## UNIVERSITY OF PETROLEUM AND ENERGY STUDIES

### End Semester Examination, May 2021

**Course:  Machine Learning**
**Program: M. Tech - CSE**
**Course Code:  CSAI7007P**

**Semester:  II**
**Time 03 hrs.**
**Max. Marks: 100**

### SECTION A

1. Each Question will carry 5 Marks
2. Instruction: Multiple choice Questions.

| S. No. | | Marks | CO |
| --- | --- | --- | --- |
| Q 1 | Imagine you are working on a project which is a binary classification problem. You trained a model on training dataset and get the below confusion matrix on validation dataset.<br><br>| n=165 | Predicted: NO | Predicted: YES |<br>\| Actual: NO \| 50 \| 10 \|<br>\| Actual: YES \| 5 \| 100 \|<br><br>Based on the above confusion matrix, choose which option(s) below will give you correct predictions?<br>  1.  Accuracy is ~0.91<br>  2.  Misclassification rate is ~ 0.91<br>  3.  False positive rate is ~0.95<br>  4.  True positive rate is ~0.95<br><br>A) 1 and 3<br>B) 2 and 4<br>C) 1 and 4<br>D) 2 and 3 | 5 | CO1 |
| Q2 | Suppose we have a dataset which can be trained with 100% accuracy with help of a decision tree of depth 6. Now consider the points below and choose the option based on these points.<br>Note: All other hyper parameters are same and other factors are not affected.<br>  1.  Depth 4 will have high bias and low variance<br>  2.  Depth 4 will have low bias and low variance<br><br>A) Only 1<br>B) Only 2<br>C) Both 1 and 2<br>D) None of the above | 5 | CO2 |

| Q3 | Which of the following options is/are true for K-fold cross-validation? <br>    1. Increase in K will result in higher time required to cross validate the result. <br>    2. Higher values of K will result in higher confidence on the cross-validation result as compared to lower value of K. <br>    3. If K=N, then it is called Leave one out cross validation, where N is the number of observations. <br><br> A) 1 and 2 <br> B) 2 and 3 <br> C) 1 and 3 <br> D) 1,2 and 3 | 5 | CO3 |
|---|---|---|---|
| Q4 | Imagine, you are solving a classification problems with highly imbalanced class. The majority class is observed 99% of times in the training data. <br> Your model has 99% accuracy after taking the predictions on test data. Which of the following is true in such a case? <br>    1. Accuracy metric is not a good idea for imbalanced class problems. <br>    2. Accuracy metric is a good idea for imbalanced class problems. <br>    3. Precision and recall metrics are good for imbalanced class problems. <br>    4. Precision and recall metrics aren't good for imbalanced class problems. <br><br> A) 1 and 3 <br> B) 1 and 4 <br> C) 2 and 3 <br> D) 2 and 4 | 5 | CO4 |
| Q5 | For which of the following hyperparameters, higher value is better for decision tree algorithm? <br>    1. Number of samples used for split <br>    2. Depth of tree <br>    3. Samples for leaf <br><br> A)1 and 2 <br> B) 2 and 3 <br> C) 1 and 3 <br> D) 1, 2 and 3 <br> E) Can't say | 5 | CO1 |
| Q6 | Which of the following options can be used to get global minima in k-Means Algorithm? <br>    1. Try to run algorithm for different centroid initialization <br>    2. Adjust number of iterations <br>    3. Find out the optimal number of clusters <br><br> A) 2 and 3 <br> B) 1 and 3 <br> C) 1 and 2 <br> D) All of above | 5 | CO2 |

## SECTION B

1. Each question will carry 10 marks
2. Instruction: Write short / brief notes.

| | | | |
|---|---|---|---|
| **Q 7** | a) Why is Nave Bayes classifier so powerful for text classification?<br>b) Why Normalization is required in machine learning? | **6+4** | **CO1** |
| **Q8** | a) In which algorithm, Ginni index is used. Explain the algorithm in detail with suitable example.<br>b) Why does the decision tree suffer often with overfitting problem? | **6+4** | **CO3** |
| **Q9** | a) What is the goal of SVM? How to select the margin?<br>b) Given the following data for the sales (in million dollars) of Car of an Automobile Company for 6 consecutive years.<br><br>Year \| 2013 \| 2014 \| 2015 \| 2016 \| 2017 \| 2018<br>Sales \| 110 \| 100 \| 250 \| 275 \| 230 \| 300<br><br>Based on the above data, predict the sales for next three consecutive years.<br><br>**OR**<br><br>A data set is given to you about utilities fraud detection. You have built a classifier model and achieved a performance score of 98.5%. Is this a good model? If yes, justify. If not, what can you do about it? | **6+4** | **CO3** |
| **Q10** | a) Which algorithm can be used to fit the data over a linear line? Is that algorithm supervised or unsupervised? And how would you calculate the cost for that algorithm?<br>b) Which is more important to you- model accuracy or model performance? Support with suitable example. | **6+4** | **CO2** |
| **Q11** | a) How could you divide the 'training Set' and 'test Set' in a Machine Learning Model? How much data will you allocate for training, validation, and test Sets?<br>b) Explain why k-fold cross validation does not work well with time series model? What can you do about it? | **6+4** | **CO2** |

The Q9 table reproduced properly:

| Year | 2013 | 2014 | 2015 | 2016 | 2017 | 2018 |
|------|------|------|------|------|------|------|
| Sales | 110 | 100 | 250 | 275 | 230 | 300 |

## SECTION-C

Note: Attempt any one question from two options.

| | | | |
|---|---|---|---|
| **Q12** | Differentiate between<br><br>a) Supervised, unsupervised and reinforcement learning<br>b) Bagging and boosting.<br>c) Linear Regression and Logistic Regression<br>d) Overfitting and under fitting<br><br>**OR**<br><br>Consider a medical diagnosis problem in which there are two alternative hypotheses: | **10+10** | **CO4** |

| | (1) that the patient has a particular form of COVID19 (+) and<br>(2) That the patient does not (-).<br>A patient takes a lab test and the result comes back positive. The test returns a correct positive result in only 98% of the cases in which the disease is actually present, and a correct negative result in only 97% of the cases in which the disease is not present. Furthermore, .008 of the entire population have this COVID19. Determine whether the patient has COVID19 or not using the MAP hypothesis. | | |