

First Sem 20-21

Comprehensive Exam Solution

CS F469 Information Retrieval

Q1. A)  $q_{23} = \frac{q_{21}}{2} + \frac{q_{22}}{2}$

B)

$$\begin{matrix} & 1 & 2 & 3 & 4 \\ 1 & 0 & \frac{1}{2} & 0 & 0 \\ 2 & \frac{1}{2} & 0 & 0 & 0 \\ 3 & \frac{1}{2} & \frac{1}{2} & 0 & 0 \\ 4 & 0 & 0 & 1 & 0 \end{matrix}$$

C) The matrix generated is not suitable for pagerank since the node 4 was a dead end and there is no way a user can exit once landed on this page. Hence the matrix is not column stochastic.

D)

$$0.8 \begin{bmatrix} 0 & \frac{1}{2} & 0 & 0 \\ \frac{1}{2} & 0 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} + 0.2 \begin{bmatrix} \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \\ \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \\ \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \\ \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \end{bmatrix}$$

$$\begin{bmatrix} 0.05 & 0.45 & 0.05 & 0.05 \\ 0.45 & 0.05 & 0.05 & 0.05 \\ 0.45 & 0.45 & 0.05 & 0.05 \\ 0.05 & 0.05 & 0.85 & 0.05 \end{bmatrix}$$

$$B) \quad \pi_{t+1} = M \pi_t$$

$$\begin{bmatrix} \pi_1 \\ \pi_2 \\ \pi_3 \\ \pi_4 \end{bmatrix} = \begin{bmatrix} 0.05 & 0.45 & 0.05 & 0.05 \\ 0.45 & 0.05 & 0.05 & 0.05 \\ 0.45 & 0.45 & 0.05 & 0.05 \\ 0.05 & 0.05 & 0.85 & 0.05 \end{bmatrix} \begin{bmatrix} \frac{1}{4} \\ \frac{1}{4} \\ \frac{1}{4} \\ \frac{1}{4} \end{bmatrix}$$

$$\begin{bmatrix} \pi_1 \\ \pi_2 \\ \pi_3 \\ \pi_4 \end{bmatrix} = \begin{bmatrix} \frac{1}{4} \\ \frac{1}{4} \\ \frac{1}{4} \\ \frac{1}{4} \end{bmatrix} \begin{bmatrix} 0.15 \\ 0.15 \\ 0.25 \\ 0.25 \end{bmatrix}$$

F) Street graph construction - 3M

Teleportation / dead ends etc - 2M

There is no single solution hence the evaluation  
is based on the above two Explanations.

Q2.

a)  $c_d(:, i) = A(:, i) / \sqrt{C_P(j)}$

$$c_3 = \frac{1}{\sqrt{76/171}} [1 \ 0 \ 5 \ 5 \ 5]^T$$

$$= [1.05 \ 0 \ 7.6 \ 7.6 \ 7.6]^T$$

$\left. \right\} 2M$

$$R_4 = \frac{1}{\sqrt{51/171}} [0 \ 1 \ 5 \ 5]^T$$

$$= [0 \ 1.85 \ 9.3 \ 9.3]$$

$\left. \right\} 2M$

b) Intersection of  $c_d$  &  $r$  is 5.  $\rightarrow 1M$

c) If  $v = \frac{1}{25}$  rating is computed using  $r = C \cup R$

$$\begin{bmatrix} 1.5 \\ 0 \\ 7.6 \\ 7.6 \\ 7.6 \end{bmatrix} \left[ \frac{1}{25} \right] \begin{bmatrix} 0 & 1.85 & 9.3 & 9.3 \end{bmatrix}$$

$\left. \right\} 2M$

Rating for Ria for the movie Newton is 2.8

(3)

D) i) since the algorithms use different approaches for predictions they are bound to give different results and considering the strength & weakness of each of the algorithms we can take a linear combination of the all predictions using the following approaches.

(4M)

1. Take the avg prediction of all algorithms.
2. learn weights for each algorithm using a log linear model & then compute a weighted avg of all predictions.
3. Apply an Ensemble methodology.

ii) The RMSE of the final rating using any of the above shall be lower than a single algorithm.

(2M)

E) Shyam's model with 5 latent factors would be superior since the no of latent factors are higher.

(3M)

(Q3)

A)  $9 - 4 - 1 = 4$

The reordering distance for the phrase "for hunting"  
is 4.

- B) Since the order of the long clauses are different if the phrases do not move the sentences generated may not be semantically correct hence movement of the phrases should not be penalised. This can be controlled by a factor  $\alpha$  which is generally in a range of 0-1. Since the reordering is modeled as  $\alpha^{|d|}$  an  $\alpha$ -value close to 1 will not penalise the movement.
- C) Phrase based SMT models are better than word based models since it reduces lexical ambiguities. We can also control the movement of the phrases using the parameter  $\alpha$ . Since we translate phrases the words are seen together in both languages which is an advantage for shorter sentences since we are learning the whole sentence translation in a few cases.

Q6. A) Since the odds ratio is a monotonic function which behaves in the same way as computing the probability it is used to derive the ranking function. BM

B) Since Step considers that there are  $M$  terms in a document - the products runs from 1 to  $M$ . Each term may be present or absent in the document and if the document does not contain a term this parameter computation can be avoided - hence we split it into two factors. BM

$$O(R|B') = \prod_{t: x_t=1} \frac{P(x_t=1 | R=1, \vec{q})}{P(x_t=1 | R=0, \vec{q})},$$

$$\prod_{t: x_t=0} \frac{P(x_t=0 | R=1, \vec{q})}{P(x_t=0 | R=0, \vec{q})}$$

c)

$$\delta(R|\vec{x}, \vec{q}) = \prod_{\substack{t: x_{td}=1 \\ t: x_{tq}=1}} \frac{P(x_t=1|R=1, \vec{q})}{P(x_t=1|R=0, \vec{q})} \cdot \prod_{\substack{t: x_{td}=0 \\ t: x_{tq}=1}} \frac{P(x_t=1|R=1, \vec{q})}{P(x_t=1|R=0, \vec{q})} \cdot$$

$$\prod_{\substack{t: x_{td}=1 \\ t: x_{tq}=0}} \frac{P(x_t=1|R=1, \vec{q})}{P(x_t=1|R=0, \vec{q})} \cdot \prod_{\substack{t: x_{td}=0 \\ t: x_{tq}=0}} \frac{P(x_t=1|R=1, \vec{q})}{P(x_t=1|R=0, \vec{q})}$$

Basically the idea is to have all combinations of the presence/absence of the words in the document & the query. From those four combinations we know that if the word is not present in the query & the doc and also not present in the query & present in the doc they need not be estimated hence we only consider the combinations where by default  $t: x_{tq}=1$  and find their presence absence ( $t: x_{td}=1$  or  $t: x_{td}=0$ ) in the document.

- D) Use the BM25 formula to add the scores for "health" and  
 i) "Insurance" terms:

$$RSV(D_1) = \log_{10}\left(\frac{2000}{20}\right) \frac{(0.8+1) \times 6}{0.8((1-0.7)+0.7 \times \frac{52}{156})+6} +$$

$$\log_{10}\left(\frac{2000}{200}\right) \frac{(0.8+1) \times 6}{0.8((1-0.7)+0.7 \times \frac{52}{156})+6}$$

$$RSV_1(D_1) = 5.04$$

$$RSV(D_2) = \log_{10}\left(\frac{2000}{20}\right) \frac{(0.8+1) \times 8}{0.8((1-0.7)+0.7 \times \frac{38}{156})+8} +$$

$$\log_{10}\left(\frac{2000}{200}\right) \frac{(0.8+1) \times 4}{0.8((1-0.7)+0.7 \times \frac{38}{156})+4}$$

$$RSV_1(D_2) = 4.97$$

Since  $RSV(D_1) > RSV(D_2)$ ,  $D_1$  is considered more relevant to the query.

- i) When  $b=0$  the document length is not taken into account and as  $b$  increases the relative document length matters more. The higher the document length, the term frequency is discounted.
- ii) As  $df$  goes higher, the effect of  $tf$  is cancelled and hence RSV approximates to Scallded version of IDF.