

# COMP6714 (16S2) ASSIGNMENT 2

DUE ON 23:59 31 OCT, 2016 (MON)

Q1 (25 marks)

Consider using the maxscore algorithm to find top-2 results for a query with three different terms  $\{A, B, C\}$ . The scoring function is the BM25 function with  $k_1 = k_3 = 2.0$  and  $b = 0$ .

$$\text{score}(d, Q) = \sum_{t \in Q} \text{idf}_t \cdot \frac{(k_1 + 1) \text{tf}_{t,d}}{k_1((1-b) + b \frac{\text{Ld}_d}{\text{Ld}_{\text{ave}}}) + \text{tf}_{t,d}} \cdot \frac{(k_3 + 1) \text{tf}_{t,Q}}{k_3 + \text{tf}_{t,Q}}$$

Answer the following questions. You need to show major steps. The posting lists are shown below. Each posting consists of document ID and tf.

term	idf	postings
A	6	(D <sub>1</sub> :1), (D <sub>2</sub> :8), (D <sub>5</sub> :3), (D <sub>8</sub> :10)
B	2	(D <sub>1</sub> :1), (D <sub>5</sub> :4), (D <sub>6</sub> :1), (D <sub>7</sub> :4)
C	1	(D <sub>1</sub> :1), (D <sub>2</sub> :2), (D <sub>4</sub> :1), (D <sub>5</sub> :2), (D <sub>6</sub> :3), (D <sub>8</sub> :1), (D <sub>9</sub> :1), (D <sub>10</sub> :3), (D <sub>11</sub> :7)

TABLE 1. Posting Lists

- Show that the maxscore for each keyword can be computed without examining the postings list.
- Using the maxscore obtained above, determine the postings that are accessed for scoring by the algorithm. You need to assume that each skipTo(x) call "magically" moves the cursor directly to the first posting with document ID at least x (i.e., it does not access any other postings).

Hint 1. Calculate the maxscore if you know that the maximum tf is 1, 10, 100, and 1000, respectively.

Q2. (25 marks)

The cluster pruning method is introduced in Chap 7.1.6 of [MRS08].

- Consider the basic method (i.e., only using the closest leader to the query q). Justify the choice of choosing  $\sqrt{N}$  leaders in the preprocessing step. (Hint: try to design a simple model to estimate the query processing cost)

a. N个词, 找  $\sqrt{N}$  leader

b<sub>1</sub> = follower 的 最近 leader

b<sub>2</sub> = 最近 leader.

$$\text{cost} = \left( \frac{N}{\sqrt{N}} \right) + \left( \frac{N}{\sqrt{N}} - 1 \right) \cdot \text{follower} \cdot \sqrt{N} + b_2 \left[ \frac{b_1 N}{\sqrt{N}} - 1 \right]$$

= ...  $\therefore \sqrt{N}$  时, cost 最小



$(0, 1, 0, 0)$   
 $(0, 0, 1, 0)$   
 $(1, 0, 0, 0)$   
 $(0, 0, 0, 1)$

- (2) Find a minimal example where the basic method fails to return the closest document vector to the query  $q$ . You only need to give the document vectors and the query vector, and list the document returned by the cluster pruning method and the correct answer. Will the variation of the basic method (i.e.,  $b_1, b_2 > 1$ ) eliminate such problem (and guaranteed to return the correct answer)?
- (3) Can you propose some modification to this method such that it guarantees returning the closest vector for any query? Describe your method and illustrate it with a small example.

## Q3. (25 marks)

The following list of Rs and Ns represents relevant (R) and nonrelevant (N) returned documents in a ranked list of 20 documents retrieved in response to a query from a collection of 10,000 documents. The top of the ranked list is on the left of the list. This list shows 6 relevant documents. Assume that there are 8 relevant documents in total in the collection.

R R N N N N N N N R N R N N N R N N N N R

(Note that spaces above are just added to make the list easier to read)

- What is the precision of the system on the top-20?
- What is the  $F_1$  on the top-20?
- What is/are the uninterpolated precision(s) of the system at 25% recall?
- What is the interpolated precision at 33% recall?
- Assume that these 20 documents are the complete result set of the system. What is the MAP for the query?

Assume, now, instead, that the system returned the entire 10,000 documents in a ranked list, and these are the first 20 results returned.

- What is the largest possible MAP that this system could have?
- What is the smallest possible MAP that this system could have?
- In a set of experiments, only the top-20 results are evaluated by hand. The result in (5) is used to approximate the range (6) to (7). For this example, how large (in absolute terms) can the error for the MAP be by calculating (5) instead of (6) and (7) for this query?

## Q4. (25 marks)

Consider the documents below.

docID	document text
$D_1$	I don't want to go A groovy king of love You can't hurry love This must be love Take me with you
$D_2$	All out of love Here i am I remember love Love is all Don't tell me

- build a unigram query likelihood language model (LM) for each document. Assume that (i) the only preprocessing done before tokenization is to transform all letters to lower cases, and (ii) we use the Jelinek-Mercer smoothing method with  $\lambda = 0.5$ .
- show which document will be ranked first for the queries:

$$\frac{1}{44} + \frac{1}{76}$$

$$\frac{76+44}{44 \times 76}$$

$$\frac{19+11}{11 \times 76}$$

$$\frac{1}{44} + \frac{1}{76}$$

$$\frac{76+44}{44 \times 76}$$

$$\frac{11+19}{11 \times 76}$$

- $Q_1$ : i remember you
  - $Q_2$ : don't want you to love me
- (3) assume that we have a prior probability distribution over the two documents as  $p(D_1) = 0.7$  and  $p(D_2) = 0.3$ . Will this change the ranking results of the two previous queries?

## SUBMISSION INSTRUCTIONS

You need to write your solutions to the questions in a pdf file named `ass2.pdf`. You must

- include your **name and student ID** in the file, and
- the file can be opened correctly on CSE machines.

You need to show the key steps to get the full mark.

**Note:** Collaboration is allowed. However, each person must independently write up his/her own solution.

You can then submit the file by give `cs6714 ass2 ass2.pdf`.

**Late Penalty:** -10% for the first two days, and -30% for the following days.

$$P(Q_1/d_1) = \prod_{t \in Q_1} P(t|d) = \left(\frac{52}{836}\right) \times \frac{11}{836} \times \frac{62}{836}$$

$$P(Q_1/d_2) = \prod_{t \in Q_1} P(t|d_2) = \frac{62}{608} \times \frac{27}{608} \times \frac{16}{608}$$

$$\underline{P(d_1|Q_1)} = P(d_1|Q_1) \cdot P(Q_1) = \underline{P(Q_1/d_1) \cdot p(d_1) \cdot P(Q_1)}$$

$$\frac{1}{2} \times \frac{1}{22} + \frac{1}{2} \times \frac{1}{38} = \frac{\cancel{19} \times 38 + \cancel{11} \times 22}{22 \times 38} = \frac{19 + 11}{22 \times 38}$$

## COMP6714 (16S2) ASSIGNMENT 2 SAMPLE SOLUTION

### Q1. (25 marks)

- (1) The BM25 formula essentially limits the impact of  $tf$ s (the value converges when  $tf \rightarrow \infty$ ). In our case, the scoring function is

$$score(d) \leq 6f(tf_1) + 2f(tf_2) + f(tf_3)$$

where  $f(x) = \frac{3x}{2+x}$ . Since  $\lim_{x \rightarrow \infty} f(x) = 3$ , we can find the maxscores for the terms are 18, 6, and 3.

- (2) We first consider  $D_1$ , with score

$$score(D_1) = 6f(1) + 2f(1) + f(1) = 9$$

Then we consider  $D_2$

$$score(D_2) = 6f(8) + 2f(0) + f(2) = 15.90$$

At this stage, both of them become the current top-2 results, and  $\tau' = 9$ . Since  $3 + 6 \leq \tau'$ , we only need to consider  $A$ . (hence no need to score  $D_4$ )

Driven by  $A$ , the next document to score is  $D_5$ . We need to probe the lists of  $B$  and  $C$  for  $D_5$ , and compute its score as

$$score(D_5) = 6f(3) + 2f(4) + f(2) = 16.30$$

Similarly, since now  $\tau' = 15.90$ ,

The next document to consider is  $D_8$

$$score(D_8) = 6f(10) + 2f(0) + f(1) = 16.00$$

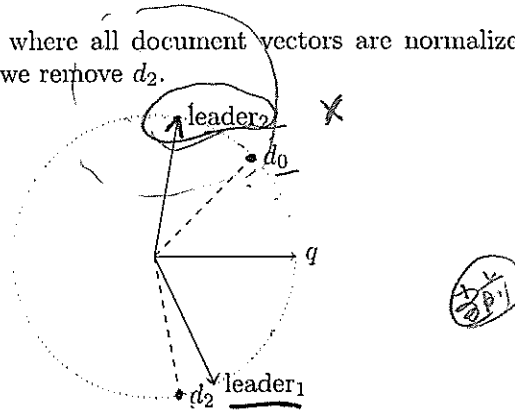
Since  $A$ 's postings list is now exhausted, we conclude that the final top-2 documents are  $D_5$  and  $D_8$ . The algorithm scored 4 documents, and accessed 10 postings.

### Q2. (25 marks)

The *cluster pruning* method is introduced in Chap 7.1.6 of [MRS08].

- (1) Let the number of leaders be  $x$ . Each leader has  $\frac{b_1 N}{x} - 1$  followers on average. During the query processing, we use linear scan to find the  $b_2$  closest leader and then find at most  $b_2(\frac{b_1 N}{x})$  candidates. (We ignore the cost of finding the top- $k$  results from these candidates) The total query processing cost (in terms of distance calculation) is  $f(x) = x + b_2(\frac{b_1 N}{x})$ .  $\frac{d}{dx}f(x) = 1 - \frac{b_2 b_1 N}{x^2}$ . Hence when  $x = \sqrt{b_2 b_1 N}$ , the overall query processing cost is minimized. In the basic model,  $b_2 = b_1 = 1$ , hence  $x = \sqrt{N}$ .

- (2) See the following example where all document vectors are normalized to a unit vector. It is also correct if we remove  $d_2$ .



The query  $q$  is closer to  $leader_1$  than  $leader_2$ . But the correct answer is  $d_0$  which is a follower of  $leader_2$ . Even with  $b_1, b_2 > 1$ , we can find counter-examples (omitted).

- (3) Necessary modifications:

- For each cluster  $c_i$ , calculate the maximum angle between any of the followers and the leader (denoted as  $\theta_i$ ).
- Assume  $k = 1$ . In the query processing, we first calculate the angles between all the leaders and the query. We iterate through the clusters identified by the leaders in increasing order of the angle. When visiting a cluster, we explore all its members by calculating the cosine distance to the query (essentially the angle). The stopping criteria is that the current best result has a smaller angle than  $\alpha_{next}$ , where  $\alpha_{next} = \text{angle}(c_i, q) - \theta_i$  is a lower bound of the angles between a document in  $c_i$  and the query  $q$ . This method can be easily extended to deal with top- $k$  queries.

(The performance of the method might be heavily affected by how well documents form clusters)

Note this is just one of the correct modification methods.

### Q3. (25 marks)

$k$	1	2	3	4	5	6	7	8	9	10
precision (%)	100.00	100.00	66.67	50.00	40.00	33.33	28.57	25.00	33.33	30.00
recall (%)	12.50	25.00	25.00	25.00	25.00	25.00	25.00	25.00	37.50	37.50
$k$	11	12	13	14	15	16	17	18	19	20
precision (%)	36.36	33.33	30.77	28.57	33.33	31.25	29.41	27.78	26.32	30.00
recall (%)	50.00	50.00	50.00	50.00	62.50	62.50	62.50	62.50	62.50	75.00

(1) precision@20 is  $\frac{6}{20}$ .

(2) recall@20 is  $\frac{6}{8}$ .  $F_1 = \frac{2 \cdot \frac{3}{10} \cdot \frac{3}{4}}{(\frac{3}{10} + \frac{3}{4})} = 0.4286$

- (3) 25% recall corresponds to uninterpolated precisions of 100%, 66.67%, 50.00%, 40.00%, 33.33%, 28.57%, 25.00%.
- (4) the interpolated precision for (33%) recall is the maximum precision achieved for  $k \geq 9$ . Obviously, the maximum value is  $\frac{4}{11} = 0.3636$ .
- (5) MAP is  $\frac{1}{8} \cdot (\frac{1}{1} + \frac{2}{2} + \frac{3}{9} + \frac{4}{11} + \frac{5}{15} + \frac{6}{20}) = 0.4163$ .
- (6) The largest possible MAP is  $\frac{1}{8} \cdot (\frac{1}{1} + \frac{2}{2} + \frac{3}{9} + \frac{4}{11} + \frac{5}{15} + \frac{6}{20} + \frac{7}{21} + \frac{8}{22}) = 0.5034$ .
- (7) The smallest possible MAP is  $\frac{1}{8} \cdot (\frac{1}{1} + \frac{2}{2} + \frac{3}{9} + \frac{4}{11} + \frac{5}{15} + \frac{6}{20} + \frac{7}{9999} + \frac{8}{10000}) = 0.4165$ .
- (8)  $0.5034 - 0.4163 = 0.0871$

## Q4. (25 marks)

- (1) The probability distributions for each document model and the background model are:

Model		a	all	am	be	can't	don't	go
background		1/38	2/38	1/38	1/38	1/38	2/38	1/38
doc1	raw	1/22	0	0	1/22	1/22	1/22	1/22
	smoothed	30/836	22/836	11/836	30/836	30/836	41/836	30/836
doc2	raw	0	2/16	1/16	0	0	1/16	0
	smoothed	8/608	54/608	27/608	8/608	8/608	35/608	8/608
Model		groovy	here	hurry	i	is	king	love
background		1/38	1/38	1/38	3/38	1/38	1/38	6/38
doc1	raw	1/22	0	1/22	1/22	0	1/22	3/22
	smoothed	30/836	11/836	30/836	52/836	11/836	30/836	123/836
doc2	raw	0	1/16	0	2/16	1/16	0	3/16
	smoothed	8/608	27/608	8/608	62/608	27/608	8/608	105/608
Model		me	must	of	out	remember	take	tell
background		2/38	1/38	2/38	1/38	1/38	1/38	1/38
doc1	raw	1/22	1/22	1/22	0	0	1/22	0
	smoothed	41/836	30/836	41/836	11/836	11/836	30/836	11/836
doc2	raw	1/16	0	1/16	1/16	1/16	0	1/16
	smoothed	35/608	8/608	35/608	27/608	27/608	8/608	27/608
Model		this	to	want	with	you		
background		1/38	1/38	1/38	1/38	2/38		
doc1	raw	1/22	1/22	1/22	1/22	2/22		
	smoothed	30/836	30/836	30/836	30/836	60/836		
doc2	raw	0	0	0	0	0		
	smoothed	8/608	8/608	8/608	8/608	16/608		

(2)

$$P(Q_1|D_1) = 52/836 * 11/836 * 60/836 = 0.0000587$$

$$P(Q_1|D_2) = 62/608 * 27/608 * 16/608 = 0.000119$$

$$P(Q_2|D_1) = 41/836 * 30/836 * 60/836 * 30/836 * 123/836 * 41/836 = 0.0000000327$$

$$P(Q_2|D_2) = 35/608 * 8/608 * 16/608 * 8/608 * 105/608 * 35/608 = 0.0000000261$$

Thus,  $D_2$  will be ranked first for  $Q_1$  and  $D_1$  will be ranked first for  $Q_2$ .

(3)

$$\frac{P(Q_1|D_1) * P(D_1)}{P(Q_1|D_2) * P(D_2)} = 0.0000587 * 0.7 = 0.0000411$$

$$\frac{P(Q_1|D_2) * P(D_2)}{P(Q_2|D_1) * P(D_1)} = 0.000119 * 0.3 = 0.0000358$$

$$\frac{P(Q_2|D_1) * P(D_1)}{P(Q_2|D_2) * P(D_2)} = 0.0000000327 * 0.7 = 0.0000000229$$

$$\frac{P(Q_2|D_2) * P(D_2)}{P(Q_1|D_1) * P(D_1)} = 0.0000000261 * 0.3 = 0.00000000782$$

Thus,  $D_1$  will be ranked first for both queries by taking into consideration the prior.

$$P(Q_1|D_1)$$

$$P(Q_1|D_1) = P(Q_1|D_1) * P(D_1)$$