



Daffodil International University
Faculty of Science & Information Technology
Department of Computer Science & Engineering
Mid Examination, Spring 2025

Course Code: CSE325, Course Title: Data Mining and Machine Learning
Level: 3 Term: 2 Batch: 62

Time: 01:30 Hrs

Marks: 25

Answer ALL Questions

[The figures in the right margin indicate the full marks and corresponding course outcomes. All portions of each question must be answered sequentially.]

1.	<p>Predict whether a student will Pass an exam based on Hours Studied and Attendance.</p> <table><tr><th>Example</th><th>Hours Studied</th><th>Attendance</th><th>Pass</th></tr><tr><td>1</td><td>High</td><td>Good</td><td>Yes</td></tr><tr><td>2</td><td>High</td><td>Good</td><td>Yes</td></tr><tr><td>3</td><td>High</td><td>Poor</td><td>Yes</td></tr><tr><td>4</td><td>High</td><td>Poor</td><td>No</td></tr><tr><td>5</td><td>Low</td><td>Good</td><td>Yes</td></tr><tr><td>6</td><td>Low</td><td>Good</td><td>No</td></tr><tr><td>7</td><td>Low</td><td>Poor</td><td>No</td></tr><tr><td>8</td><td>Low</td><td>Poor</td><td>No</td></tr></table> <p>(i) Calculate the overall entropy of the dataset for the Pass attribute. (ii) For each attribute (Hours Studied and Attendance):<ul style="list-style-type: none">• Calculate the entropy of each subset.• Compute the weighted average entropy for the attribute.(iii) Identify which attribute gives the highest information gain and explain why it should be chosen as the root of the decision tree. (iv) Draw the final decision tree using the selected attribute, showing the branches and final classifications.</p>	Example	Hours Studied	Attendance	Pass	1	High	Good	Yes	2	High	Good	Yes	3	High	Poor	Yes	4	High	Poor	No	5	Low	Good	Yes	6	Low	Good	No	7	Low	Poor	No	8	Low	Poor	No	[10]	CO2
Example	Hours Studied	Attendance	Pass																																				
1	High	Good	Yes																																				
2	High	Good	Yes																																				
3	High	Poor	Yes																																				
4	High	Poor	No																																				
5	Low	Good	Yes																																				
6	Low	Good	No																																				
7	Low	Poor	No																																				
8	Low	Poor	No																																				
2.	<p>A hospital has developed an automated system to detect diabetic retinopathy from patient images. The system was tested on 800 images. Out of these: (a) 250 images are known to be positive for diabetic retinopathy, (b) The system correctly identified 200 positive cases, and (c) It incorrectly flagged 70 healthy images as positive. Now answer the below questions. (i) Construct the confusion matrix for this binary classification problem. (ii) Calculate the precision and recall for the positive (disease) class.</p>	[5]	CO1																																				

3.	<p>A retail company wants to predict whether a customer will make a purchase ("Yes" or "No") based on two features: Age and Annual Income (in thousands). Using $K=3$, classify a new customer with Age = 33 and Annual Income = 52 using the following dataset:</p> <table border="1"><thead><tr><th>Customer ID</th><th>Age</th><th>Annual Income (k)</th><th>Purchase?</th></tr></thead><tbody><tr><td>1</td><td>25</td><td>40</td><td>Yes</td></tr><tr><td>2</td><td>30</td><td>50</td><td>Yes</td></tr><tr><td>3</td><td>22</td><td>20</td><td>No</td></tr><tr><td>4</td><td>35</td><td>60</td><td>Yes</td></tr><tr><td>5</td><td>40</td><td>70</td><td>No</td></tr><tr><td>6</td><td>28</td><td>30</td><td>No</td></tr><tr><td>7</td><td>32</td><td>55</td><td>Yes</td></tr><tr><td>8</td><td>45</td><td>80</td><td>No</td></tr></tbody></table> <p>(i) Calculate the Euclidean distances between the new customer (33,52) and each customer in the dataset. (ii) Identify the three nearest neighbors and determine the predicted class by majority vote.</p>	Customer ID	Age	Annual Income (k)	Purchase?	1	25	40	Yes	2	30	50	Yes	3	22	20	No	4	35	60	Yes	5	40	70	No	6	28	30	No	7	32	55	Yes	8	45	80	No	[5]	CO2
Customer ID	Age	Annual Income (k)	Purchase?																																				
1	25	40	Yes																																				
2	30	50	Yes																																				
3	22	20	No																																				
4	35	60	Yes																																				
5	40	70	No																																				
6	28	30	No																																				
7	32	55	Yes																																				
8	45	80	No																																				
4.	<p>Consider the following dataset with two features, X_1 and X_2:</p> <table border="1"><thead><tr><th>Data Point</th><th>X_1</th><th>X_2</th></tr></thead><tbody><tr><td>1</td><td>1.0</td><td>2.0</td></tr><tr><td>2</td><td>1.5</td><td>1.8</td></tr><tr><td>3</td><td>5.0</td><td>8.0</td></tr><tr><td>4</td><td>8.0</td><td>8.0</td></tr><tr><td>5</td><td>1.0</td><td>0.6</td></tr><tr><td>6</td><td>9.0</td><td>11.0</td></tr></tbody></table> <p>Using $K=2$ and the following initial centroids: (a) Centroid 1: Data Point 1 (1.0,2.0) and (b) Centroid 2: Data Point 4 (8.0,8.0). Now answer the below questions.</p> <p>(i) Compute the Euclidean distance between each data point and both centroids. Then, assign each data point to the nearest centroid. (ii) Calculate the new centroids for each cluster by taking the average of the X_1 and X_2 values of the data points assigned to that cluster.</p>	Data Point	X_1	X_2	1	1.0	2.0	2	1.5	1.8	3	5.0	8.0	4	8.0	8.0	5	1.0	0.6	6	9.0	11.0	[5]	CO2															
Data Point	X_1	X_2																																					
1	1.0	2.0																																					
2	1.5	1.8																																					
3	5.0	8.0																																					
4	8.0	8.0																																					
5	1.0	0.6																																					
6	9.0	11.0																																					