

University of Mumbai
Examinations Summer 2022

Time: 2 hour 30 minutes

Max. Marks: 80

T.E.(Information Technology Engineering)(SEM-VI)(Choice Base Credit Grading System) (R-2020-21) (C Scheme) 789381 - Data Mining & Business Intelligence DATE: 18/5/2022 QP CODE: 91760

Q.1	Choose the correct option for following questions. All the Questions are compulsory and carry equal marks (2 marks each)																											
1.	If dimensionality reduction is performed on a record data matrix, the transformed data matrix_____																											
Option A:	has reduced number of rows																											
Option B:	has reduced number of columns																											
Option C:	has reduced number of both rows and columns																											
Option D:	has same number of rows and columns																											
2.	Consider the following data: 4, 8, 9, 15, 21, 21, 24, 25, 26, 28, 29, 34. Partition the given data with Bin size: 4. What is the output obtained after smoothing the data by Bin Boundaries.																											
Option A:	Bin 1: 4, 4, 4, 15 Bin 2: 21, 21, 25, 25 Bin 3: 26, 26, 26, 34																											
Option B:	Bin 1: 4, 4, 15, 15 Bin 2: 21, 21, 21, 25 Bin 3: 26, 26, 34, 34																											
Option C:	Bin 1: 4, 15, 15, 15 Bin 2: 21, 25, 25, 25 Bin 3: 26, 26, 26, 34																											
Option D:	Bin 1: 4, 4, 4, 15 Bin 2: 21, 25, 25, 25 Bin 3: 26, 26, 26, 34																											
3.	Knowledge discovery in databases is referred to																											
Option A:	Non Trivial process of choosing dataset																											
Option B:	Non Trivial process for identifying useful patterns in data																											
Option C:	Non Trivial process for identifying invalid patterns in data																											
Option D:	Non Trivial process of creating patterns in data																											
4.	For the given confusion matrix compute recall <table border="1" style="margin-left: auto; margin-right: auto;"> <tr> <td colspan="2"></td> <td colspan="3" style="text-align: center;">Predicted data</td> </tr> <tr> <td colspan="2"></td> <td style="text-align: center;">Cancer Classes</td> <td style="text-align: center;">Yes</td> <td style="text-align: center;">No</td> <td style="text-align: center;">Total</td> </tr> <tr> <td rowspan="3" style="text-align: center;">Actual data</td> <td style="text-align: center;">Yes</td> <td></td> <td style="text-align: center;">90</td> <td style="text-align: center;">210</td> <td style="text-align: center;">300</td> </tr> <tr> <td style="text-align: center;">No</td> <td></td> <td style="text-align: center;">140</td> <td style="text-align: center;">9560</td> <td style="text-align: center;">9700</td> </tr> <tr> <td style="text-align: center;">Total</td> <td></td> <td style="text-align: center;">230</td> <td style="text-align: center;">9770</td> <td style="text-align: center;">10000</td> </tr> </table>			Predicted data					Cancer Classes	Yes	No	Total	Actual data	Yes		90	210	300	No		140	9560	9700	Total		230	9770	10000
		Predicted data																										
		Cancer Classes	Yes	No	Total																							
Actual data	Yes		90	210	300																							
	No		140	9560	9700																							
	Total		230	9770	10000																							
Option A:	20%																											
Option B:	30%																											
Option C:	40%																											
Option D:	45%																											
5.	You are given reviews of food quality of few restaurants as Good, Average or Poor. Finding reviews of a new restaurant is an example of_____																											
Option A:	Classification																											
Option B:	Regression																											
Option C:	Clustering																											
Option D:	Association mining																											

6.	BIRCH falls under which clustering approach												
Option A:	Partitioning approach												
Option B:	Hierarchical approach												
Option C:	Density-based approach												
Option D:	Distribution based approach												
7.	Given {2,4,3,10,11,12,20,25,30}, Assume k=2 and initial means are m1=4, m2=11. Apply k -means clustering technique and find its output after 1st iteration												
Option A:	K1= {2,3,4,10,11,12} K2= {20,30,25}												
Option B:	K1= {2,3,4} K2= {10,11,12,20,30,25}												
Option C:	K1= {2,3} K2={4,10,11,12,20,30,25}												
Option D:	K1= {2,3,4,10} K2={11,12,20,30,25}												
8.	In one of the frequent item-set examples, it is observed that if milk and bread are bought then eggs are also purchased by the customers. After generating an association rule among the given set of items, it is inferred												
Option A:	{Milk} is antecedent and {eggs} is consequent												
Option B:	{Milk} is antecedent and the item set {bread, eggs} is consequent												
Option C:	The item set {milk, bread} is consequent and {eggs} is antecedent												
Option D:	The item set {milk, bread} is antecedent and {eggs} is consequent												
9.	For the given transactional database compute confidence for the rule Milk \Rightarrow Beer												
	<table border="1"> <thead> <tr> <th>TID</th> <th>Items</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>Bread, Milk</td> </tr> <tr> <td>2</td> <td>Bread, Diaper, Beer, Eggs</td> </tr> <tr> <td>3</td> <td>Milk, Diaper, Beer, Coke</td> </tr> <tr> <td>4</td> <td>Bread, Milk, Diaper, Beer</td> </tr> <tr> <td>5</td> <td>Bread, Milk, Diaper, Coke</td> </tr> </tbody> </table>	TID	Items	1	Bread, Milk	2	Bread, Diaper, Beer, Eggs	3	Milk, Diaper, Beer, Coke	4	Bread, Milk, Diaper, Beer	5	Bread, Milk, Diaper, Coke
TID	Items												
1	Bread, Milk												
2	Bread, Diaper, Beer, Eggs												
3	Milk, Diaper, Beer, Coke												
4	Bread, Milk, Diaper, Beer												
5	Bread, Milk, Diaper, Coke												
Option A:	20%												
Option B:	50%												
Option C:	40%												
Option D:	60%												
10.	_____ is an interactive computer-based application that combines data and mathematical models to help decision makers solve complex problems faced in managing the public and private enterprises and organizations.												
Option A:	Data Mining												
Option B:	Data dredging												
Option C:	Decision support system												
Option D:	Artificial Intelligence system												

Q.2 Solve any Two Questions out of Three

Marks

- A** Define data warehouse. Describe different OLAP operations in detail **10**
- B** Apply Naive Bayes classifier algorithm to the dataset given below, and classify the unknown data sample? **10**
Given all the previous patients I've seen(below are their symptoms and their diagnosis)

chills	runny nose	headache	fever	flu ?
Y	N	Mild	Y	N
Y	Y	No	N	Y
Y	N	Strong	Y	Y
N	Y	Mild	Y	Y
N	N	No	N	N
N	Y	Strong	Y	Y
N	Y	Strong	N	N
Y	Y	Mild	Y	Y

Do I believe that patient with following symptoms has the flu?

chills	runny nose	headache	fever	flu ?
Y	N	Mild	Y	?

- C** Explain multi-level and multidimensional association rules with example **10**

Q.3 Solve any Two Questions out of Three

- A** Suppose we have six objects with name A, B, C, D, E and F. Apply single linkage clustering and draw dendrogram for the given data. **10**

	X	Y
A	1	1
B	1.5	1.5
C	5	5
D	3	4
E	4	4
F	3	3.5

- B** Suppose the data for analysis includes the attribute age. The age values for data tuples are (in increasing order):
13,15,16,16,19,20,20,21,22,22,25,25,25,25,30,33,33,35,35,35,35,36,40,45,46,52,70 **10**
 - i) What is mean of data? What is median of data?
 - ii) What is mode of data? Comment on data's modality.
 - iii) What is mid range of data?
 - iv) Give the five point summary of the data.
 - v) Show box plot of the data
- C** What is Business Intelligence (BI)? Explain BI architecture in detail **10**

Q.4 Solve any Two Questions out of Three

- A** Briefly explain Bagging and Boosting of classifiers **10**
- B** For the table given, apply Apriori algorithm and show frequent item set and strong association rules. Assume Minimum Support of 30% and Minimum confidence of 70%. **10**

TID	Items
01	1,3,4,6
02	2,3,5,7
03	1,2,3,5,8
04	2,5,9,10
05	1,4

- C** What is an outlier? Describe methods used for outlier analysis. **10**
