

Name:

Enrolment No:



UNIVERSITY OF PETROLEUM AND ENERGY STUDIES

End Semester Examination, December 2020

Course: Data Mining & Predictive Analytics

Program: B.Tech CSE ECRA

Course Code: CSBA 4005

Semester: VII

Time : 03 hrs.

Max. Marks: 100

Instructions:

1. In Section A, you have to write one word/one sentence answers, no explanation, no calculation is required to be furnished.
2. Section B and C are the sections in which you will be writing the answers on A4 sheets and after clicking the picture, upload as per the direction.

SECTION A

S. No.		Marks	CO
Q 1	Discuss whether or not each of the following activities is a data mining task: a) Establishing a cricket player into All-rounder category by analyzing his batting and bowling statistics. b) Monitoring the oxygen level of a patient for abnormalities. c) Predicting the outcome of a hockey match based on past performances of two teams. d) Predicting the outcome of tossing pair of dice. e) Monitoring climatic conditions for tsunami.	5	CO1
Q 2	Name THREE data mining task which are <i>descriptive</i> in nature.	5	CO2
Q 3	State True/False a) ROC curve is drawn between TP and TN values. b) Accuracy of a classifier = 1- Error rate c) Accuracy of a Rule-based classifier is written as $\text{accuracy}(R) = \frac{n_{\text{correct}}}{n_{\text{covers}}}$ d) DIANA and AGNES are two classification algorithms. e) Full form of ID in ID3 is	5	CO3
Q 4	In k-fold validation method, what does <i>k</i> denote? What is the value of <i>k</i> which is taken as standard?	5	CO4
Q 5	Name FIVE important characteristics of structured data.	5	CO1
Q 6	a) If there are N data objects, each having P attributes, what will be the dimension of Proximity Matrix? b) In a data set, a variable is having values in the range 1000 to 9000. We wish to normalize these to a new range 0-5. What will be the equivalent of 5000 while we use min-max normalization? c) What will be the Euclidean distance between two data points (5,7,10) and (6,8,2)? d) If we do the partitioning of dataset, and pick up the proportional volume from each partition, which type of sampling this is called? e) Name THREE data visualization techniques.	5	CO2

SECTION B

	For a given dataset – Classify whether the following student will buy or not, using Naïve Bayes Classifier <i>(Senior, high, No, excellent,?)</i>	10	CO3
--	--	----	-----

Q 7	<table border="1"> <thead> <tr> <th>RID</th> <th>age</th> <th>income</th> <th>student</th> <th>credit_rating</th> <th>Class: buys_computer</th> </tr> </thead> <tbody> <tr><td>1</td><td>youth</td><td>high</td><td>no</td><td>fair</td><td>no</td></tr> <tr><td>2</td><td>youth</td><td>high</td><td>no</td><td>excellent</td><td>no</td></tr> <tr><td>3</td><td>middle_aged</td><td>high</td><td>no</td><td>fair</td><td>yes</td></tr> <tr><td>4</td><td>senior</td><td>medium</td><td>no</td><td>fair</td><td>yes</td></tr> <tr><td>5</td><td>senior</td><td>low</td><td>yes</td><td>fair</td><td>yes</td></tr> <tr><td>6</td><td>senior</td><td>low</td><td>yes</td><td>excellent</td><td>no</td></tr> <tr><td>7</td><td>middle_aged</td><td>low</td><td>yes</td><td>excellent</td><td>yes</td></tr> <tr><td>8</td><td>youth</td><td>medium</td><td>no</td><td>fair</td><td>no</td></tr> <tr><td>9</td><td>youth</td><td>low</td><td>yes</td><td>fair</td><td>yes</td></tr> <tr><td>10</td><td>senior</td><td>medium</td><td>yes</td><td>fair</td><td>yes</td></tr> <tr><td>11</td><td>youth</td><td>medium</td><td>yes</td><td>excellent</td><td>yes</td></tr> <tr><td>12</td><td>middle_aged</td><td>medium</td><td>no</td><td>excellent</td><td>yes</td></tr> <tr><td>13</td><td>middle_aged</td><td>high</td><td>yes</td><td>fair</td><td>yes</td></tr> <tr><td>14</td><td>senior</td><td>medium</td><td>no</td><td>excellent</td><td>no</td></tr> </tbody> </table>	RID	age	income	student	credit_rating	Class: buys_computer	1	youth	high	no	fair	no	2	youth	high	no	excellent	no	3	middle_aged	high	no	fair	yes	4	senior	medium	no	fair	yes	5	senior	low	yes	fair	yes	6	senior	low	yes	excellent	no	7	middle_aged	low	yes	excellent	yes	8	youth	medium	no	fair	no	9	youth	low	yes	fair	yes	10	senior	medium	yes	fair	yes	11	youth	medium	yes	excellent	yes	12	middle_aged	medium	no	excellent	yes	13	middle_aged	high	yes	fair	yes	14	senior	medium	no	excellent	no		
RID	age	income	student	credit_rating	Class: buys_computer																																																																																								
1	youth	high	no	fair	no																																																																																								
2	youth	high	no	excellent	no																																																																																								
3	middle_aged	high	no	fair	yes																																																																																								
4	senior	medium	no	fair	yes																																																																																								
5	senior	low	yes	fair	yes																																																																																								
6	senior	low	yes	excellent	no																																																																																								
7	middle_aged	low	yes	excellent	yes																																																																																								
8	youth	medium	no	fair	no																																																																																								
9	youth	low	yes	fair	yes																																																																																								
10	senior	medium	yes	fair	yes																																																																																								
11	youth	medium	yes	excellent	yes																																																																																								
12	middle_aged	medium	no	excellent	yes																																																																																								
13	middle_aged	high	yes	fair	yes																																																																																								
14	senior	medium	no	excellent	no																																																																																								
Q 8	Discuss various important and desirous characteristics of data which are really required for a reasonable good outcome of any data mining task. Here, in between you can discuss the various types of data for which data mining is applied.	10	CO2																																																																																										
Q 9	<p>For what purpose <i>k</i>-Means algorithm is used. Write various standard steps of this algorithm. How the values of <i>k</i> are decided?</p> <p style="text-align: center;">OR</p> <p>Perform KNN- Classification algorithm on following dataset and predict the class for X (Height=1633 and Weight=57). Given K=3.</p> <table border="1" style="margin-left: auto; margin-right: auto;"> <thead> <tr> <th>Height</th> <th>Weight</th> <th>Class</th> </tr> </thead> <tbody> <tr><td>157</td><td>45</td><td>Underweight</td></tr> <tr><td>162</td><td>76</td><td>Normal</td></tr> <tr><td>153</td><td>48</td><td>Underweight</td></tr> <tr><td>161</td><td>56</td><td>Underweight</td></tr> <tr><td>155</td><td>51</td><td>Normal</td></tr> <tr><td>172</td><td>78</td><td>Normal</td></tr> </tbody> </table>	Height	Weight	Class	157	45	Underweight	162	76	Normal	153	48	Underweight	161	56	Underweight	155	51	Normal	172	78	Normal	10	CO3																																																																					
Height	Weight	Class																																																																																											
157	45	Underweight																																																																																											
162	76	Normal																																																																																											
153	48	Underweight																																																																																											
161	56	Underweight																																																																																											
155	51	Normal																																																																																											
172	78	Normal																																																																																											
Q 10	<p>For the given confusion matrix of a binary classification problem, find out Accuracy, Precision, Recall and F-Score for the classifier.</p> <table border="1" style="margin-left: auto; margin-right: auto;"> <thead> <tr> <th>Actual Class ↓ \ Predicted Class →</th> <th>Covid 19 = Yes</th> <th>Covid 19 = No</th> </tr> </thead> <tbody> <tr> <th>Covid 19 = Yes</th> <td>110</td> <td>190</td> </tr> <tr> <th>Covid 19 = No</th> <td>240</td> <td>9460</td> </tr> </tbody> </table> <p style="text-align: center;">OR</p> <p>Write a short note on the following. Short notes should essentially include the explanations with examples</p> <p>a) Bagging c) Boosting b) Precision d) ROC</p>	Actual Class ↓ \ Predicted Class →	Covid 19 = Yes	Covid 19 = No	Covid 19 = Yes	110	190	Covid 19 = No	240	9460	10	CO4																																																																																	
Actual Class ↓ \ Predicted Class →	Covid 19 = Yes	Covid 19 = No																																																																																											
Covid 19 = Yes	110	190																																																																																											
Covid 19 = No	240	9460																																																																																											
Q11	Explain various steps in KDD process. Also comment on the statement “Data preprocessing takes 60% of the efforts in entire process”.	10	CO1																																																																																										

SECTION-C

Create a complete decision tree of the following data set using C 4.5 algorithm (based on the parameter *Gain Ratio*)

Customer ID	Gender	Car Type	Shirt Size	Class
1	M	Family	Small	C0
2	M	Sports	Medium	C0
3	M	Sports	Medium	C0
4	M	Sports	Large	C0
5	M	Sports	Extra Large	C0
6	M	Sports	Extra Large	C0
7	F	Sports	Small	C0
8	F	Sports	Small	C0
9	F	Sports	Medium	C0
10	F	Luxury	Large	C0
11	M	Family	Large	C1
12	M	Family	Extra Large	C1
13	M	Family	Medium	C1
14	M	Luxury	Extra Large	C1
15	F	Luxury	Small	C1
16	F	Luxury	Small	C1
17	F	Luxury	Medium	C1
18	F	Luxury	Medium	C1
19	F	Luxury	Medium	C1
20	F	Luxury	Large	C1

OR

Q 12

20

CO3

a) Find out all frequent item sets in the following data set using *FP Algorithm* with $\text{min_sup} = 25\%$.

b) Find all association rules with more than 80% confidence.

Transaction ID	Items Bought
1	Chips, Cookies, Regular Soda, Ham
2	Chips, Ham, Boneless Chicken, Diet Soda
3	Ham, Bacon, Whole Chicken, Regular Soda
4	Chips, Ham, Boneless Chicken, Diet Soda
5	Chips, Bacon, Boneless Chicken
6	Chips, Ham, Bacon, Whole Chicken, Regular Soda
7	Chips, Cookies, Boneless Chicken, Diet Soda

Q 12

20

CO3