

Name - Ayesha Ansab

Roll no - 17K-3867

Section - GR2

Course - Information Retrieval

①

Date: 11/09/2020

## QUESTION 1

### Part (a)

- ① Limitations are: ① It doesn't allow more flexible matching operation. The exact matching either retrieves too many or less documents. ② All terms are equally weighted which gives flat results & doesn't allow ranked retrieval. ③ It is difficult to formulate a query & process large data.
- ② We assume theoretically assume that the terms are statistically independent which is not always true in practice. The order in which the term appears in document is lost in NBM representation.
- ③ Other important factors are: ① The keywords that are searched should precisely match the document term, else the substring words can result in false negative match. ② The synonyms (i.e. semantically sensitive words) that have same context (meaning) but different spelling (terms) might also result in false negative as they don't match the vocabulary.

- ③ A few important factors are: ① The keywords that are searched should precisely match the document term, else the substring words can result in false negative match. ② The synonyms (i.e. semantically sensitive words that have same context (meaning) but different spelling (terms)) might also result in false negative as they don't match the vocabulary.
- ④ In VSM, the features are spaces. It can not deal with lexical (vocabulary) ambiguity therefore if we want to work on domain specific words, similarity would be close to null resulting in meaningless data.
- ⑤ Since/we assume that terms are independent of each other & if we do not need to estimate as we are only ranking documents. Naive Bayes independence assumption makes it easier that absence & presence of one word is independent to another. & since each term is either present or not therefore & we assume & it is likely to happen that there's an equal chance.

## QUESTION 1 (part b)

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

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

$$d_2 = \langle 0, 0.11, 0.16, 0, 0, 0.09, 0 \rangle$$

$$= 0.43$$

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

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

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

$$= 0.21 \times 1.73$$

$$\cos(d_3, q) = \frac{0.03}{0.09 \times 1.73} = 0.24$$

Time complexity =  $O(n^2)$

rank:  $d_2, d_1, d_3$ .

Page No.

Bright

Teacher's Signature \_\_\_\_\_

Name - Ayesha Asraf

Roll no - 17K-3867

Section - GR2

course - Information Retrieval

Date \_\_\_\_\_

## QUESTION 2

### Part (a)

An IR system can give 100% precision when a single related doc is retrieved & it can also give 100% recall when all docs are returned for some query. Whereas, when used ~~together~~ together, they can alter the values.

### Part (b)

relationship

Breakeven point is used to define equality, holding. It is evaluate retrieval system that returns a list of docs ordered by their supposed relevance need.

Precision = recall iff  $fp = fn$  or  $tp = 0$ . If the highest ranked element is not relevant then  $tp = 0$  and that is trivial breakeven point.

If highest ranked element is relevant then the no. of false +ve increases as you go down the list & the no. of false -ve decreases. As we go down the list, if item is R, then  $fn$  decrease &  $fp$  does not change. If item is N, then  $fp$  increase

&  $fn$  does not change. At the beginning of list  $fp < fn$  & at the end  $fp > fn$ ; Thus, there has to be a break even point

It is related to F1 because it can be used to evaluate classification model

Name - Ayesha Asraf

Roll no - 17k-3867

Section - GR2

Course - Information Retrieval

Date - 20 \_\_\_\_\_

## QUESTION 2

Part (c)

$$\textcircled{a} \text{ Precision} = 5/12 \approx 0.416$$

$$\textcircled{b} \text{ Recall} = 5/8 \approx 0.625$$

$$F_1 = \frac{2PR}{P+R} \approx 0.499$$

$$\textcircled{c} \text{ MAP} = \frac{1}{8} \times (1 + 2/2 + 3/5 + 4/9 + 5/12) \\ \approx 0.692$$

$$\textcircled{d} \text{ MAP}_{\text{largest}} = \frac{1}{8} \times (1 + 1 + 3/5 + 4/9 + 5/12 + 6/13 + 7/14 + 8/15) \\ = 0.61$$

$$\textcircled{e} \text{ MAP}_{\text{smallest}} = \frac{1}{8} \times (1 + 1 + 3/5 + 4/9 + 5/12 + \underline{6/998} + 7/999 + 8/1000) \\ = 0.43$$

Name - Hyesha Ausab

Roll No - 17K-3867

Section - GR2

Course - Information Retrieval

20

### QUESTION 3

Part

(i) KNN — k=3.

Dictionary = {Japan, Macao, Osaka, Sapporo, Shanghai, Taipei, Taiwan}

$$d_1 = \langle 0, 0, 0, 0, 0, 1, 1 \rangle$$

$$d_2 = \langle 0, 1, 0, 0, 1, 0, 1 \rangle$$

$$d_3 = \langle 1, 0, 0, 1, 0, 0, 0 \rangle$$

$$d_4 = \langle 0, 0, 1, 1, 0, 0, 1 \rangle$$

$$d_5 = \langle 0, 0, 0, 1, 0, 0, 2 \rangle$$

$$\begin{aligned} \cos(d_1, d_5) &= (d_1 \cdot d_5) \div |d_1| \times |d_5| \\ &= (0+0+0+0+0+0+2) \div \sqrt{|1+1|} \times \sqrt{|1+4|} \\ &= 2 / \sqrt{2} * \sqrt{5} \\ &\approx 0.547 \quad 0.632 \end{aligned}$$

$$\begin{aligned} \cos(d_2, d_5) &= (d_2 \cdot d_5) \div |d_2| \times |d_5| \\ &= 2 \div (\sqrt{3} * \sqrt{5}) \\ &\approx 0.516 \end{aligned}$$

$$\begin{aligned} \cos(d_3, d_5) &= 1 \div \sqrt{2} \sqrt{5} \\ &\approx 0.316 \end{aligned}$$

$$\begin{aligned} \cos(d_4, d_5) &= (1+2) \div \sqrt{3} \sqrt{5} \\ &\approx 0.77 \end{aligned}$$

for 3-NN : 0.77, 0.632, 0.516

since  $d_1, d_2$  belong to yes ∴

$d_5$  belongs to yes

Name — Ayesha Ausaf

Roll no — 17K-3867

Section — GR2

Course — Information Retrieval

— 20 —

part (iii)

$$\hat{P}(c) = \frac{N_c}{N} ; \hat{P}(w|c) = \frac{\text{count}(w, c) + 1}{\text{count}(c) + 1 \vee 1}$$

① Prior Probability

$$P(\text{yes}) = \frac{2}{4} \\ = \frac{1}{2}$$

$$P(\text{no}) = \frac{2}{4} \\ = \frac{1}{2}$$

② Conditional probability

$$P(\text{Japan}|\text{yes}) = 0+1/5+7 = \frac{1}{12}$$

$$P(\text{Macao}|\text{yes}) = 1+1/5+7 = \frac{1}{6}$$

$$P(\text{Osaka}|\text{yes}) = 0+1/5+7 = \frac{1}{12}$$

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

$$P(\text{Shanghai}|\text{yes}) = 1+1/5+7 = \frac{1}{6}$$

$$P(\text{Taipei}|\text{yes}) = 1+1/5+7 = \frac{1}{6}$$

$$P(\text{Taiwan}|\text{yes}) = 2+1/5+7 = \frac{1}{4}$$

$$P(\text{Japan}|\text{no}) = \frac{0+1/5+7}{1} = \frac{1}{6}$$

$$P(\text{Macao}|\text{no}) = \frac{0+1/5+7}{1} = \frac{1}{12}$$

$$P(\text{Osaka}|\text{no}) = \frac{1+1/5+7}{1} = \frac{1}{6}$$

$$P(\text{Sapporo}|\text{no}) = \frac{2+1/5+7}{1} = \frac{1}{4}$$

$$P(\text{Shanghai}|\text{no}) = \frac{0+1/5+7}{1} = \frac{1}{12}$$

$$P(\text{Taipei}|\text{no}) = \frac{0+1/5+7}{1} = \frac{1}{12}$$

$$P(\text{Taiwan}|\text{no}) = \frac{1+1/5+7}{1} = \frac{1}{6}$$

③ Choosing class

$$P(\text{yes}|\text{ds}) = \frac{1}{4} \times \frac{1}{4} \times \frac{1}{12} \times \frac{1}{2} \\ = 2.6 \times 10^{-3}$$

$$P(\text{no}|\text{ds}) = \frac{1}{6} \times \frac{1}{6} \times \frac{1}{4} \times \frac{1}{2} \\ = 3.47 \times 10^{-3}$$

Since  $3.47 \times 10^{-3}$  is greater, this doc5 belongs to class  
no

QUESTION 4 (part a)

let  $c_1 = d_2 \quad \& \quad c_2 = d_5$ .

$c_1 \& c_2$  are initial clusters. Deciding to which cluster each doc belongs to.

For  $d_1$ :

$$\text{Distance}(c_1, d_1) < \text{Distance}(c_2, d_1)$$

$\therefore d_1$  belongs to  $c_1$

For  $d_3$ :

$$\text{dist}(c_1, d_3) < \text{dist}(c_2, d_3)$$

$\therefore d_3$  belongs to  $c_1$

For  $d_4$ :  $\text{dist}(c_1, d_4) > \text{dist}(c_2, d_4)$

$\therefore d_4$  belongs to  $c_2$

For  $d_6$ :  $\text{dist}(c_1, d_6) > \text{dist}(c_2, d_6)$

$\therefore d_6$  belongs to  $c_2$

$\therefore d_1, d_2$  and  $d_3$  are in  $c_1$ .  $d_4, d_5, d_6$  are in  $c_2$ .

K mean converge to local maximum, in order to find optimal clustering one need to produce all possible clustering arrangements

Name — Ayesha Ansab  
Roll no — 17K-3867  
Section — GR2  
Course — Information Retrieval

Date \_\_\_\_\_ 20 \_\_\_\_\_

### QUESTION 5 (part a)

There are ~~4~~ type of queries

- Informational query
- Navigational query
- Transactional query
- Faceted query (extra)

Navigational → some user might put youtube or instagram in search bar rather than adding a URL.

Informational → some user searching about IR or AI to learn about the domains of CS.

Transactional → some user might search to download some app, or know price of some iPhone X.

Faceted → A query on drugs for prevention of osteoporosis would include

- osteoporosis OR bone loss
- prevention OR cure

## QUESTION 5 (b)

### (i) Politeness

Web servers have both implicit & explicit <sup>polices</sup> regulating the rate at which a crawler can visit them. These politeness policies must be respected.

Sol: ~~Yield/ St/ specify~~

challenge: URL frontier - containing URL yet to be fetched in current crawl

solution: decrease speed by which it visits sites.

### (ii) Freshness

In many app, crawler should operate in continuous mode. It should obtain fresh copies of previously fetched pages.

challenge: Refreshing page & bringing new page is a challenge.

Sol: don't overload server & use network to maximum extent.

### (iii) Extensible

Crawler should be designed to be extensible in many ways - to cope with new data format, new fetch protocol etc.

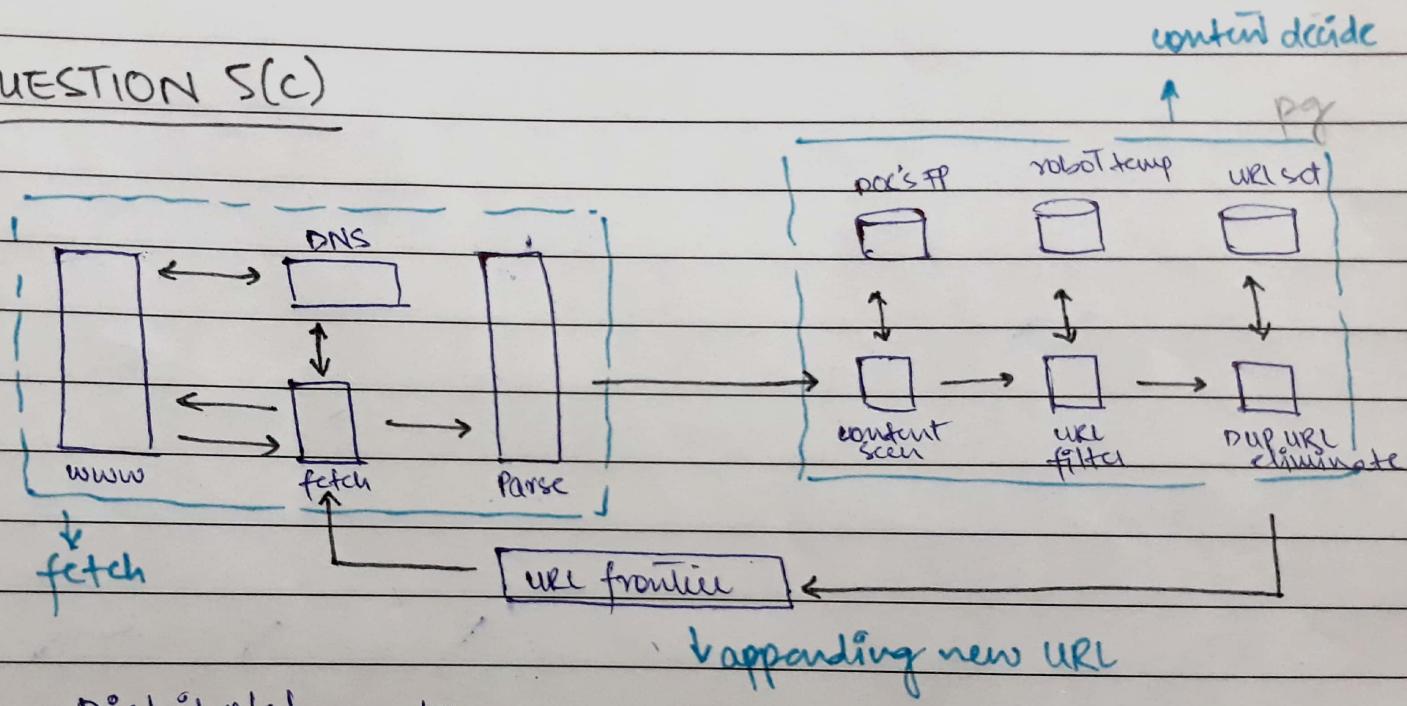
Crawler architect should be modular.

challenge: Data with unknown format

Sol: Decode the format to extract data.

content decide

### QUESTION 5(c)



→ Distributed crawler

Threads in a crawler can run under different process each at a node of a distributed crawling sys. It is essential for scaling & can also be good for geographically know where each node crawler hosts 'near it'. Partitioning can be done by hash function.

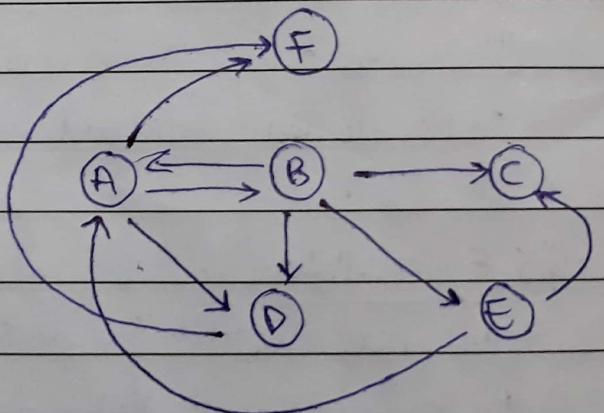
→ ~~Distributed crawler for han~~

Name — Ayesha Ausab  
Roll no — 17K-3867  
Section —  
Course —

Date \_\_\_\_\_

### QUESTION 6 (part a)

	A	B	C	D	E	F
A	0	1	0	1	0	1
B	1	0	1	1	1	0
C	0	0	0	0	0	0
D	0	0	0	0	0	1
E	1	0	1	0	0	0
F	0	0	0	0	0	0



Name - Ayesha Asraf

Roll no - 17K-3867

Section -

Course -

Date \_\_\_\_\_

## QUESTION 6 (part b)

$$a^* / a^T = A^T \cdot h^{(0)}$$

$$\left[ \begin{array}{cccccc|ccc} 0 & 1 & 0 & 1 & 0 & 1 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 & 1 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ 1 & 0 & 1 & 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 \end{array} \right]$$

$$\left[ \begin{array}{cccccc|c} 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 1 & 0 & 1 \\ 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 & 1 \end{array} \right]$$

$$\left[ \begin{array}{c|c} 2 & \\ 1 & \\ 2 & \\ 2 & \\ 1 & \\ 2 & \end{array} \right]$$

$$= \sqrt{2^2 + 1^2 + 2^2 + 2^2 + 1^2 + 2^2}$$

$$= \sqrt{4 + 1 + 4 + 4 + 1 + 4}$$
$$= \sqrt{18}$$

$$= \left[ \begin{array}{c} 0.471 \\ 0.236 \\ 0.471 \\ 0.471 \\ 0.236 \\ 0.471 \end{array} \right]$$

Name - Ayushma Arora  
Roll no - 17E-3267  
Section -  
Course -

Date \_\_\_\_\_  
20

$$u^{(0)} = A \cdot a^0$$

$$\begin{bmatrix} 0 & 1 & 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0.571 \\ 0.236 \\ 0.471 \\ 0.471 \\ 0.236 \\ 0.471 \end{bmatrix}$$

$$\begin{bmatrix} 1.0178 \\ 1.0169 \\ 0 \\ 0.471 + 0.942 \\ 0 \end{bmatrix}$$

$$a^2 = \begin{bmatrix} 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0.31 & 0 & 0 & 1 & 0 & 0 \\ 1.4 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

by

$$u^{(1)} \Rightarrow \begin{bmatrix} 1.025 \\ 1.08 \\ 0 \\ 0.31 \\ 0.31 \\ 0.94 \end{bmatrix}$$

$$A = \begin{bmatrix} 0.47 \\ 0.23 \\ 0.47 \\ 0.45 \\ 0.24 \\ 0.31 \end{bmatrix}$$