



PES University, Bangalore
(Established under Karnataka Act No. 16 of 2013)

UE14CS414
(Prq. Ananthanarayan)

**END SEMESTER ASSESSMENT (ESA) B.TECH. VII SEMESTER- Dec.
2017**

UE14CS414 – Algorithms for Information Retrieval

Time: 3 Hrs

Answer All Questions

Max Marks: 100

General instructions:

- In case of ambiguity, feel free to make reasonable assumptions. Clearly state your assumptions.
- Answer the questions in the order
- Answer neatly and legibly. Use the margin only for writing question number and not for numbering any lists you have as part of your answer.

1.	a)	I intend to use boolean matching to evaluate your answer paper. Is this acceptable? Justify your answer.	3												
	b)	How do you optimize the execution of boolean queries?	4												
	c)	Name the methods used to support phrase queries	3												
2.	a)	Explain how dynamic indexing works.	7												
	b)	What is a parametric index? What is a zone, how does weighted zone scoring work? Note: you do not have to derive the weight scores.	7												
	c)	State Zipf's law. You have a collection of documents containing exactly 4 terms a, b, c, and d, with frequencies $a > b > c > d$. If the number of tokens in the collection is 5000, what are the frequencies of the individual terms.	6												
3.	a)	Your search system retrieved the 5 documents for a query and the relevance scores are captured below. Assuming a total of 10 relevant documents, draw the PR curve. <table border="1"><thead><tr><th>Document</th><th>Relevant</th></tr></thead><tbody><tr><td>1</td><td>Y</td></tr><tr><td>2</td><td>Y</td></tr><tr><td>3</td><td>N</td></tr><tr><td>4</td><td>N</td></tr><tr><td>5</td><td>Y</td></tr></tbody></table>	Document	Relevant	1	Y	2	Y	3	N	4	N	5	Y	6
	Document	Relevant													
1	Y														
2	Y														
3	N														
4	N														
5	Y														
b)	Explain the following methods used to speed up cosine similarity calculations: 1. Champion lists 2. Cluster pruning	5													

3	c)	Correct the statements below. If no change is required, indicate 'no change' 1. The probabilistic IR model doesn't consider relevance. 2. An empty document has the same probability of relevance and not relevance. 3. Principal Component Analysis is a subset of Latent Semantic Indexing 4. LSI is a method of soft clustering 5. LSI can retrieve documents when the query and document don't share common terms. 6. Low-rank approximations of SVD maximise the Frobenius norm	6x1 =6																																	
	d)	Calculate the kappa statistic for the following confusion matrix. Use pooled marginals in your calculation. <table><tr><td></td><td>Yes</td><td>No</td></tr><tr><td>Yes</td><td>100</td><td>20</td></tr><tr><td>No</td><td>10</td><td>70</td></tr></table>		Yes	No	Yes	100	20	No	10	70	6																								
	Yes	No																																		
Yes	100	20																																		
No	10	70																																		
4.	a)	What are the mandatory and the recommended features of a web crawler?	5																																	
	b)	Explain the working of the URL frontier of a web crawler.	6																																	
	c)	Calculate the pairwise document similarity for the following documents using 3-shingles. D1: The sky is clear and the moon is shining D2: The moon is shining so are the stars D3: Where are the stars and the moon today?	6																																	
	d)	How are hub and authority scores defined? Calculate the hub and authority scores using the following adjacency matrix. Stop your calculation after the first iteration. <table><tr><td></td><td></td><td colspan="4">To</td></tr><tr><td></td><td></td><td>1</td><td>2</td><td>3</td><td>4</td></tr><tr><td rowspan="4">From</td><td>1</td><td>0</td><td>1</td><td>0</td><td>0</td></tr><tr><td>2</td><td>1</td><td>0</td><td>1</td><td>0</td></tr><tr><td>3</td><td>1</td><td>0</td><td>0</td><td>0</td></tr><tr><td>4</td><td>1</td><td>0</td><td>1</td><td>0</td></tr></table>			To						1	2	3	4	From	1	0	1	0	0	2	1	0	1	0	3	1	0	0	0	4	1	0	1	0	5
		To																																		
		1	2	3	4																															
From	1	0	1	0	0																															
	2	1	0	1	0																															
	3	1	0	0	0																															
	4	1	0	1	0																															

5	a)	What is the Rocchio classification algorithm? Calculate the centroids for the following document classes. Use raw tf (no IDF) for the vector.	8																																				
		<table><tr><td></td><td></td><td>Class label</td></tr><tr><td>D1</td><td>hybrid battery efficient</td><td>electric</td></tr><tr><td>D2</td><td>lithium battery recharge</td><td>electric</td></tr><tr><td>D3</td><td>energy recharge mileage</td><td>electric</td></tr><tr><td>D4</td><td>'Fossil fuel' pollution hybrid</td><td>fuel</td></tr><tr><td>D5</td><td>Non-renewable mileage energy</td><td>fuel</td></tr></table>			Class label	D1	hybrid battery efficient	electric	D2	lithium battery recharge	electric	D3	energy recharge mileage	electric	D4	'Fossil fuel' pollution hybrid	fuel	D5	Non-renewable mileage energy	fuel																			
		Class label																																					
D1	hybrid battery efficient	electric																																					
D2	lithium battery recharge	electric																																					
D3	energy recharge mileage	electric																																					
D4	'Fossil fuel' pollution hybrid	fuel																																					
D5	Non-renewable mileage energy	fuel																																					
	b)	Use k-NN to classify the following document against the document classes in 5a. Use Euclidean distance calculation. Test document: energy efficient recharge	8																																				
	c)	Merge the following clusters using complete link clustering. Draw the dendrogram. Use $\text{dist}(xy, i) = \max(\text{dist}(y, i), \text{dist}(x, i))$ where xy is formed by merging clusters x & y and i is any other cluster.	8																																				
		<table><tr><td></td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td></tr><tr><td>1</td><td>0</td><td></td><td></td><td></td><td></td></tr><tr><td>2</td><td>9</td><td>0</td><td></td><td></td><td></td></tr><tr><td>3</td><td>3</td><td>7</td><td>0</td><td></td><td></td></tr><tr><td>4</td><td>6</td><td>5</td><td>9</td><td>0</td><td></td></tr><tr><td>5</td><td>11</td><td>10</td><td>2</td><td>8</td><td>0</td></tr></table>		1	2	3	4	5	1	0					2	9	0				3	3	7	0			4	6	5	9	0		5	11	10	2	8	0	
	1	2	3	4	5																																		
1	0																																						
2	9	0																																					
3	3	7	0																																				
4	6	5	9	0																																			
5	11	10	2	8	0																																		
	d)	You deserve one really easy question so here it is: Draw a smiley face.	1																																				