

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

Autonomous Institute Affiliated to VTU

## January / February 2025 Semester End Main Examinations

**Programme: B.E.**

**Semester: V**

**Branch: Artificial Intelligence and Machine Learning**

**Duration: 3 hrs.**

**Course Code: 24AM5PCSL**

**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.

<b>UNIT - I</b>			<b>CO</b>	<b>PO</b>	<b>Marks</b>												
1	a)	Define Simple Linear Regression and explain how it models relationships between variables, including a brief example.	<i>CO1</i>	<i>PO1</i>	<b>06</b>												
	b)	Derive the least squares estimates for the parameters of a simple linear regression model.	<i>CO1</i>	<i>PO2</i>	<b>07</b>												
	c)	Explain the steps involved in testing the hypothesis for the significance of slope parameter in a simple linear regression model.	<i>CO2</i>	<i>PO2</i>	<b>07</b>												
<b>OR</b>																	
2	a)	What is the coefficient of determination ( $R^2$ )? How is it calculated, and what does it signify in a regression model?	<i>CO1</i>	<i>PO2</i>	<b>06</b>												
	b)	<p>The data regarding sales and advertisement expenditure of a firm is as follows:</p> <table border="1" style="margin-left: auto; margin-right: auto;"> <thead> <tr> <th><b>Measure</b></th> <th><b>Sales</b> (in crores)</th> <th><b>Advertisement expenditure</b> (in crores)</th> </tr> </thead> <tbody> <tr> <td>Means</td> <td>40</td> <td>6</td> </tr> <tr> <td>Standard deviations</td> <td>10</td> <td>1.5</td> </tr> <tr> <td>Correlation coefficient</td> <td colspan="2">0.9</td> </tr> </tbody> </table> <p>If the firm targets sales of 60 crores, what should be the required advertisement expenditure? Use linear regression to calculate and explain.</p>	<b>Measure</b>	<b>Sales</b> (in crores)	<b>Advertisement expenditure</b> (in crores)	Means	40	6	Standard deviations	10	1.5	Correlation coefficient	0.9		<i>CO1</i>	<i>PO2</i>	<b>07</b>
<b>Measure</b>	<b>Sales</b> (in crores)	<b>Advertisement expenditure</b> (in crores)															
Means	40	6															
Standard deviations	10	1.5															
Correlation coefficient	0.9																
	c)	Perform ANOVA for a simple linear regression model. Explain the decomposition of the total sum of squares and its relation to the regression and error sum of squares.	<i>CO2</i>	<i>PO2</i>	<b>07</b>												

**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.

		UNIT - II																										
3	a)	List and explain the assumptions of the Multiple Linear Regression model. Why are these assumptions important?			CO1	PO1	<b>06</b>																					
	b)	Derive the least squares estimates for the parameters of a Multiple Linear Regression model and obtain the same for the following:			CO1	PO2	<b>10</b>																					
		<table border="1"> <thead> <tr> <th>Execution Time in milliseconds (y)</th> <th>Number of Elements (X1)</th> <th>Input Complexity (X2)</th> </tr> </thead> <tbody> <tr><td>78.5</td><td>7</td><td>26</td></tr> <tr><td>74.3</td><td>1</td><td>29</td></tr> <tr><td>104.3</td><td>11</td><td>56</td></tr> <tr><td>87.6</td><td>11</td><td>31</td></tr> <tr><td>95.9</td><td>7</td><td>52</td></tr> <tr><td>109.2</td><td>11</td><td>55</td></tr> </tbody> </table>			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			
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																										
	c)	State Gauss-Markov theorem.			CO1	PO1	<b>04</b>																					
		<b>OR</b>																										
4	a)	What is multicollinearity in a regression model? Explain how the Variance Inflation Factor (VIF) is used to detect multicollinearity. Suggest remedies to address multicollinearity.			CO2	PO2	<b>07</b>																					
	b)	Define heteroscedasticity in regression analysis. How can it be detected using residual plots? Suggest ways to address heteroscedasticity.			CO2	PO2	<b>07</b>																					
	c)	Given the following dataset: <table border="1"> <thead> <tr> <th>Observation</th> <th>1</th> <th>2</th> <th>3</th> <th>4</th> <th>5</th> <th>6</th> </tr> </thead> <tbody> <tr><td>Rainfall (cms)</td><td>30</td><td>23</td><td>34</td><td>31</td><td>17</td><td>36</td></tr> <tr><td>Yield (tons)</td><td>65</td><td>62</td><td>70</td><td>64</td><td>52</td><td>68</td></tr> </tbody> </table> Calculate the Durbin-Watson $d$ statistic to test positive autocorrelation and conclude. ( $d_L = 0.61$ and $d_U = 1.40$ )			Observation	1	2	3	4	5	6	Rainfall (cms)	30	23	34	31	17	36	Yield (tons)	65	62	70	64	52	68	CO2	PO2	<b>06</b>
Observation	1	2	3	4	5	6																						
Rainfall (cms)	30	23	34	31	17	36																						
Yield (tons)	65	62	70	64	52	68																						
		<b>UNIT - III</b>																										
5	a)	Explain the importance of model diagnostics in regression analysis. How are added variable plots used to diagnose model issues?			CO2	PO1	<b>06</b>																					
	b)	What are Hat Matrix Leverage values? How are they used to identify leverage data points in regression analysis?			CO2	PO2	<b>07</b>																					
	c)	The DFBETAS matrix for a regression model with $n=10$ observations and 2 predictors are given below:			CO1	PO3	<b>07</b>																					

		$\begin{bmatrix} 0.02 & -0.15 \\ -0.08 & 0.10 \\ 0.05 & -0.02 \\ -0.30 & 0.25 \\ 0.12 & -0.05 \\ 0.07 & 0.04 \\ -0.50 & 0.60 \\ 0.20 & -0.10 \\ -0.05 & 0.08 \\ 0.03 & -0.04 \end{bmatrix}$ <p>Compute the threshold for identifying influential observations and determine which data points are influential for each predictor.</p>																												
		<b>OR</b>																												
6	a)	Discuss the use of Studentized Deleted Residuals in identifying outliers. How do they differ from regular residuals?	CO2	PO2	<b>06</b>																									
	b)	Explain model validation criteria such as AIC, BIC, and Mallows' Cp. How are they used in model selection?	CO3	PO1	<b>06</b>																									
	c)	<p>For the given dataset:</p> <table border="1" style="display: inline-table; vertical-align: middle;"> <tr><td><b>i</b></td><td><b>1</b></td><td><b>2</b></td><td><b>3</b></td><td><b>4</b></td></tr> <tr><td><math>y_i</math></td><td>301</td><td>327</td><td>246</td><td>187</td></tr> <tr><td><math>x_{i1}</math></td><td>14</td><td>19</td><td>12</td><td>11</td></tr> <tr><td><math>x_{i2}</math></td><td>25</td><td>32</td><td>22</td><td>15</td></tr> <tr><td><math>h_{ii}</math></td><td>0.28</td><td>0.33</td><td>0.85</td><td>0.68</td></tr> </table> <p>The fitted regression model and the error mean sum of square are given as  <math>\hat{y} = 80.93 - 5.84X_1 + 11.32X_2</math> and <math>MSE = 574.9</math>.</p> <p>i. Compute the fitted values and residuals  ii. Identify the outliers in Y by computation of studentized residuals</p>	<b>i</b>	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	$y_i$	301	327	246	187	$x_{i1}$	14	19	12	11	$x_{i2}$	25	32	22	15	$h_{ii}$	0.28	0.33	0.85	0.68	CO3	PO3	<b>08</b>
<b>i</b>	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>																										
$y_i$	301	327	246	187																										
$x_{i1}$	14	19	12	11																										
$x_{i2}$	25	32	22	15																										
$h_{ii}$	0.28	0.33	0.85	0.68																										
		<b>UNIT - IV</b>																												
7	a)	Define a contingency table and explain how it is used in categorical data analysis.	CO1	PO1	<b>06</b>																									
	b)	Explain the concept of an odds ratio. How is it calculated and interpreted?	CO1	PO2	<b>07</b>																									
	c)	Describe the chi-squared test for independence. Provide the formula and explain its application.	CO2	PO1	<b>07</b>																									
		<b>OR</b>																												
8	a)	Explain the logistic regression model and its importance in analyzing categorical data.	CO1	PO2	<b>06</b>																									
	b)	Discuss how sensitivity, specificity, and ROC curves are used to summarize the predictive power of a logistic regression model.	CO3	PO1	<b>07</b>																									
	c)	How are parameters interpreted in logistic regression? Provide an example.	CO3	PO1	<b>07</b>																									

		UNIT - V			
	9	a) Let $\{X_n, n \geq 0\}$ be a discrete time Markov Chain with state space $\{1,2,3,4\}$ and the transition probability matrix A given as: $A = \begin{pmatrix} 0.4 & 0.3 & 0.3 \\ 0.5 & 0.0 & 0.5 \\ 0.1 & 0.2 & 0.7 \end{pmatrix}$ Given the initial probabilities $\pi = (0.25 \quad 0.15 \quad 0.6)$ Compute the following: i. $P(X_2 = 2   X_0 = 1)$ ii. $P(X_3 = 3, X_2 = 2)$ iii. $P(X_2 = 2, X_1 = 1   X_0 = 1)$	CO3	PO2	<b>06</b>
		b) Describe the Forward-Backward algorithm and its applications in Hidden Markov Models.	CO3	PO2	<b>07</b>
		c) Discuss the importance of Gaussian mixture models with Hidden Markov Models and their applications.	CO3	PO2	<b>07</b>
		<b>OR</b>			
	10	a) Given a Hidden Markov Model (HMM) with the following parameters: Transition Probability Matrix (A), Emission Probability matrix (B) and initial probabilities ( $\pi$ ) as: $p \quad q \quad 0 \quad 1$ $A = \begin{bmatrix} p & 0.7 & 0.3 \\ q & 0.4 & 0.6 \end{bmatrix}, B = \begin{bmatrix} 0.5 & 0.5 \\ 0.1 & 0.9 \end{bmatrix} \text{ and } \pi = [0.6 \quad 0.4]$ Compute the likelihood probability of the sequence $\{0, 1, 0\}$ using backward probabilities obtained from the backward Algorithm.	CO3	PO2	<b>06</b>
		b) Compare generative and discriminative classifiers. Provide examples where each is preferred.	CO3	PO1	<b>07</b>
		c) Discuss the challenges in choosing the number of hidden states for an HMM and describe how smoothing and filtering techniques are applied.	CO3	PO2	<b>07</b>

\*\*\*\*\*