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

| PAPER CODE | V313 - 232 (RE) |
|---|---|

### December 2023 (REEXAM)

### TY (SEMESTER - I)

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

### (PATTERN 2020)

**Time: [2 Hrs]**                                                                 **[Max. Marks: 60]**

() 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 two sub questions each from each Question 1 ,2, 3,4,5,and 6 respectively**

| Q. No. | Question Description | Max. Marks | CO mapped | BT Level |
|---|---|---|---|---|
| Q.1 | a) Describe the difference between supervised and unsupervised learning and their relevance in data science projects. | [5] | 1 | Understand |
| | b) Discuss the difference between Business Intelligence (BI) and Data Science (DS) in the context of analytics | [5] | 1 | Understand |
| | c) Explain Data analytics life cycle with all phases of life cycle for any example | [5] | 1 | Understand |
| Q2 | a) Determine the significance of discretization in machine learning. | [5] | 2 | Apply |
| | b) Illustrate the use of data reduction techniques in data preprocessing. | [5] | 2 | Apply |
| | c) Suppose you're working on a data science project of "Predictive Analysis for Student Placement in VIIT". How would you apply data cleaning and preprocessing methods to handle missing values and outliers in the Data Preparation phase? | [5] | 2 | Apply |
| Q3. | a) Examine the advantages of K medoids algorithm over K means algorithm | [5] | 3 | Analyze |
| | b) We have given a collection of 8 points. P1=[0.1,0.6] P2=[0.15,0.71],P3=[0.08,0.9],P4=[0.16,0.85],P5=[0.2,0.3] P6=[0.25,0.5] ,P7=[0.24,0.1] ,P8=[0.3,0.2]. Perform the k-mean clustering with initial centroids as m1=P1 =Cluster#1=C1 and m2=P8=cluster#2=C2. Answer the following. 1] Which cluster does P6 belong to? 2] What is the population of cluster around m2? 3] What is updated value of m1 and m2? | [5] | 3 | Analyze |
| | c) Apply k-medoid for the given dataset where the number of clusters will be 2. | [5] | 3 | Analyze |

| | x | y |
|---|---|---|
| 0 | 5 | 6 |

|  |  |  |  |  |  |  |  |
| --- | --- | --- | --- | --- | --- | --- | --- |

| 1 | 4 | 5 |
| --- | --- | --- |
| 2 | 4 | 6 |
| 3 | 6 | 7 |
| 4 | 7 | 8 |

I. Randomly select 2 medoids M1(4,6) and M2(6,7) and do cluster assignment.
II. Calculate cost.
III. Randomly select 2 medoids M1(4,5) and M2(6,7) and do cluster assignment.
IV. Calculate cost.

Analyze the two costs and decide whether the algorithm should keep running or stop

| Q.4 | a) Illustrate naive Bayes Classifier algorithm | [5] | 4 | Apply |
| | b) Illustrate how does the Binary tree classification 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? Find out Accuracy, Precision, Recall and F1 Score for the given problem. | [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. Which 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) In Spam email classifier, which of the two false predictions would you care about more:<br>i. Falsely classifying a non-spam email as spam<br>ii. Falsely claiming a spam email as non-span.<br>Identify the evaluation metric that can be used for evaluating this. | [5] | 5 | Apply |
| Q.6) | a) Compare and contrast the fundamental principles of simplicity and clarity in data visualization. How does each principle contribute to effective communication of information? Provide examples of visualizations that successfully apply each principle | [5] | 6 | Analyze |
| | b) Evaluate the strengths and weaknesses of two different types of data visualizations (e.g., bar charts vs. pie charts) for representing categorical data | [5] | 6 | Analyze |
| | c) Elaborate the effects of poor data cleaning and wrangling on quality of visual representation. | [5] | 6 | Analyze |