

Bansilal Ramnath Agarwal Charitable Trust's  
**VISHWAKARMA INSTITUTE OF TECHNOLOGY, PUNE – 411037.**  
 (An Autonomous Institute Affiliated to Savitribai Phule Pune University)  
**Examination: ESE**

Year: S.Y. Common

Branch:

Subject: Data science

Subject Code: MD 2201

Max. Marks:60

Total Pages of Question Paper: 1

Day &amp; Date: Wed : 22/11/23

Time: 10.30 am -12.30 pm

### Instructions to Candidate

1. All questions are compulsory.
2. Neat diagrams must be drawn wherever necessary.
3. Figures to the right indicate full marks.

Q.No.	CO No	BT No		Max marks																																										
Q.1.	1	1	<p>Part of an annual transparency report published by a leading multinational technology company is as shown below –</p> <table border="1"> <thead> <tr> <th>Country</th> <th>CR_req</th> <th>CR_compl in %</th> <th>UD_req</th> <th>UD_compl in %</th> <th>Hemi</th> <th>HDI</th> </tr> </thead> <tbody> <tr> <td>Austria</td> <td>21</td> <td>100</td> <td>134</td> <td>32</td> <td>Southern</td> <td>High</td> </tr> <tr> <td>Belgium</td> <td>10</td> <td>33</td> <td>361</td> <td>73</td> <td>Northern</td> <td>High</td> </tr> <tr> <td>Brazil</td> <td>224</td> <td>67</td> <td>703</td> <td>82</td> <td>Southern</td> <td>Medium</td> </tr> <tr> <td>Somalia</td> <td>104</td> <td>31</td> <td>227</td> <td>61</td> <td>Southern</td> <td>Poor</td> </tr> <tr> <td>USA</td> <td>92</td> <td>63</td> <td>5950</td> <td>93</td> <td>Northern</td> <td>High</td> </tr> </tbody> </table> <p>Where, the variables are as follows -            CR_req: Content removal requests            CR_comp: Content removal compliance in %            UD_req: User Data requests            UD_compl: User Data compliance in %            Hemi: Hemisphere            HDI: Human Development Index</p> <p>Identify each variables as Discrete Numerical, Continuous Numerical, Ordinal Categorical or Regular Categorical with justification</p>	Country	CR_req	CR_compl in %	UD_req	UD_compl in %	Hemi	HDI	Austria	21	100	134	32	Southern	High	Belgium	10	33	361	73	Northern	High	Brazil	224	67	703	82	Southern	Medium	Somalia	104	31	227	61	Southern	Poor	USA	92	63	5950	93	Northern	High	12
Country	CR_req	CR_compl in %	UD_req	UD_compl in %	Hemi	HDI																																								
Austria	21	100	134	32	Southern	High																																								
Belgium	10	33	361	73	Northern	High																																								
Brazil	224	67	703	82	Southern	Medium																																								
Somalia	104	31	227	61	Southern	Poor																																								
USA	92	63	5950	93	Northern	High																																								
Q. 2. (A)	2	2	What are type-I and type-II errors in hypothesis testing? Which error should be minimized while giving a judgement in the court of law? Justify	4																																										
(B)	2	2	A sample of 50 news paper readers were asked about the total hours they spend per week in reading the newspaper. The group in the sample had an average of 3.2 hours with a standard deviation of 1.74. Calculate the 95% confidence interval range, based on this data.	6																																										
Q. 3. (A)	3	1	Calculate the distance between points A (2,7,4) and B (3,2,4,8,5,8) using i. Manhattan Distance metric and ii. Euclidean Distance metric	4																																										
(B)	3	2	For a given univariate function $f(x) = 2x^4 - 8x^3 - 112x^2 + 1717$ find out the optimal local minimum and global minimum	4																																										
Q. 4. (A)	4	4	As an outcome of a linear regression process, following performance values are obtained. SSR also known as sum of squares due to regression = 92.48; SSE also known as sum of squares due to error = 12.87. Calculate the value of $R^2$ and comment if this regression fit is a good fit or not.	4																																										

(B)	4	1	<p>A group of 10 customers, having yearly savings between 0 to 6 lacs are considered in the following logistic regression example. Loan status as 0 indicates a loan defaulter and 1 indicates non-defaulter.</p> <table><thead><tr><th>Amount in Savings (in lacs)</th><th>Loan Status</th></tr></thead><tbody><tr><td>1.00</td><td>0</td></tr><tr><td>1.25</td><td>0</td></tr><tr><td>1.5</td><td>0</td></tr><tr><td>1.8</td><td>1</td></tr><tr><td>2.25</td><td>0</td></tr><tr><td>2.4</td><td>0</td></tr><tr><td>3.4</td><td>0</td></tr><tr><td>4.6</td><td>0</td></tr><tr><td>5.2</td><td>1</td></tr><tr><td>5.8</td><td>1</td></tr></tbody></table> <p>The logistic regression function after Maximum Likelihood estimation is given as <math>\Pi = (1 / 1 + e^{-(4.07778 + 1.5046 * Savings)})</math>.</p> <p>Answer the following –</p> <ol style="list-style-type: none"><li>1. If a loan applicant with annual savings of Rs. 2.5 Lacs approaches the bank, should the application be processed for loan disbursement or rejected?</li><li>2. Based on the predictions made using logistic regression for all 10 data, how many times the predictions would match with the actual loan status? What will be the classification accuracy of the logistic regression as a classifier?</li></ol>	Amount in Savings (in lacs)	Loan Status	1.00	0	1.25	0	1.5	0	1.8	1	2.25	0	2.4	0	3.4	0	4.6	0	5.2	1	5.8	1	8																																																																										
Amount in Savings (in lacs)	Loan Status																																																																																																			
1.00	0																																																																																																			
1.25	0																																																																																																			
1.5	0																																																																																																			
1.8	1																																																																																																			
2.25	0																																																																																																			
2.4	0																																																																																																			
3.4	0																																																																																																			
4.6	0																																																																																																			
5.2	1																																																																																																			
5.8	1																																																																																																			
Q. 5. (A)	5	4	<p>The training data for a supervised classification is as follows – <math>X_1 = (1.8, 1.6, 1)</math>, <math>X_2 = (2.1, 8, 1)</math>, <math>X_3 = (3.2, 2.4, 1)</math>, <math>X_4 = (2.4, 2.6, 1)</math>, <math>X_5 = (6.5, 4.2, 2)</math>, <math>X_6 = (7.3, 4.3, 2)</math>, <math>X_7 = (6.5, 4.2, 2)</math>, <math>X_8 = (7.0, 4.8, 2)</math>. The test datapoint is at (4.6, 3.2). Use</p> <ol style="list-style-type: none"><li>1. Nearest Neighbor assign appropriate class to the test point</li><li>2. K - Nearest Neighbor with <math>k = 3</math> to assign appropriate class to the test point</li><li>3. Weighted / Modified K - Nearest Neighbor with <math>k = 3</math> to assign appropriate class to the test point</li></ol>	6																																																																																																
Q. 5. (B)	5	4	<p>Consider 50 patterns that are split in 3 classes with 12, 28 and 10 samples respectively. Calculate the Entropy impurity, Gini impurity and Misclassification impurity at the node.</p>	4																																																																																																
OR alternate option for Q.5 (A) AND Q. 5 (B) together as Q. 5 (C)																																																																																																				
Q. 5 (C)	5	4	<p>The table below gives the following training data –</p> <table><thead><tr><th>Sr. No.</th><th>Worker</th><th>Mood</th><th>Job</th><th>Time</th><th>Good Work Quality?</th></tr></thead><tbody><tr><td>1</td><td>Sam</td><td>Bad</td><td>Painting</td><td>Morning</td><td>Yes</td></tr><tr><td>2</td><td>Sam</td><td>Good</td><td>Plumbing</td><td>Evening</td><td>Yes</td></tr><tr><td>3</td><td>Ashwin</td><td>Bad</td><td>Painting</td><td>Morning</td><td>No</td></tr><tr><td>4</td><td>Ashwin</td><td>Bad</td><td>Plumbing</td><td>Evening</td><td>No</td></tr><tr><td>5</td><td>Ashwin</td><td>Good</td><td>Washing</td><td>Morning</td><td>Yes</td></tr><tr><td>6</td><td>Sam</td><td>Good</td><td>Washing</td><td>Evening</td><td>Yes</td></tr><tr><td>7</td><td>Sham</td><td>Good</td><td>Painting</td><td>Morning</td><td>Yes</td></tr><tr><td>8</td><td>Sham</td><td>Bad</td><td>Plumbing</td><td>Evening</td><td>No</td></tr><tr><td>9</td><td>Ashwin</td><td>Bad</td><td>Washing</td><td>Morning</td><td>No</td></tr><tr><td>10</td><td>Ashwin</td><td>Good</td><td>Washing</td><td>Evening</td><td>Yes</td></tr><tr><td>11</td><td>Sam</td><td>Good</td><td>Painting</td><td>Evening</td><td>Yes</td></tr><tr><td>12</td><td>Sham</td><td>Bad</td><td>Washing</td><td>Morning</td><td>No</td></tr><tr><td>13</td><td>Sham</td><td>Good</td><td>Washing</td><td>Evening</td><td>Yes</td></tr><tr><td>14</td><td>Sam</td><td>Bad</td><td>Plumbing</td><td>Morning</td><td>Yes</td></tr><tr><td>15</td><td>Sham</td><td>Good</td><td>Painting</td><td>Morning</td><td>No</td></tr></tbody></table> <p>Using Naïve Bays Classification, estimate the class for 'Good Work Quality' for the given feature vector <math>X = \{ \text{Worker} = \text{Sham}, \text{Mood} = \text{Bad}, \text{Job} = \text{Painting}, \text{Time} = \text{Morning} \}</math></p>	Sr. No.	Worker	Mood	Job	Time	Good Work Quality?	1	Sam	Bad	Painting	Morning	Yes	2	Sam	Good	Plumbing	Evening	Yes	3	Ashwin	Bad	Painting	Morning	No	4	Ashwin	Bad	Plumbing	Evening	No	5	Ashwin	Good	Washing	Morning	Yes	6	Sam	Good	Washing	Evening	Yes	7	Sham	Good	Painting	Morning	Yes	8	Sham	Bad	Plumbing	Evening	No	9	Ashwin	Bad	Washing	Morning	No	10	Ashwin	Good	Washing	Evening	Yes	11	Sam	Good	Painting	Evening	Yes	12	Sham	Bad	Washing	Morning	No	13	Sham	Good	Washing	Evening	Yes	14	Sam	Bad	Plumbing	Morning	Yes	15	Sham	Good	Painting	Morning	No	10
Sr. No.	Worker	Mood	Job	Time	Good Work Quality?																																																																																															
1	Sam	Bad	Painting	Morning	Yes																																																																																															
2	Sam	Good	Plumbing	Evening	Yes																																																																																															
3	Ashwin	Bad	Painting	Morning	No																																																																																															
4	Ashwin	Bad	Plumbing	Evening	No																																																																																															
5	Ashwin	Good	Washing	Morning	Yes																																																																																															
6	Sam	Good	Washing	Evening	Yes																																																																																															
7	Sham	Good	Painting	Morning	Yes																																																																																															
8	Sham	Bad	Plumbing	Evening	No																																																																																															
9	Ashwin	Bad	Washing	Morning	No																																																																																															
10	Ashwin	Good	Washing	Evening	Yes																																																																																															
11	Sam	Good	Painting	Evening	Yes																																																																																															
12	Sham	Bad	Washing	Morning	No																																																																																															
13	Sham	Good	Washing	Evening	Yes																																																																																															
14	Sam	Bad	Plumbing	Morning	Yes																																																																																															
15	Sham	Good	Painting	Morning	No																																																																																															
Q. 6. (A)	6	2	<p>Explain the 'Hold-out' approach used in planning the training and testing data</p>	4																																																																																																
(B)	6	3	<p>A Confusion matrix as shown below, is observed for a machine learning algorithm, used to detect authentic messages and spam messages. Consider the positive class as – Authentic and Spam as negative class</p> <table><thead><tr><th></th><th>Predicted authentic</th><th>Predicted Spam</th></tr></thead><tbody><tr><td>Actual authentic</td><td>1202</td><td>5</td></tr><tr><td>Actual Spam</td><td>29</td><td>54</td></tr></tbody></table> <p>Calculate Accuracy, precision, recall and f-score</p>		Predicted authentic	Predicted Spam	Actual authentic	1202	5	Actual Spam	29	54	4																																																																																							
	Predicted authentic	Predicted Spam																																																																																																		
Actual authentic	1202	5																																																																																																		
Actual Spam	29	54																																																																																																		