## UNIVERSITY OF PETROLEUM AND ENERGY STUDIES
### End Semester Examination, May 2023

**Course:** Predictive Modelling                              **Semester:** IV
**Program:** MBA (BA)                                          **Time** : 03 hrs.
**Course Code:** DSBA 8003                                     **Max. Marks: 100**

**Instructions: Attempt all sections**

### SECTION A
### 10Qx2M=20Marks

| S. No. | | Marks | CO |
|---|---|---|---|
| Q 1 | Attempt all multiple choice questions | | **CO1** |
| a. | The purpose of applying data reduction is<br>   a) to generate a larger set of variables<br>   b) to remove negative values<br>   c) to use a smaller set of variables that capture maximum information<br>   d) None of the above | **2** | **CO1** |
| b. | Which of the following is required by K-means clustering?<br>   a) defined distance metric<br>   b) number of clusters<br>   c) initial guess as to cluster centroids<br>   d) all of the mentioned | **2** | **CO1** |
| c. | Movie recommendation systems are an example of:<br>   1. Classification<br>   2. Clustering<br>   3. Reinforcement Learning<br>   4. Regression<br><br>   a) 2 only is correct<br>   b) 1 and 2 are correct<br>   c) 1 and 3 are correct<br>   d) 1, 2 and, 3 are correct<br>   e) All are correct | **2** | **CO1** |
| d. | The main benefit of standardizing a dataset is<br>   a) it makes multiple variables of a dataset come to a common scale.<br>   b) eliminates negative data values<br>   c) makes data interpretation easier. | **2** | **CO1** |
| e. | What is an outlier?<br>   a) data point most proximal to mean<br>   b) data point that falls outside the overall pattern. | **2** | **CO1** |

| | | | |
|---|---|---|---|
| | c) data point above or below 3 standard deviations of the mean. | | |
| f. | Which of the following IS NOT a component for a time series plot?<br><br>a) Seasonality<br>b) Trend<br>c) Cyclical<br>d) Noise<br>e) None of the above | 2 | CO1 |
| g. | Financial fraud detection is an example of:<br>a) Prediction problem<br>b) Clustering problem<br>c) Outlier detection problem<br>d) None of these | 2 | CO1 |
| h. | What kind of target variables are we dealing with in simple linear regression?<br>a) continous<br>b) binary<br>c) categorical | 2 | CO1 |
| i. | On what stage of data exploration are the missing values handled?<br>a) Data transformation<br>b) Data reduction<br>c) Data cleaning<br>d) All of the above | 2 | CO1 |
| j. | Which of the following is not a necessary condition for weakly stationary time series?<br>a) Mean is constant and does not depend on time<br>b) Autocovariance function depends on s and t only through their difference \|s-t\| (where t and s are moments in time)<br>c) The time series under considerations is a finite variance process<br>d) Time series is Gaussian | 2 | CO1 |
| **SECTION B**<br>**4Qx5M= 20 Marks** | | | |
| Q2. | What do you understand by data clustering? Explain the difference between different types clustering methods. | 5 | CO2 |
| Q3. | What is dimensionality reduction? Explain the difference between feature extraction and feature extraction. | 5 | CO2 |
| Q4. | What is curse of dimensionality? | 5 | CO2 |
| Q5. | Explain the difference between feature extraction and feature selection. | 5 | CO2 |

| | SECTION-C<br>3Qx10M=30 Marks | | |
|---|---|---|---|
| Q6. | Explain in detail the steps in Principle component Analysis. | **10** | **CO3** |
| Q7. | What do you understand by a time series? What is stationarity? How do you know if a given time series is stationary or not? | **10** | **CO3** |
| Q8. | A. What do you understand by CART and CHAID? What is the difference between the two?<br><br>OR<br><br>B. What is data mining? What are the different techniques used in data mining? | **10** | **CO3** |
| | SECTION-D<br>2Qx15M= 30 Marks | | |
| Q9. | Considering the following confusion matrix, define and compute the following:<br>   a) Accuracy<br>   b) Precision<br>   c) Recall<br>   d) F1 score<br>   e) Sensitivity | **15** | **CO4** |
| Q10. | A data scientist is trying to build a Breathalyzer. A breathalyzer registers someone's blood alcohol content to tell if they are "over the limit" or "under the influence" of alcohol. It is typically used at roadside police stops to determine if someone is legally able to drive. Answer the following questions:<br><br>   a. What is the positive class?<br>   b. What would a recall of 70% mean?<br>   c. What would a precision of 90% mean?<br>   d. If the recall is 80% and the precision is 75%, what is the FPR? | **15** | **CO4** |

Confusion matrix for Q9:

| N=165 | | Predicted | | |
|---|---|---|---|---|
| | | No | Yes | |
| Actual | No | 50 | 10 | 60 |
| | Yes | 5 | 100 | 105 |
| | | 55 | 110 | |