

COURSE CODE (CREDITS): 18B11BI611 (3)

MAX. MARKS: 35

COURSE NAME: Machine Learning for Bioinformatics

COURSE INSTRUCTORS: D. Gupta

MAX. TIME: 2 Hours

Note: (a) All questions are compulsory. (b) The candidate can make suitable numeric assumptions wherever required to solve problems. (c) **Be concise.** (d) Use of a calculator is permitted.

Q. No.	Question	CO	Marks																		
Q. 1	<p>a) Explain how high bias causes underfitting and low bias causes overfitting in predicting student exam scores. Also, discuss how including relevant features like hours studied, attendance, and previous grades can help balance bias and improve the model's performance.</p> <p>b) Briefly differentiate between stochastic, mini-batch, and batch gradient descent. Provide real-world examples to illustrate each method.</p> <p>c) Explain the differences between Gini impurity and split information. Consider real-world examples to illustrate why these measures are important for making effective splits in machine learning applications.</p>	1, 2, 4	[3]																		
Q. 2	<p>a) Given the following observations from the survey with two attributes (acid durability and strength) to classify whether a special paper tissue is good or not. Find the class of the test sample using k-NN algorithm. Take k=3 and use L2 norm for distance computations.</p> <table border="1"><thead><tr><th>X1= Acid durability (seconds)</th><th>X2= Strength (kg/ square meter)</th><th>Y= Classification</th></tr></thead><tbody><tr><td>7</td><td>7</td><td>Bad</td></tr><tr><td>7</td><td>4</td><td>Bad</td></tr><tr><td>3</td><td>4</td><td>Good</td></tr><tr><td>1</td><td>4</td><td>Good</td></tr></tbody></table> <p>Test Sample</p> <table border="1"><tbody><tr><td>3</td><td>7</td><td>?</td></tr></tbody></table>	X1= Acid durability (seconds)	X2= Strength (kg/ square meter)	Y= Classification	7	7	Bad	7	4	Bad	3	4	Good	1	4	Good	3	7	?	2, 3, 4	[4]
X1= Acid durability (seconds)	X2= Strength (kg/ square meter)	Y= Classification																			
7	7	Bad																			
7	4	Bad																			
3	4	Good																			
1	4	Good																			
3	7	?																			

	b) What is the kernel trick in SVM? Explain how does the RBF kernel help in non-linear classification taking a real-world example.		[2]																								
Q. 3	<p>a) Using the following dataset predict the class for the record (Confident=Yes, Sick=No) applying Naïve Bayes algorithm</p> <table border="1"> <thead> <tr> <th>Confident</th><th>Studied</th><th>Sick</th><th>Result</th></tr> </thead> <tbody> <tr> <td>Yes</td><td>No</td><td>No</td><td>Fail</td></tr> <tr> <td>Yes</td><td>No</td><td>Yes</td><td>Pass</td></tr> <tr> <td>No</td><td>Yes</td><td>Yes</td><td>Fail</td></tr> <tr> <td>No</td><td>Yes</td><td>No</td><td>Pass</td></tr> <tr> <td>Yes</td><td>Yes</td><td>Yes</td><td>Pass</td></tr> </tbody> </table>	Confident	Studied	Sick	Result	Yes	No	No	Fail	Yes	No	Yes	Pass	No	Yes	Yes	Fail	No	Yes	No	Pass	Yes	Yes	Yes	Pass	4	[5]
Confident	Studied	Sick	Result																								
Yes	No	No	Fail																								
Yes	No	Yes	Pass																								
No	Yes	Yes	Fail																								
No	Yes	No	Pass																								
Yes	Yes	Yes	Pass																								
	b) Using the Laplace smoothing formula, explain the problem it solves in text classification with Naive Bayes and mention one limitation of this method?		[2]																								
Q. 4	<p>a) Using the K-Means clustering algorithm with K=2, cluster the following data points: (2,3), (3,2), (4,4), (8,7), (7,8).</p> <p>Use Euclidean distance as the similarity measure. Initialize the algorithm by selecting the first two points as the initial centroids. Perform two iterations and clearly show the formation of clusters after each iteration.</p> <p>b) What is a dendrogram in hierarchical agglomerative clustering? Explain different linkage criteria affect the formation of clusters.</p>	5	[4]																								
Q. 5	<p>a) What is a genetic algorithm? Explain its main steps using the binary population: 1010, 1100, 1001, 0110 Perform one complete iteration of the algorithm - including fitness evaluation, selection, crossover (1-point), mutation (bit flip mutation), and show the resulting new population.</p> <p>b) Explain the Silhouette Coefficient and Dunn Index along with their formulas. How do these metrics help in determining the optimal number of clusters in clustering algorithms?</p>	5	[4]																								
			[3]																								