

PRN No.	
---------	--

PAPER CODE	V313-232-ESB
---------------	--------------

December 2023 (ENDSEM) EXAM

TY (SEMESTER - I)

COURSE NAME: DATA SCIENCE AND MACHINE LEARNING Branch: **COMPUTER ENGINEERING** COURSE CODE: **CSUA31202**

(PATTERN 2020)

Time: [1Hr. 30 Min]

[Max. Marks: 40]

(i) Instructions to candidates:

- 1) Figures to the right indicate full marks.
- 2) Use of scientific calculator is allowed
- 3) Use suitable data wherever required
- 4) All questions are compulsory. Solve any one sub question from Question 3 and any two sub questions each from Questions 4,5 and 6 respectively.

Q. No.	Question Description	Max. Marks	CO mapped	BT Level																		
Q.1	a) Explain in detail how the model building phase is built by team in data analytics life cycle?	[2]	1	Understand																		
Q2	a) Explain categorical data and numerical data with example	[2]	2	Understand																		
Q3.	<p>a) Using k-medoid for the given dataset where the number of clusters will be 2.</p> <table><tr><td></td><td>x</td><td>y</td></tr><tr><td>0</td><td>5</td><td>6</td></tr><tr><td>1</td><td>4</td><td>5</td></tr><tr><td>2</td><td>4</td><td>6</td></tr><tr><td>3</td><td>6</td><td>7</td></tr><tr><td>4</td><td>7</td><td>8</td></tr></table> <p>I. Randomly select 2 medoids M1(4,6) and M2(6,7) and do cluster assignment.</p> <p>II. Calculate cost.</p> <p>III. Randomly select 2 medoids M1(4,5) and M2(6,7) and do cluster assignment.</p> <p>IV. Calculate cost.</p> <p>Analyze the two costs and decide whether the algorithm should keep running or stop.</p> <p>b)Examine the advantages of K medoids algorithm over K means algorithm.</p>		x	y	0	5	6	1	4	5	2	4	6	3	6	7	4	7	8	[6]	3	Analyze
	x	y																				
0	5	6																				
1	4	5																				
2	4	6																				
3	6	7																				
4	7	8																				
		[6]	3	Analyze																		

Q.4	a) Illustrate the concept of entropy and information gain? Explain how to calculate with suitable steps.	[5]	4	Apply									
	b) Executing a binary classification tree algorithm is a simple task. But how does tree split take place? Identify how does the tree determine which variable to break at the root node and which at its child nodes?	[5]	4	Apply									
	c) Demonstrate use of Decision tree and Naïve bayes classifier using real time application	[5]	4	Apply									
Q.5	a) Calculate Accuracy, Precision, Recall and F1 Score for the following Confusion Matrix on Water Shortage in Schools; Also suggest which metric would not be a good evaluation parameter here and why? <table border="1"><thead><tr><th>The Confusion Matrix</th><th>Reality: 1</th><th>Reality: 0</th></tr></thead><tbody><tr><td>Prediction: 1</td><td>75</td><td>5</td></tr><tr><td>Prediction: 0</td><td>5</td><td>15</td></tr></tbody></table> Find out Accuracy, Precision, Recall and F1 Score for the given problem.	The Confusion Matrix	Reality: 1	Reality: 0	Prediction: 1	75	5	Prediction: 0	5	15	[5]	5	Apply
	The Confusion Matrix	Reality: 1	Reality: 0										
	Prediction: 1	75	5										
Prediction: 0	5	15											
b) During the treatment of cancer patients, the doctor needs to be very careful about which patients need to be given chemotherapy. Choose correct evaluation metric among following should be used to decide whom to give chemotherapy? Justify your answer. Metrics: Precision, Recall, Accuracy, F1 score.	[5]	5	Apply										
c) Choose suitable the performance metric of a recommendation system?	[5]	5	Apply										
Q.6)	a) Evaluate the impact of outliers on numerical data visualizations. Using a line chart and a scatter plot as examples, discuss how outliers can affect the perception of trends and correlations.	[5]	6	Analyze									
	b) Estimate the use of pivot tables in handling granularity in visual representation of data.	[5]	6	Analyze									
	c) Relate the effects of poor data cleaning and wrangling on quality of visual representation.	[5]	6	Analyze									