



December 2020:END SEMESTER ASSESSMENT (ESA), CSE, VII SEMESTER

UE17CS412–ALGORITHMS FOR INFORMATION RETRIEVAL

Time: 3 Hrs.

Answer All Questions

Max Marks: 100

Provide full calculation for the numerical problems. No partial marking will be done.

Provide brief answers in bullets for all theory questions

1	a)	<p>The question below has one correct answer. In your answer script, write your <u>chosen correct answer with reason in 1-3 sentences</u>:</p> <p>i. In minimum Edit Distance, for common errors :</p> <ul style="list-style-type: none">• cost of such errors is lower• cost of such errors is higher <p>ii. For a large corpus that is static</p> <ul style="list-style-type: none">• Hash table is the best choice for the dictionary and array for posting list• B tree is the best choice for the dictionary and linked list for posting list <p>iii. Soundex algorithm as a query term error correction</p> <ul style="list-style-type: none">• Uses one index• Uses two indexes	6 (3x2)
	b)	<p>i. Write one strength and one weakness of Boolean Query based Information retrieval.</p> <p>ii. What is a possible tradeoff issue in skip pointer implementation in Boolean retrieval?</p>	4 (2+2)
	c)	Write a modified INTERSECT algorithm for XOR operation on two posting lists.	6
	d)	<p>We have a two-word query. For one term the postings list consists of the following 16 entries: [4,6,10,12,14,16,18,20,22,32,47,81,120,122,157,180] and for the other it is the one entry postings list: [47]</p> <p>How many comparisons would be done to intersect the two postings lists using postings lists stored with skip pointers, with a skip length of \sqrt{P} where P is the length of the posting list.</p>	4
2	a)	<p>The question below has one correct answer. In your answer script, write your <u>chosen correct answer with reason in 1-3 sentences</u>:</p> <p>i. If pre-processing is done on corpus using Porter Stemmer even before any compression technique is applied, it is an example of</p> <ol style="list-style-type: none">1. Lossy compression2. Lossless compression	6 (3x2)

		<p>II. Distributed Indexing is for a scenario when</p> <ol style="list-style-type: none"> 1. The index is very large but can be fitted into a large hard disk attached to the server 2. The index is too big to be fitted into a single server <p>III. In TF IDF, the document frequency is used instead of collection frequency as</p> <ol style="list-style-type: none"> 1. collection frequency does not capture frequency of the term well 2. collection frequency does not capture rarity of the term well 	
	b)	<p>i. What is a common limitation of BSBI and SPIMI algorithm ?</p> <p>ii. Why very large block is not advisable in the blocking method to save space on term pointers in case of dictionary compression ?</p>	4 (2+2)
	c)	<p>i. What is the possible significance of IDF (inverse document frequency) for a single term query ? Give reasons.</p> <p>ii. Why Jaccard coefficient is not a good scoring framework for term ? Give any two reasons.</p>	5 (3+2)
	d)	<p>Using the tf-idf model of ranking, find out which document will be listed in front of the user for the information given below.</p> <p>For the terms t1, t2, t3, t4, t5, and t6, the vector of document 1 is (2, 0, 1, 2, 0, 0), the vector of document 2 is (0,3,2,1, 1, 1) and the vector of query is (0, 0, 0,3, 0, 4).</p> <p><i>Assume that these vectors already contain the tf-idf values so you do not have to do any further term weighting.</i></p>	5
3	a)	<p>The question below has one correct answer. In your answer script, write your <u>chosen correct answer with reason in 1-3 sentences</u>:</p> <p>I. In the language model approach of information retrieval, we are finding</p> <ol style="list-style-type: none"> 1. the most relevant language model of query to match the documents 2. the most relevant language model of documents to match the query <p>II. Language model approach and vector space model approach of Information retrieval :</p> <ol style="list-style-type: none"> 1. Vector space model rooted in probability theory and language model rooted in linear algebra 2. Vector space model rooted in linear algebra and language model rooted in probability theory <p>III. If the search engine instruments the mechanism of relevance feedback, the main idea behind that is</p> <ol style="list-style-type: none"> 1. same query should match more relevant documents 2. modified query should match more relevant documents 	6 (3x2)

3

4