

which will re...  
who.

①

22<sup>nd</sup> June, 2020.

17K-3730.

E GR-1.

## INFORMATION RETRIEVAL FINAL EXAM - 2020

### Question # 01:

(a)

#### 1. Limitations of boolean retrieval model:

Some limitations of boolean retrieval model are

- Exact matching → boolean retrieval model only allows exact matching and doesn't support the retrieval based on partial matching.
- The results are not based on scoring → all documents have equal importance which is not a very good approach in some cases.
- Hard query formation/translation into boolean expression.

#### 2. Major problems observed in VSM:

Words are ~~not~~ kept in a "bag of word" which means ~~they~~ their order is lost and they are treated in statistically independent manner. There can be a lot of zero elements in the document vectors which can lead to large amount of data sparsity in vectors, which wastes large amount of

REB12.

②

problem of synonymy

17K-3730

memory - The problem of synonymy ambiguity may occur. Longer documents may have higher frequency.

### 3- False negatives in a VSM

False negatives are where the document model incorrectly predicts negative class. In a VSM, sometimes the word substrings/stemmed similar words (i.e.: doll<sup>ed</sup> up or doll become doll after stemming) might result in a false positive match. The search keywords should, therefore, match precisely to doc. terms.

### 4- Drawback of VSM from human standpoint:

From a human perspective, the sentence and documents have a context and the words are connected with each other having some context while, in a VSM terms are treated statistically independent. The problem of synonymy and polysemy may also occur which is also a drawback.

Statistically independent refers that order of the terms is lost.

(3)

(a)

5.

(a)

It's one of the assumptions of PRP that if the word is not present in query it's equally likely to occur in R and NR populations. Assuming this, we calculate the probabilities of the common words of query and docs.

(b)

$$d_1 = \langle 0.04, 0.04, 0, 0.06, 0, 0, 0 \rangle$$

$$d_2 = \langle 0, 0.11, 0.16, 0, 0, 0.07, 0 \rangle$$

$$d_3 = \langle 0.01, 0.07, 0, 0, 0, 0, 0.03 \rangle$$

Query vector :

$$q = \langle 0, 0, 1, 1, 0, 0, 1 \rangle$$

$$\cos(d_1, q) = \frac{d_1 \cdot q}{|d_1| \cdot |q|} = \frac{0.06}{(0.08)(1.73)} = 0.43.$$

$$\cos(d_2, q) = \frac{d_2 \cdot q}{|d_2| \cdot |q|} = \frac{0.16}{(0.21)(1.73)} = 0.44$$

$$\cos(d_3, q) = \frac{d_3 \cdot q}{|d_3| \cdot |q|} = \frac{0.03}{(0.07)(1.73)} = 0.24$$

Document ranking :

D<sub>2</sub>, D<sub>1</sub>, D<sub>3</sub>.

(4)

## Question #02

(a)

Getting The two quantities, precision and recall clearly trade off against one another. If we talk about the IR

system development perspective, there can be two types of applications; One system that is precision sensitive i.e. web search and one that is recall sensitive i.e. legal/patent search. In these type of IR systems having precision or recall greater than each other is tolerable to a certain percentage in order to serve the demand/need. F-measure can be used, which is HM of precision and recall.

(b)

Break even point in IR systems are evaluation measure that returns a list of documents based on their order of relevance score according to the user's requirement/need.

There has to be a break even point between precision and recall. This can be explained as

If ranked Element is NR, then  $tp=0$

which will result in trivial break-even point where precision = recall. If highest relevant element is ranked and Relevant docs are found in list, then fr decreases and false positive doesn't change. Therefore the document list is ordered based on the false positives and false negatives i.e

$$\uparrow fp < fn$$

$$\downarrow fn < fp$$

so, there has to be a break-even point in the list.

(c)

1) Precision =  $\frac{tp}{tp + fp}$

$tp = 5, fp = 7, fn = 8 - 5 = 3$

$$\text{precision} = \frac{5}{7+5} = \frac{5}{12}$$

precision = 0.41667

2) finding recall to get the F<sub>1</sub> score.

$$\text{Recall} = \frac{tp}{fn + fp} = \frac{5}{3+5} = \frac{5}{8}$$

Recall = 0.625

now, finding f<sub>1</sub>:

PFB11.

$$F_1 = \frac{2PR}{P+R} = \frac{2 \times 0.4166 \times 0.625}{0.41667 \times 0.625}$$

F<sub>1</sub> = 0.5

3) MAP =  $\frac{1}{5} \left( \frac{1}{1} + \frac{2}{2} + \frac{3}{5} + \frac{4}{9} + \frac{5}{12} \right) = 0.6922$

MAP = 0.6922

4) Largest possible MAP would be :

$$MAP = \frac{1}{8} \left( \frac{1}{1} + \frac{2}{2} + \frac{3}{5} + \frac{4}{9} + \frac{5}{12} + \frac{6}{13} + \frac{7}{14} + \frac{8}{15} \right)$$

MAP = 0.619

5) Smallest possible MAP :

→ when remaining 3 docs retrieved from collection  
 $(8 - 5 = 3)$ .

$$MAP = \frac{1}{8} \left( \frac{1}{1} + \frac{2}{2} + \frac{3}{5} + \frac{4}{9} + \frac{5}{12} + \frac{6}{998} + \frac{7}{999} + \frac{8}{1000} \right)$$

MAP = 0.435

## QUESTION # 03

|          | D <sub>1</sub> | D <sub>2</sub> | D <sub>3</sub> | D <sub>4</sub> | D <sub>5</sub> | Test |
|----------|----------------|----------------|----------------|----------------|----------------|------|
| Taipei   | 1              | 0              | 0              | 0              | 0              |      |
| Taiwan   | 1              | 1              | 0              | 1              | 2              |      |
| Macao    | 0              | 1              | 0              | 0              | 0              |      |
| Shanghai | 0              | 1              | 0              | 0              | 0              |      |
| Japan    | 0              | 0              | 1              | 0              | 0              |      |
| Sapporo  | 0              | 0              | 1              | 1              | 1              |      |
| Osaka    | 0              | 0              | 0              | 1              | 0              |      |

a)  $\cos(d_1, d_5) = \frac{\vec{d}_1 \cdot \vec{d}_5}{|d_1||d_5|}$

$$= \frac{\langle 1, 1, 0, 0, 0, 0, 0 \rangle \langle 0, 2, 0, 0, 0, 1, 0 \rangle}{\sqrt{2} \cdot \sqrt{5}}$$

$$= \frac{2}{\sqrt{10}} = 0.6324$$

b)

$$= \frac{\langle 0, 1, 1, 1, 0, 0, 0 \rangle \langle 0, 2, 0, 0, 0, 1, 0 \rangle}{\sqrt{3} \cdot \sqrt{5}}$$

$$= \frac{2}{\sqrt{15}} = 0.5164$$

c)

$$= \frac{<0,0,0,0,1,1,0> <0,2,0,0,0,1,0>}{\sqrt{2} \cdot \sqrt{5}}$$

$$= \frac{1}{\sqrt{10}}$$

$$= 0.3162$$

d)

$$= \frac{<0,1,0,0,0,1,1> <0,2,0,0,0,1,0>}{\sqrt{3} \cdot \sqrt{5}}$$

$$= \frac{3}{\sqrt{15}}$$

$$= 0.7746$$

(i) For 3-NN, the query doc ( $D_5$ ) belongs to class "Yes".

(ii) Rocchio's Algorithm:

Centroids :

$$\begin{aligned} a) M_{\text{china}=\text{yes}} &= \frac{1}{D_{\text{china}=\text{yes}}} \sum_{d_j \in D_{\text{china}=\text{yes}}} \vec{v}(d_j) \\ &= \frac{1}{2} (<1,1,0,0,0,0,0> <0,1,1,1,0,0,0>) \\ &= \frac{1}{2} <1,2,1,1,0,0,0> \end{aligned}$$

$$= \langle 0.5, 1, 0.5, 0.5, 0, 0, 0 \rangle$$

b)  $\vec{M}_{\text{china}=\text{no}} = \frac{1}{D_{\text{china}=\text{no}}} \sum_{d_j \in D_{\text{china}=\text{no}}} \vec{v}(d_j)$

$$= \frac{1}{2} (\langle 0, 1, 0, 0, 0, 1, 1 \rangle + \langle 0, 0, 0, 0, 1, 1, 0 \rangle)$$

$$= \frac{\langle 0, 1, 0, 0, 1, 2, 1 \rangle}{2}$$

$$= \langle 0, 0.5, 0, 0, 0.5, 1, 0.5 \rangle$$

for Doc 5 :

$$\|\vec{M}_{\text{yes}} - \vec{d}_5\| = \sqrt{(0.5)^2 + (1-2)^2 + 0.5^2 + 0.5^2 + 0 + 1 + 0}$$

$$= \sqrt{0.25 + 1 + 0.25 + 0.25 + 0 + 1 + 0}$$

$$= 1.658.$$

$$\|\vec{M}_{\text{no}} - \vec{d}_5\| = \sqrt{0 + (0.5-2)^2 + 0 + 0 + 0.5^2 + 0 + 0.5^2}$$

$$= \sqrt{0 + 2.25 + 0 + 0 + 0.25 + 0 + 0.25}$$

$$= 1.6583$$

$$\min(1.65, 1.65)$$

Result can be 'Yes' or 'no'

→ If the result is same the higher prior class is usually considered.

### (iii) Multinomial NB

Priors would be:

$$P(\text{Yes}) = 2/4 = 0.5$$

$$P(\text{No}) = 2/4 = 0.5$$

$$P(w|c) = \frac{\text{count}(w, c) + 1}{\text{count}(c) + |V|}$$

$$\bullet - P(\text{Taiwan}|\text{Yes}) = 2+1/5+7 = 3/12$$

$$P(\text{Sapporo}|\text{Yes}) = 0+1/5+7 = 1/12$$

$$\bullet - P(\text{Taiwan}|\text{No}) = 1+2/12 = 2/12$$

$$P(\text{Sapporo}|\text{No}) = (1+3)/12 = 3/12$$

$$\text{so } P(y|d_5) = 0.5 \times \frac{3}{12} \times \frac{3}{12} \times \frac{1}{12} \\ = 2.60 \times 10^{-3}$$

$$P(\text{no}|d_5) = 0.5 \times \frac{2}{12} \times \frac{2}{12} \times \frac{3}{12}$$

$$= 0.264$$

$d_5$  has the label 'no'

Q4

(a)

let,  $c_1 = d_2$  and  $c_2 = d_5$   
taking the  $c_1$  and  $c_2$  as initial cluster  
so

for  $d_1$  :

$$\text{distance } (c_1, d_1) = \sqrt{(1-2)^2 + (4-4)^2} = 1$$

$$\therefore (c_2, d_1) = \sqrt{(1-2)^2 + (4-1)^2} = 3.16$$

$c_1 < c_2$  so it belongs to  $c_1$

for  $d_3$  :

$$\text{distance } (c_1, d_3) = (4-2)^2 + (4-4)^2 = 2$$

$$\therefore (c_2, d_3) = (4-2)^2 + (4-1)^2 = 3.6$$

$c_1 < c_2$  so it belongs to  $c_1$

for  $d_4$  :

$$d(c_1, d_4) = (1-2)^2 + (1-4)^2 = 3.16$$

$$d(c_2, d_4) = (1-2)^2 + (1-1)^2 = 1$$

$c_2 < c_1$  so belongs to  $c_2$

for  $d_6$  :

$$d(c_1, d_6) = 3.6$$

$$d(c_2, d_6) = 2$$

$c_2 < c_1$

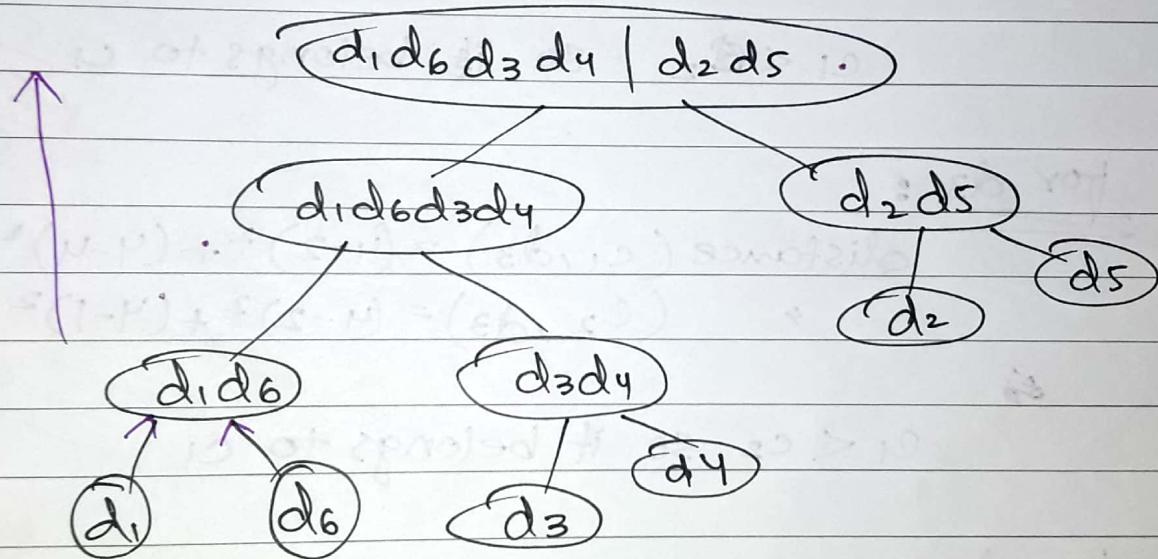
so  $c_2$ .

REB312.

for k-mean the time complexity is  $O(n^2)$   
 which can be improved through various  
 techniques i.e.: choosing the right stopping  
 criteria.

(b)

### HAC approach



|                | d <sub>1</sub> | d <sub>2</sub> | d <sub>3</sub> | d <sub>4</sub> | d <sub>5</sub> | d <sub>6</sub> |     |
|----------------|----------------|----------------|----------------|----------------|----------------|----------------|-----|
| d <sub>1</sub> | 0              | 1              | 3              | 3              | 3.16           | 4.24           | max |
| d <sub>2</sub> | 1              | 0              | 2              | 3.16           | 3              | 3.6            |     |
| d <sub>3</sub> | 3              | 2              | 0              | 4.24           | 3.6            | 3              |     |
| d <sub>4</sub> | 3              | 3.16           | 4.24           | 0              | 1              | 3              |     |
| d <sub>5</sub> | 3.16           | 3              | 3.6            | 1              | 0              | 2              |     |
| d <sub>6</sub> | 4.2            | 3.6            | 3              | 3              | 2              | 0              |     |

|   | $d_1d_6$ | $d_2$ | $d_3$ | $d_4$ | $d_5$ |
|---|----------|-------|-------|-------|-------|
| 1 | 2.12     | 2.3   | 3     | 3     | 2.55  |
| 2 | 2.3      | 0     | 2     | 3.16  | 3     |
| 3 | 3        | 2     | 0     | 4.24  | 3.6   |
| 4 | 3        | 3.16  | 4.24  | 0     | 1     |
| 5 | 2.56     | 3     | 3.6   | 1     | 0     |

max.

|          | $d_1d_6$ | $d_2$ | $d_3d_4$ | $d_5$ |
|----------|----------|-------|----------|-------|
| $d_1d_6$ | 2.1      | 2.3   | 3        | 2.58  |
| $d_2$    | 2.3      | 0     | 2.58     | 3     |
| $d_3d_4$ | 3        | 2.58  | 2.12     | 2.12  |
| $d_5$    | 2.58     | 3     | 2.12     | 0     |

max

|                | $d_1d_6d_3d_4$ | $d_2$ | $d_5$ |
|----------------|----------------|-------|-------|
| $d_1d_6d_3d_4$ | 2.56           | 2.4   | 2.56  |
| $d_2$          | 2.4            | 0     | 3     |
| $d_5$          | 2.5            | 3     | 0     |

max

Complexity :  $n^3$ .

(c) No, we do not expect same result in part a, b - As in ~~mean centroid~~

## Q5

(a)

Type of web queries:

a) Navigational → user search for single entity  
~~or directly go to the page~~

b) Transactional → user go for online payments and transactions.

c) ~~Navigatio~~ Informational → user search for specific topic or general information of specific topic.

(b)

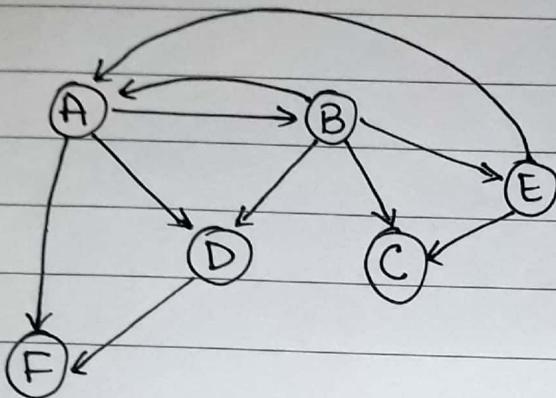
Politeness: No. of searched web pages without permission

Solution: Policies for regulating the crawler

Freshness: The crawler should crawl page with the rate of change of that page.

Extensible: Crawler should be extensible to handle newly fetched protocols, new format etc.

Q6



In the given web graph edge C and F have no outgoing edge. Therefore, it's not always possible to redirect ~~from the all page F and C~~ <sup>from</sup> the pages (due to F and C), also from D there is no way to go to any page other than F.

|   | A             | B             | C | D | E             | F |        |
|---|---------------|---------------|---|---|---------------|---|--------|
| A | 0             | $\frac{1}{4}$ | 0 | 0 | $\frac{1}{2}$ | 0 | $0.85$ |
| B | $\frac{1}{3}$ | 0             | 0 | 0 | 0             | 0 | $0.85$ |
| C | 0             | $\frac{1}{4}$ | 0 | 0 | $\frac{1}{2}$ | 0 | $0.85$ |
| D | $\frac{1}{3}$ | $\frac{1}{4}$ | 0 | 0 | 0             | 0 | $0.85$ |
| E | 0             | $\frac{1}{4}$ | 0 | 0 | 0             | 0 | $0.85$ |
| F | $\frac{1}{3}$ | 0             | 0 | 1 | 0             | 0 | $0.85$ |