

17/5/22

K. J. Somaiya Institute of Engineering and Information Technology, Sion, Mumbai
(An Autonomous Institute Affiliated to the University of Mumbai)

End Semester Exam
April – May 2022

B.Tech. (Information Technology)

Examination: TY - Semester VI

Course Code: IUITC601 Course Name: Data Mining and Business Intelligence

Duration: 03 Hours

Max. Marks: 60

Instructions:

- (1) All questions are compulsory.
- (2) Draw neat diagrams wherever applicable.
- (3) Assume suitable data, if necessary.

Ques. No.	Question	Max. Marks	CO	BT Level																		
Q1.	Solve any six questions out of eight:	12																				
i)	Explain Drill Down operation with example.	2	CO1	U																		
ii)	Compute the Dissimilarity Matrix for 04 instances of an attribute 'Ratings' represented as: (Good, Fair, Average, Good).	2	CO2	U																		
iii)	List different approaches for handling missing data.	2	CO3	U																		
iv)	Normalize data 15, 250, 250000 in the range [0, 1] using Decimal Scaling.	2	CO3	U																		
v)	Explain two-step process of Classification.	2	CO4	U																		
vi)	Explain Market Basket Analysis.	2	CO4	U																		
vii)	Explain Support And Confidence with example.	2	CO5	U																		
viii)	Explain Decision Support Systems.	2	CO6	U																		
Q2.	Solve any four questions out of six:	16																				
i)	Explain Data Warehousing.	4	CO1	U																		
ii)	Sketch a boxplot and identify the outliers in the data of <i>age</i> of participants in an event: 12, 51, 56, 60, 65, 70, 77.	4	CO2	A																		
iii)	Sketch a Concept Hierarchy for <i>Computer Accessories</i> .	4	CO3	A																		
iv)	Consider data of 5 objects characterized by a single continuous feature <i>score</i> : <table border="1" style="margin-left: auto; margin-right: auto;"> <thead> <tr> <th>Player 1</th> <th>Player 2</th> <th>Player 3</th> <th>Player 4</th> <th>Player 5</th> </tr> </thead> <tbody> <tr> <td>10</td> <td>20</td> <td>40</td> <td>50</td> <td>60</td> </tr> </tbody> </table> Assume that there are two clusters: C1: {a, b}, and C2: {c, d, e}. Calculate the distance matrix. Also calculate single link, complete link, and average distance between C1 and C2 to evaluate clusters' quality.	Player 1	Player 2	Player 3	Player 4	Player 5	10	20	40	50	60	4	CO4	A								
Player 1	Player 2	Player 3	Player 4	Player 5																		
10	20	40	50	60																		
v)	State the Apriori principle. Consider the below 2-itemsets in vertical data format and identify the acceptable 3-itemsets: <table border="1" style="margin-left: auto; margin-right: auto;"> <thead> <tr> <th>Itemset</th> <th>TID_set</th> </tr> </thead> <tbody> <tr> <td>{A, B}</td> <td>{T1, T4, T8, T9}</td> </tr> <tr> <td>{A, C}</td> <td>{T5, T7, T8, T9}</td> </tr> <tr> <td>{A, D}</td> <td>{T4}</td> </tr> <tr> <td>{A, E}</td> <td>{T1, T8}</td> </tr> <tr> <td>{B, C}</td> <td>{T3, T6, T8, T9}</td> </tr> <tr> <td>{B, D}</td> <td>{T2, T4}</td> </tr> <tr> <td>{B, E}</td> <td>{T1, T8}</td> </tr> <tr> <td>{C, E}</td> <td>{T8}</td> </tr> </tbody> </table>	Itemset	TID_set	{A, B}	{T1, T4, T8, T9}	{A, C}	{T5, T7, T8, T9}	{A, D}	{T4}	{A, E}	{T1, T8}	{B, C}	{T3, T6, T8, T9}	{B, D}	{T2, T4}	{B, E}	{T1, T8}	{C, E}	{T8}	4	CO5	A
Itemset	TID_set																					
{A, B}	{T1, T4, T8, T9}																					
{A, C}	{T5, T7, T8, T9}																					
{A, D}	{T4}																					
{A, E}	{T1, T8}																					
{B, C}	{T3, T6, T8, T9}																					
{B, D}	{T2, T4}																					
{B, E}	{T1, T8}																					
{C, E}	{T8}																					
vi)	Explain the cycle of Business Intelligence analysis.	4	CO6	U																		

Q3. Solve any two questions out of three:		16																																										
i)	Sketch and explain the process of Knowledge Discovery from Data.	8	CO1	U																																								
ii)	Compute the highest similarity between two financial articles with term frequency vectors as follows using Cosine Similarity:	8	CO2	A																																								
	<table border="1"> <thead> <tr> <th>Article No.</th> <th>Digital Rupee</th> <th>Budget</th> <th>Blockchain</th> <th>Digital Assets</th> <th>Income</th> <th>GDP</th> <th>Disinvestments</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>2</td> <td>4</td> <td>3</td> <td>2</td> <td>3</td> <td>2</td> <td>1</td> </tr> <tr> <td>2</td> <td>5</td> <td>0</td> <td>3</td> <td>4</td> <td>0</td> <td>0</td> <td>2</td> </tr> <tr> <td>3</td> <td>3</td> <td>0</td> <td>6</td> <td>3</td> <td>0</td> <td>0</td> <td>3</td> </tr> <tr> <td>4</td> <td>1</td> <td>3</td> <td>4</td> <td>2</td> <td>7</td> <td>4</td> <td>2</td> </tr> </tbody> </table>				Article No.	Digital Rupee	Budget	Blockchain	Digital Assets	Income	GDP	Disinvestments	1	2	4	3	2	3	2	1	2	5	0	3	4	0	0	2	3	3	0	6	3	0	0	3	4	1	3	4	2	7	4	2
	Article No.				Digital Rupee	Budget	Blockchain	Digital Assets	Income	GDP	Disinvestments																																	
	1				2	4	3	2	3	2	1																																	
	2				5	0	3	4	0	0	2																																	
3	3	0	6	3	0	0	3																																					
4	1	3	4	2	7	4	2																																					
A survey was conducted to analyse if <i>discount sale</i> has a correlation with the ratings given by customers. Apply Chi-square test to analyze whether the attributes <i>Discount_Sale</i> and <i>Customer_Ratings</i> are correlated or not:		8	CO3	A																																								
<table border="1"> <thead> <tr> <th rowspan="2">Discount_Sale</th> <th colspan="2">Customer_Ratings</th> <th rowspan="2">Total</th> </tr> <tr> <th>Good</th> <th>Average</th> </tr> </thead> <tbody> <tr> <th>Yes</th> <td>250</td> <td>300</td> <td>550</td> </tr> <tr> <th>No</th> <td>50</td> <td>1100</td> <td>1150</td> </tr> <tr> <th>Total</th> <td>300</td> <td>1400</td> <td>1700</td> </tr> </tbody> </table>					Discount_Sale	Customer_Ratings		Total	Good	Average	Yes	250	300	550	No	50	1100	1150	Total	300	1400	1700																						
Discount_Sale	Customer_Ratings					Total																																						
	Good	Average																																										
Yes	250	300	550																																									
No	50	1100	1150																																									
Total	300	1400	1700																																									
Consider that for 1 degrees of freedom, the chi-square value to reject the hypothesis at 0.001 significance level is 10.828.																																												
Q4. Solve any two questions out of three:		16																																										
i)	Suppose that the data mining task is to cluster points A1(2, 10), A2(2, 5), A3(8, 4), B1(5, 8), B2(7, 5), B3(6, 4), C1(1, 2), C2(4, 9) representing coordinates of location (x, y) into 3 clusters. Suppose initially A1, B1, and C1 are represented as the centers of each clusters respectively. Use Euclidean distance and apply the <i>k-means</i> algorithm to show: a. The three cluster centers after the first round of execution. b. The final three clusters.	8	CO4	A																																								
ii)	Consider an <i>Electronics</i> store that surveyed customers likely to buy computer from them, based on the data of past customers. Apply the given Decision Tree classifier on the below mentioned <u>test</u> dataset and calculate Accuracy, Precision, and Recall, Sensitivity, Specificity.	8	CO5	A																																								
	<table border="1"> <thead> <tr> <th>Age</th> <th>Income</th> <th>Student</th> <th>Credit Rating</th> <th>Buys Computer</th> </tr> </thead> <tbody> <tr> <td>Youth</td> <td>High</td> <td>No</td> <td>Fair</td> <td>No</td> </tr> <tr> <td>Senior</td> <td>High</td> <td>No</td> <td>Excellent</td> <td>No</td> </tr> <tr> <td>Middle</td> <td>High</td> <td>No</td> <td>Fair</td> <td>No</td> </tr> <tr> <td>Senior</td> <td>Medium</td> <td>No</td> <td>Fair</td> <td>Yes</td> </tr> <tr> <td>Senior</td> <td>Low</td> <td>Yes</td> <td>Excellent</td> <td>Yes</td> </tr> <tr> <td>Middle</td> <td>Low</td> <td>Yes</td> <td>Excellent</td> <td>Yes</td> </tr> </tbody> </table>				Age	Income	Student	Credit Rating	Buys Computer	Youth	High	No	Fair	No	Senior	High	No	Excellent	No	Middle	High	No	Fair	No	Senior	Medium	No	Fair	Yes	Senior	Low	Yes	Excellent	Yes	Middle	Low	Yes	Excellent	Yes					
Age	Income	Student	Credit Rating	Buys Computer																																								
Youth	High	No	Fair	No																																								
Senior	High	No	Excellent	No																																								
Middle	High	No	Fair	No																																								
Senior	Medium	No	Fair	Yes																																								
Senior	Low	Yes	Excellent	Yes																																								
Middle	Low	Yes	Excellent	Yes																																								
iii)	Consider the case of classifying credit card transactions as legitimate or fraudulent. Apply KDD process to derive Business Intelligence. Clearly explain each phase / operation in the KDD process with respect to the stated application.	8	CO6	A																																								
