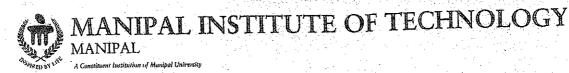
				**							1	
1							1	2.5	1000		1	, ,
- 1			1 1				7.4			1	1 '	ŧ. I
- 1			- I - i		1. 1						•	1 1
•	Reg. I	NΛ	1 1				1 1		l.		1	
	i Neu. i	WO.	. 1	10.00				1.0				1. 1
			11. 1			1			5	5	1 1 1	



## V SEMESTER B.TECH. (INFORMATION TECHNOLOGY / COMPUTER AND COMMUNICATION ENGINEERING) END SEMESTER EXAMINATIONS, NOV/DEC 2016

SUBJECT: PROGRAM ELECTIVE I - INFORMATION RETRIEVAL [ICT 4006]

## REVISED CREDIT SYSTEM (05/12/2016)

MAX. MARKS: 50 Time: 3 Hours

## Instructions to Candidates:

Answer ALL the questions

Α.	Write an algorithm for posting list intersection with skip pointers.  Consider a postings intersection between this postings list, with skip pointers:	5
	pl: 3 5 9 15 24 39 60 68 75 81 84 89 92 96 97 100 115	
	and the following intermediate result postings list (which hence has no skip	
	pointers):	
	$2 \le 80 \ 05 \ 97 \ 99 \ 100 \ 101$	
	Trace through the postings intersection algorithm and answer the following	
	anarios	
	(i) How often is a skin pointer followed (i.e., p) is advanced to skip $(p^1)$ ):	
	(i) How many postings comparisons will be made by this algorithm while	
	interpretation the two lists?	
	(iii) How many postings comparisons would be made if the postings lists are	
	intercapted without the use of skip pointers?	,
1B.	. Describe the Boolean retrieval model. Consider the following document concentric.	•
	Doc 1: new home sales top forecasts	
	Doc 2: home sales rise in july	
	INCO M. HOMO Bears	
	Doc 3: increase in home sales in july	
	Doc 3: increase in home sales in july	
	Doc 3: increase in home sales in july  Doc 4: july new home sales rise  (i) Draw the term-document incidence matrix for this document collection.  (ii) Answer the following query using Boolean retrieval model.	

- home AND sales AND (july OR rise). Explain the following with an example for each.
  - (i) biword index
  - (ii) positional index
- Explain the various dictionary compression techniques with examples. Also discuss the limitations of each compression technique.

Page 1 of 2

2

Compute the edit distance between the strings paris and alice. 2B. 2 From the following sequence of  $\gamma$  -coded gaps, reconstruct first the gap sequence and 2C. then the postings sequence: 11011111000111010101111111011011111011 5 Explain the complete IR search system with a neat diagram. Table Q.3B shows how two human judges rate the relevance of a set of 12 3B. documents to a particular information need (0 = nonrelevant, 1 = relevant). Let us assume that an IR system has been developed which returns for a query the set of documents {4, 5, 6, 7, 8}. Table Q.3B 10 12 Doc Id 5 8 6 0 Judge 1 1 0 0 0 0 0 Judge 2 (i) Calculate the kappa measure between the two judges. (ii) Calculate precision, recall, and F<sub>1</sub> of the system if a document is considered relevant if either judge thinks it is relevant. 2 What do you understand by the term Shingling? Why is it used in web search? Consider a case insensitive document collection with a query (q) and a document collection consisting of the following three documents: q: gold silver truck dl: shipment of gold damaged in a fire d2: delivery of silver arrived in a silver truck d3: shipment of gold arrived in a truck Assume that the document vector is formed using tf-idf weighting scheme. Rank these documents based on the following similarity score formula:  $SC(q, d_i) = \frac{2\sum_{j=1}^{t} w_{qj} d_{ij}}{\sum_{j=1}^{t} (d_{ij})^2 \sum_{j=1}^{t} (w_{qj})^2}.$ 3 4B. Describe a technique for the low-rank approximation of a given matrix A such that the approximated matrix satisfies Eckart- Young theorem. If the query is: friends AND romans AND (NOT countrymen), how could we use the 4C. frequency of countrymen in evaluating the best query evaluation order? In particular, propose a way of handling negation in determining the order of query processing. 5A. Consider a web graph with three nodes 1, 2 and 3. The links are as follows:  $1\rightarrow 2$ ,  $1\rightarrow 3$ ,  $2\rightarrow 3$ ,  $3\rightarrow 2$ . Compute PageRank after three iterations, hub and authority scores for each of the three pages. Also give the relative ordering of the three nodes for each

of these scores indicating any ties. Assume that at each step of the PageRank random walk, we teleport to a random page with a probability 0.1, with a uniform distribution

Explain URL frontier components of web crawler with a neat diagram.

Write the Block Sort Based Indexing (BSBI) algorithm.

over which particular page we teleport to.

5B.