

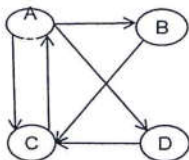
December 2019: END SEMESTER ASSESSMENT (ESA) B.TECH. VII SEMESTER
UE16CS412 ALGORITHMS FOR INFORMATION RETRIEVAL QP

Time: 3 Hours

Answer all the questions

Max Marks: 100

1.	a)	Let D be a document in a text collection. Suppose we add a copy of D to the collection. How would this affect the IDF values of all the words in the collection? Why?	4														
	b)	<div>Given the following index characteristics</div> <table><tr><td>term</td><td> #(postings)</td></tr><tr><td>sky</td><td>189 000</td></tr><tr><td>blue</td><td>230 000</td></tr><tr><td>field</td><td>32 000</td></tr><tr><td>red</td><td>453 000</td></tr><tr><td>high</td><td>345 000</td></tr><tr><td>low</td><td>21 000</td></tr></table> <div>Propose an order to process the following query: (sky OR field) AND (blue OR red) AND (high OR low)</div>	term	#(postings)	sky	189 000	blue	230 000	field	32 000	red	453 000	high	345 000	low	21 000	4
term	#(postings)																
sky	189 000																
blue	230 000																
field	32 000																
red	453 000																
high	345 000																
low	21 000																
	c)	<div>Consider these documents:</div> <ul style="list-style-type: none">• Doc1: solution found for laziness• Doc2: old laziness found• Doc3: old approach for treatment of laziness• Doc4: old hopes for laziness patients <div>i: Draw the term–document incidence matrix for this document collection. ii: Draw the inverted index representation for this collections.</div>	4 +4														
	d)	<div>i. How do you perform stopping and stemming to reduce the size of an inverted index?</div> <div>ii. What is the soundex code for the names for Allan, Allen, Alan, and even Allynn given the mapping of the numbers as (B,F,P, V -> 1) , (C,G, K,Q,S,X,Z -> 2), (D,T -> 3), L->4, (M,N ->5), R->6 (A,E, H,I, J,O,U,W,Y -> 0)</div>	2 2														
2	a)	<div>i. How would you create the dictionary in blocked sort-based indexing on the fly to avoid an extra pass through the data?</div> <div>ii. Bring out the basic principle of Single Pass In Memory Indexing (SPIMI).</div>	4+2														
	b)	<div>i. What is the basic principle of MapReduce for distributed processing? Give an example for illustration.</div> <div>ii. Write a MapReduce program to solve the following problem. Given a large number of credit card transactions, count the total number of transactions for each distinct card number.</div>	6+2														
	c)	<div>You are given the following 3 documents.</div> <div>doc1 : two two tea tea</div> <div>doc2: tea tea me you</div> <div>doc3: me me you you</div> <div>Rank these documents for the query tea me based on similarity using vector space technique.</div>	6														
3	a)	<div>Describe champion list with an example.</div> <div>What is query-term proximity?</div>	4+2														

	b)	Below is a table showing how two human judges rated the relevance of a set of 12 documents to a particular information need (0 = nonrelevant, 1 = relevant). Let us assume that you've written an IR system that for this query returns the set of documents {4, 5, 6, 7, 8}.																																								
		<table><tr><th>DocID</th><th>1</th><th>2</th><th>3</th><th>4</th><th>5</th><th>6</th><th>7</th><th>8</th><th>9</th><th>10</th><th>11</th><th>12</th></tr><tr><td>Judge1</td><td>0</td><td>0</td><td>1</td><td>1</td><td>1</td><td>1</td><td>1</td><td>1</td><td>0</td><td>0</td><td>0</td><td>0</td></tr><tr><td>Judge2</td><td>0</td><td>0</td><td>1</td><td>1</td><td>0</td><td>0</td><td>0</td><td>0</td><td>1</td><td>1</td><td>1</td><td>1</td></tr></table>	DocID	1	2	3	4	5	6	7	8	9	10	11	12	Judge1	0	0	1	1	1	1	1	1	0	0	0	0	Judge2	0	0	1	1	0	0	0	0	1	1	1	1	
DocID	1	2	3	4	5	6	7	8	9	10	11	12																														
Judge1	0	0	1	1	1	1	1	1	0	0	0	0																														
Judge2	0	0	1	1	0	0	0	0	1	1	1	1																														
		i. Calculate the kappa measure between the two judges ii. Calculate precision, recall, and F1 of your system if a document is considered relevant only if the two judges agree. iii. Calculate precision, recall, and F1 of your system if a document is considered relevant if either judge thinks it is relevant	8																																							
	c)	i. What is the intuition behind LSI ii. Correct the statements below. If no change is required, indicate 'no change' 1. The probabilistic IR model doesn't consider relevance. 2. LSI is a method of soft clustering 3. LSI can retrieve documents when the query and document don't share common terms. 4. Interpolated precision reduces the 'jaggedness' of the PR Curve	2+4																																							
4	a)	What do you mean by the term "Google bombing"? Suggest the ways in which the search engine may cope up with this problem.	4																																							
	b)	Given the following two documents and their shingles, calculate the Jaccard Coefficient of $J(d_1, d_2)$ using the concept of comparing the minimum of the hashed values. $h(x) = x \bmod 5$ and $g(x) = (2x+1) \bmod 5$ <table><tr><th>Shingle</th><th>d_1</th><th>d_2</th></tr><tr><td>s_1</td><td>1</td><td>0</td></tr><tr><td>s_2</td><td>0</td><td>1</td></tr><tr><td>s_3</td><td>1</td><td>1</td></tr><tr><td>s_4</td><td>1</td><td>0</td></tr><tr><td>s_5</td><td>0</td><td>1</td></tr></table>	Shingle	d_1	d_2	s_1	1	0	s_2	0	1	s_3	1	1	s_4	1	0	s_5	0	1	6																					
Shingle	d_1	d_2																																								
s_1	1	0																																								
s_2	0	1																																								
s_3	1	1																																								
s_4	1	0																																								
s_5	0	1																																								
	c)	i. List any two features a crawler should provide. ii. As per your understanding what are the two important functions of Mercator Web Crawler as a scalable Web Crawler?	2+2																																							
	d)	i) What are the main factors that influence Page Rank? ii) We have a small web comprising of 4 pages A,B,C and D. Assuming the damping factor as 0.5 compute the page ranks of these pages 	2+4																																							
5	a)	What is attention and briefly explain how is it implemented in text processing?	2+3																																							
	b)	What is a Knowledge Graph and what are some typical examples?	3																																							
	c)	Suppose we have a collection that consists of three documents given below. D1: Malaga Rimini Tibet Tibet Tibet --> class = Y D2: Tibet Tibet sun sun --> class = Y D3: Mexico Malaga Tibet sun --> class= N Assume that we also have the following query: D4 Tibet Malaga . Find the class of D4	6																																							
	d)	Cluster to following documents using K-means with K=2 and cosine similarity. D1: "go monster go" D2: "go karting" D3: "karting monster" D4: "monster monster" Assume D1 and D3 are chosen as initial seeds. Use tf (no idf). Show the clusters and their centroids for each iteration. The algorithm should converge after 2 iterations.	6																																							