# B.M.S. College of Engineering, Bengaluru-560019

**Autonomous Institute Affiliated to VTU**

## June / July 2025 Semester End Main Examinations

**Programme: B.E.**          **Semester: V**

**Branch: Artificial Intelligence and Machine Learning**      **Duration: 3 hrs.**

**Course Code: 24AM5PCSML**          **Max Marks: 100**

**Course:  STATISTICAL MODELING**

**Instructions**:  1. Answer any FIVE full questions, choosing one full question from each unit.
2. Missing data, if any, may be suitably assumed.

*Important Note:* Completing your answers, compulsorily draw diagonal cross lines on the remaining blank pages. Revealing of identification, appeal to evaluator will be treated as malpractice.

| | | | | CO | PO | Marks |
|---|---|---|---|---|---|---|
| | | | **UNIT - I** | *CO* | *PO* | **Marks** |
| 1 | a) | | By stating, the simple linear regression model and its assumptions, Obtain the least square estimates of the parameters. | *CO1* | *PO1* | **08** |
| | b) | | The relationship between expenditure ($) and income ($) is modelled using simple linear regression, with the following regression output, and its output is given below: | *CO1* | *PO2* | **06** |

| Coefficients: | Estimate | Std. Error |
|---|---|---|
| Intercept | 31.42 | 0.0565 |
| income ($) | 0.7138 | 0.1209 |

Obtain the 95% confidence interval for intercept ($\beta_0$) and slope ($\beta_1$) parameter.

| | | | | CO | PO | Marks |
|---|---|---|---|---|---|---|
| | c) | | Define R² and discuss its limitations, especially in the context of comparing models with different numbers of predictors. What advantages does adjust. R² offer in balancing model efficiency and complexity? | *CO1* | *PO2* | **06** |
| | | | **OR** | | | |
| 2 | a) | | As part of a waste removal project, a new compression machine for processing sewage sludge is being studied. Engineers are interested in the following variables. y - moisture control of compressed pellets (%) and x-machine filtration rate (kg-DS/m/hr). Engineers collect observations of (x, y) with a random sample of size 10 sewage specimens. | *CO1* | *PO2* | **08** |

| x | 125.3 | 98.2 | 161.2 | 178.5 | 165.3 | 159.5 | 145.8 | 159.6 | 110.7 |
|---|---|---|---|---|---|---|---|---|---|
| y | 77.9 | 76.8 | 80.1 | 80.2 | 79.0 | 79.9 | 79.0 | 79.0 | 78.6 |

Obtain the fitted values of y using simple regression model and also find the value of y, when x = 180.

| | | | | CO2 | PO2 | 06 |
|---|---|---|---|---|---|---|
| | | b) | Summarize the test procedure to test the significance of the simple linear regression model. | CO2 | PO2 | **06** |
| | | c) | Provide a brief overview of the following topics:<br>   i)     the applications of regression analysis, and<br>   ii)    how simple regression analysis is used for predicting new observations. | CO1 | PO2 | **06** |

### UNIT - II

| | | | | | | |
|---|---|---|---|---|---|---|
| 3 | a) | Provide a brief overview of Q-Q plots, including their purpose and how they are used to assess the normality of data. | CO2 | PO2 | **06** |

| | | | | | |
|---|---|---|---|---|---|
| | b) | Consider the following dataset related to the performance of sorting algorithms: | CO1 | PO2 | **07** |

| Execution Time in milliseconds (y) | Number of Elements (X1) | Input Complexity (X2) |
|---|---|---|
| 78.5 | 7 | 26 |
| 74.3 | 1 | 29 |
| 104.3 | 11 | 56 |
| 87.6 | 11 | 31 |
| 95.9 | 7 | 52 |
| 109.2 | 11 | 55 |

Provide the least squares estimates for the regression coefficients and summarize the fitted regression model.

| | | | | | |
|---|---|---|---|---|---|
| | c) | When the problem of heteroscedasticity exists in the data? Explain the role of residual graphs to detect the same. | CO2 | PO1 | **07** |

### OR

| | | | | | | |
|---|---|---|---|---|---|---|
| 4 | a) | Give a brief note on sources and detection of multicollinearity in multiple regression analysis. | CO2 | PO1 | **06** |
| | b) | Given the following dataset: | CO2 | PO2 | **06** |

| Observation | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| Rainfall (cms) | 30 | 23 | 34 | 31 | 17 | 36 |
| Yield (tons) | 65 | 62 | 70 | 64 | 52 | 68 |

Calculate the Durbin-Watson $d$ statistic to test positive autocorrelation and conclude. ($d_L = 0.61$ and $d_U = 1.40$)

| | | | | | |
|---|---|---|---|---|---|
| | c) | Compare and contrast the regularization techniques of Ridge regression and LASSO regression. What are the key differences in their methodologies and effects on model performance? | CO2 | PO2 | **08** |

### UNIT - III

| | | | | | | |
|---|---|---|---|---|---|---|
| 5 | a) | "Not all leverage points are outliers" – Justify the statements with an example. | CO1 | PO2 | **06** |
| | b) | How the DFFITs and DFBETAs are helpful to identify the influential observations. | CO2 | PO1 | **06** |

| | | c) | Given the table below with columns for Model No, p (number of parameters including the intercept), R-Squared ($R^2$), Adjusted R-Squared (Adj $R^2$), Standard Error (S), and Mallows' Cp: | | | **08** |
|---|---|---|---|---|---|---|

| Model | p | $R^2$ | Adj $R^2$ | S | Mallows' Cp |
|---|---|---|---|---|---|
| 1 | 1 | 0.64 | 0.62 | 0.50 | 9.4 |
| 2 | 1 | 0.62 | 0.60 | 0.52 | 10.7 |
| 3 | 2 | 0.80 | 0.77 | 0.39 | 1.5 |
| 4 | 2 | 0.79 | 0.76 | 0.40 | 3.3 |
| 5 | 3 | 0.81 | 0.76 | 0.40 | 3.2 |
| 6 | 3 | 0.80 | 0.76 | 0.40 | 3.4 |
| 7 | 4 | 0.81 | 0.74 | 0.41 | 5.0 |

Select the model that best fits the data and justify your answer.
Justify why other models are not suitable.

**OR**

| 6 | a) | Explain added-variable plots. | CO2 | PO1 | **06** |
|---|---|---|---|---|---|
| | b) | Describe the Cook's distance to detect the influential observations. | CO2 | PO1 | **06** |
| | c) | Define Mallow's $C_p$ and select the best model for the below given output. | CO2 | PO1 | **08** |

| Predictor Variables | P+1 | Mallows' Cp |
|---|---|---|
| Hours | 2 | 45.5 |
| Prep exams | 2 | 31.4 |
| GPA | 2 | 29.3 |
| Hours, Prep exams | 3 | 3.4 |
| Hours, GPA | 3 | 2.9 |
| Prep exams, GPA | 3 | 2.7 |
| Hours, Prep exams, GPA | 4 | 4 |

**UNIT - IV**

| 7 | a) | Identify the type of categorical scale (nominal or ordinal) used in each of the following scenarios and explain your reasoning:<br>1. **Customer satisfaction ratings**: "Very Unsatisfied," "Unsatisfied," "Neutral," "Satisfied," "Very Satisfied."<br>2. **Type of pet owned**: "Dog," "Cat," "Bird," "Fish," "Other."<br>3. **Highest level of education:** "High School," "Associate Degree," "Bachelor's," "Master's," "Doctorate." | CO1 | PO2 | **06** |
|---|---|---|---|---|---|
| | b) | Write down the test procedure to test the independence of attributes. | CO2 | PO1 | **06** |
| | c) | Elucidate the following.<br>i. Sensitivity<br>ii. Specificity<br>iii. ROC curve | CO3 | PO1 | **08** |

**OR**

| | 8 | a) | In a study of 263 adolescents evaluated for suicidal behaviour, 186 were classified as non-suicidal (NS) at a six-month follow-up. Among them, 86 were assessed as having depression at baseline. Of the 77 adolescents with persistent suicidal behaviour (SB) at follow-up, 45 had been assessed for depression at baseline.<br>   i.      Construct a contingency table based on this data.<br>   ii.     Calculate the odds ratio and provide an interpretation. | *CO3* | *PO2* | **08** |
| | | b) | The following is the data regarding family condition and examination result of 100 students test whether family conditions and results are independent (critical value is 6.63). | *CO3* | *PO3* | **06** |

| Family condition | Examination results | |
|---|---|---|
| | Pass | Fail |
| Good | 30 | 10 |
| Bad | 20 | 40 |

| | | c) | Explain how logistic regression is used for binary classification and discuss its advantages. Additionally, mention other machine learning methods that can also be applied for binary classification tasks. | *CO1* | *PO1* | **06** |
| | | | **UNIT - V** | | | |
| | 9 | a) | Discuss the different types of Markov models and their applications in various fields. | *CO3* | *PO1* | **06** |
| | | b) | Explain the forward and backward algorithm of hidden Markov models. | *CO3* | *PO1* | **08** |
| | | c) | Describe the smoothing and filtering techniques of classification using hidden Markov models. | *CO3* | *PO2* | **06** |
| | | | **OR** | | | |
| | 10 | a) | State the procedure of choosing number of hidden states in hidden Markov models. | *CO3* | *PO2* | **08** |
| | | b) | Provide pseudocode for the working of Viterbi Algorithm. | *CO3* | *PO2* | **06** |
| | | c) | Discuss the Gaussian mixture models with hidden Markov models. | *CO3* | *PO2* | **06** |

**\*\*\*\*\*\***