

M.Tech (CSE) 3rd Sem-2019
Sub: Natural Language Processing

Time: 3 Hours

Full Mark: 70

(Answer all questions and the figures in the right hand margin indicates marks)

Q1

- (a) What is Boolean retrieval model?
 (b) What is the role of Dictionary and stemmer in retrieval of documents? [7+7]

OR

- (a) How should the Boolean query x and not y be handled? Why is naive evaluation of this query normally very expensive? Write out a postings merge algorithm that evaluates this query efficiently.
 (b) How could a search system combine use of a positional index and use of stop words? What is the potential problem, and how could it be handled? [7+7]

Q2

- (a) What is the idf of a term that occurs in every document? Compare this with the use of stop word lists.
 (b) If we were to stem *jealous* and *jealousy* to a common stem before setting up the vector space, detail how the definitions of tf and idf should be modified. [7+7]

OR

- (a) Consider the four-term query *caught in the rye* and suppose that each of the query terms has five alternative terms suggested by isolated term correction. How many possible corrected phrases must we consider if we do not trim the space of corrected phrases, but instead try all six variants for each of the terms?
 (b) Compute the edit distance between *paris* and *alice*. Write down the 5×5 array of distances between all prefixes as computed by the algorithm. [7+7]

Q3

- (a) Compute variable byte and γ codes for the postings list 7, 17, 29, 31, 45, 90, 110. Use gaps instead of docIDs where possible. Write binary codes in 8-bit blocks.
 (b) From the following sequence of γ -coded gaps, reconstruct first the gap sequence and then the postings sequence: 11100011101010111111011011110110. [7+7]

OR

- (a) Consider the table of term frequencies for 3 documents denoted Doc1, Doc2, Doc3 in Table II. Compute the tf-idf weights for the terms car, auto, insurance, and best, for each document, using the idf values from Table I.

Table-I			Table-II			
term	df _t	idf _t	car	Doc1	Doc2	Doc3
car	18,165	1.65	27	4	24	
auto	6723	2.08	3	33	0	
insurance	19,241	1.62	0	33	29	
best	25,235	1.5	14	0	17	

- (b) From the following sequence of γ -coded gaps, reconstruct first the gap sequence and then the postings sequence: 111000101101010101111101101111011. [7+7]

Q4

- (a) Write down the steps of a basic crawler.
 (b) State various implementation issues of a crawler. [7+7]

OR

- (a) What is a crawler? Explain how the various fundamentals step like tokenization, linguistic modules and indexer are handled by the crawler. [14]

Q5 Write short notes on the following:

- (a) Page Rank Algorithm
 (b) HITS Algorithm [7+7]

OR

- (a) A user of a browser can, in addition to clicking a hyperlink on the page x he is currently browsing, use the back button to go back to the page from which he arrived at x . Can such a user of back buttons be modeled as a Markov chain? How would we model repeated invocations of the back button?
 (b) Consider a web graph with three nodes 1, 2, and 3. The links are as follows: 1→2, 3→2, 2→1, 2→3. Write down the transition probability matrices for the surfer's walk with teleporting, for the following three values of the teleport probability: (i) $\alpha = 0$; (ii) $\alpha = 0.5$; and (iii) $\alpha = 1$. [7+7]