+\beta)^2 \cdot (\text{Precision} \cdot \text{Recall})}{\beta^2 \cdot (\text{Precision} + \text{Recall})}$$

In this scenario, we would choose higher  $\beta$ . A common value could be  $\beta = 2$  which gives more weight to the recall in f-score.

In terms of weighing factor  $\alpha$ , which assigns a values that prioritizes recall, a reasonable choice would be  $\alpha = 0.2$  or  $\alpha = 0.3$ .

- Q8. a) Precision & recall when both judges must agree for a document to be relevant.

Documents considered relevant by both judges :-

$$\text{Judge 1} : \{3, 4, 5, 6, 7, 8, 12\}$$

$$\text{Judge 2} : \{8, 4, 9, 10, 11\}$$

$$\text{Intersection} : \{3, 4\}$$

$$\text{System returns the set} : \{2, 5, 6, 7, 8\}$$

Relevant documents retrieved (common to the system results & the set  $\{3, 4\}\} = 0$

$$\text{Total retrieved documents} = 5$$

Total relevant documents (Both judges agree) = 2

$$\text{Precision} = \frac{0}{5} = 0, \text{Recall} = \frac{0}{2} = 0$$

b) A document is considered relevant if either judge considers it relevant. So, we take the union of relevant documents.

Documents considered relevant by either judge :-

$$\text{Judge 1's relevant documents} = \{3, 4, 5, 6, 7, 8, 12\}$$

$$\text{Judge 2's relevant documents} = \{3, 4, 9, 10, 11\}$$

Union (Documents either judge considers relevant)

$$= \{3, 4, 5, 6, 7, 8, 9, 10, 11, 12\}$$

Now, system returns the set {2, 5, 6, 7, 8}

Relevant documents retrieved = {5, 6, 7, 8} (4 docs)

Total retrieved documents = 5

Total relevant documents (union of both judges) = 10

$$\text{Precision} = \frac{4}{5} = 0.8$$

$$\text{Recall} = \frac{4}{10} = 0.4$$

c) Average Precision Based on Both Relevance scenarios

i) Average

9a. RNRNNNNRRNNNNNNRNNNR

$$\text{MAP}(\Phi) = \frac{1}{8} \left[ 1 + \frac{2}{3} + \frac{3}{8} + \frac{4}{9} + \frac{5}{15} + \frac{6}{20} \right] \\ = 0.389 \approx 0.39$$

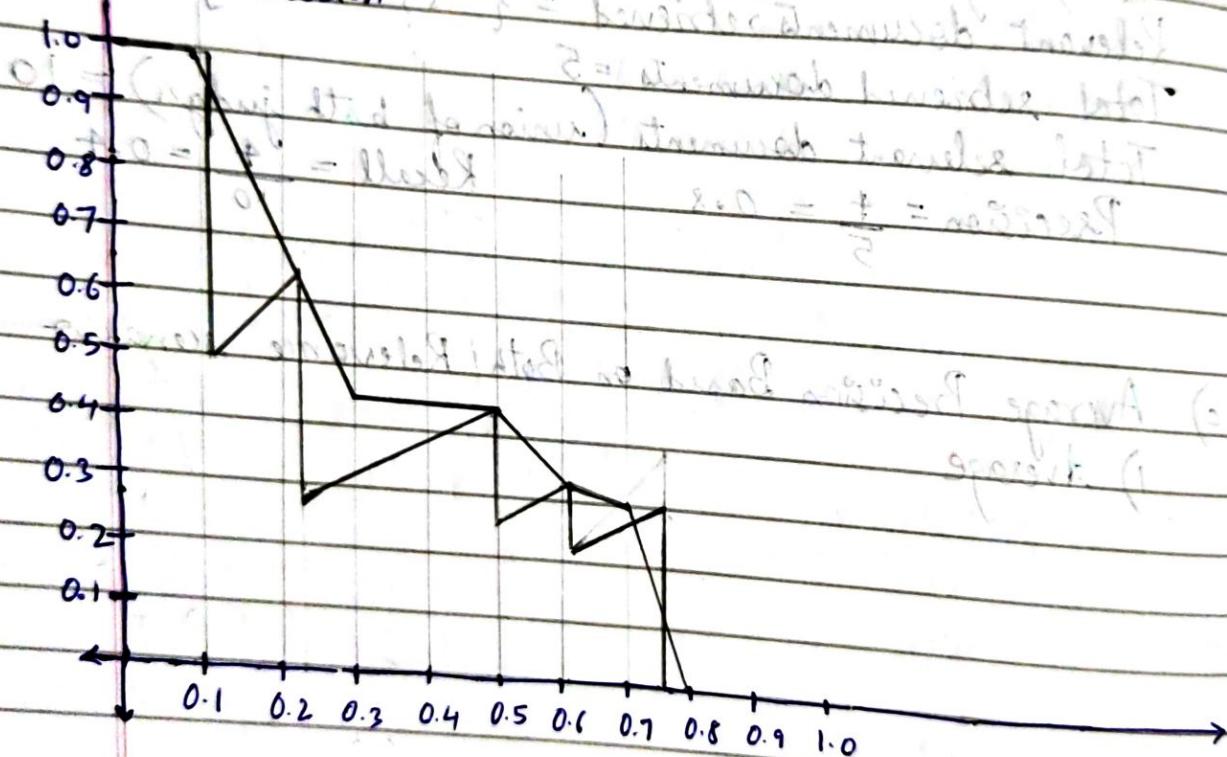
b. For R-precision, we take ratio of relevant docs in the first 8 docs (total relevant docs) to total relevant docs.

R-precision =  $\frac{3}{8} = 0.375$       R NK NNNN B

c) Precision at 10

$$\frac{4}{10} = 0.4$$

## Precision



$$10a. \cos \theta = \frac{\vec{a} \cdot \vec{b}}{|\vec{a}| |\vec{b}|}$$

$$\cos \theta = \frac{x_1y_1 + x_2y_2 + x_3y_3 + x_4y_4}{\sqrt{x_1^2 + x_2^2 + x_3^2 + x_4^2} \sqrt{y_1^2 + y_2^2 + y_3^2 + y_4^2}}$$

b. Report A  $\langle 12, 0, 3, 24 \rangle$

Report B  $\langle 10, 5, 20, 10 \rangle$

Report C  $\langle 0, 12, 9, 8 \rangle$

Query  $\langle 1, 1, 1, 1 \rangle$

$$\cos \theta_{A,q} = \frac{12 + 0 + 3 + 24}{\sqrt{12^2 + 0^2 + 3^2 + 24^2} \sqrt{1^2 + 1^2 + 1^2 + 1^2}} = \frac{13}{27 \times 2} = \frac{13}{54} = 0.722$$

$$\cos \theta_{B,q} = \frac{10 + 5 + 20 + 10}{\sqrt{10^2 + 5^2 + 20^2 + 10^2} \times 2} = \frac{45}{25 \times 2} = \frac{9}{10} = 0.9$$

$$\cos \theta_{C,q} = \frac{12 + 9 + 8}{\sqrt{12^2 + 9^2 + 8^2} \times 2} = \frac{29}{17 \times 2} = \frac{29}{34} = 0.853$$

$$\theta_{B,q} < \theta_{C,q} < \theta_{A,q} [\cos \theta_{B,q} > \cos \theta_{C,q} > \cos \theta_{A,q}]$$

B is most likely to be relevant to the query  
 $\therefore$  We choose B.

Since B is the only document containing all the keywords  
 Report C does not mention North Sea.

Report A does not mention oil at all. Superiority of Report A seems clear from cosine values, but it is not obvious from simply inspecting the numbers.