

**Department of Computer Science**  
**University of Delhi**  
**Master of Computer Application**  
**MCAC 204: Machine Learning**  
**Unique Paper Code: 223421206**

**Semester II**  
**Year of admission: 2023**

**Time: Three Hours**

**Max. Marks: 70**

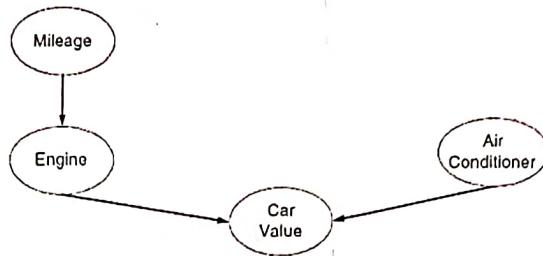
**Instructions:**

1. All questions carry equal marks.
2. Use proper notation and show complete working for full credit.

1	a. How does the choice of learning rate ( $\alpha$ ) affect the convergence and stability of the gradient descent algorithm in linear regression? Illustrate with the help of a suitable example.	[4]																							
	b. Consider the following dataset:	[4]																							
	<table border="1" style="margin: auto;"> <thead> <tr> <th>Size</th><th>Color</th><th>Shape</th><th>Class/ Label</th></tr> </thead> <tbody> <tr> <td>Big</td><td>Red</td><td>Circle</td><td>No</td></tr> <tr> <td>Small</td><td>Red</td><td>Triangle</td><td>No</td></tr> <tr> <td>Small</td><td>Red</td><td>Circle</td><td>Yes</td></tr> <tr> <td>Big</td><td>Blue</td><td>Circle</td><td>No</td></tr> <tr> <td>Small</td><td>Blue</td><td>Circle</td><td>Yes</td></tr> </tbody> </table>	Size	Color	Shape	Class/ Label	Big	Red	Circle	No	Small	Red	Triangle	No	Small	Red	Circle	Yes	Big	Blue	Circle	No	Small	Blue	Circle	Yes
Size	Color	Shape	Class/ Label																						
Big	Red	Circle	No																						
Small	Red	Triangle	No																						
Small	Red	Circle	Yes																						
Big	Blue	Circle	No																						
Small	Blue	Circle	Yes																						
2	Use the candidate elimination learning algorithm to determine the most general and the most specific hypothesis for the given training data.	[6]																							
	c. Consider the training examples $\{(x_1, y_1), (x_2, y_2), (x_3, y_3), \dots, (x_m, y_m)\}$ and the following cost function for a linear regression problem:	[6]																							
	$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (wx^{(i)} + b - y^{(i)})^2$	[6]																							
Prove that the parameters $w$ and $b$ that minimizes $J(w, b)$ can be derived as $\frac{\text{cov}(x, y)}{\text{var}(x)}$ and $\bar{y} - w\bar{x}$ , respectively.																									
2	a. What is the purpose of feature scaling in multiple linear regression? Illustrate the same with the help of a suitable example.	[4]																							
	b. How does the sigmoid function define the decision boundary in logistic regression? Also, give the intuition behind the cost function used in logistic regression with regularization.	[5]																							

	<p>c. Consider a dataset with binary class labels (+1 or -1). You are using the AdaBoost algorithm with decision stumps as the base classifiers. Initially, all samples are assigned equal weights. The first decision stump is trained and achieves an error rate of 0.2.</p> <ol style="list-style-type: none"> <li>Calculate the Amount of Say (alpha) assigned to this weak learner in the final ensemble.</li> <li>After updating the sample weights based on the performance of the first weak learner, determine the weight assigned to a sample that was misclassified by this weak learner.</li> </ol>	[5]																									
2	<p>a. Each binary classifier in an ensemble makes predictions on an input <math>\mathbf{x}</math>, as shown in the table below. Using this table, find the ensemble model's aggregated prediction for <math>\mathbf{x}</math>:</p> <table border="1"> <thead> <tr> <th></th> <th>Classifier's Confidence</th> <th>Prediction</th> </tr> </thead> <tbody> <tr> <td>Classifier 1</td> <td>0.61</td> <td>+1</td> </tr> <tr> <td>Classifier 2</td> <td>0.53</td> <td>-1</td> </tr> <tr> <td>Classifier 3</td> <td>0.88</td> <td>-1</td> </tr> <tr> <td>Classifier 4</td> <td>0.34</td> <td>+1</td> </tr> </tbody> </table> <p>b. With the help of a suitable diagram, prove that the width of the margin in the linear support vector classifier (LSVM) is <math>\frac{2}{\ \vec{w}\ }</math>, where <math>\vec{w}</math> is a vector perpendicular to the decision boundary.</p> <p>c. Find a principal component of the dataset given below (Show all the intermediate steps):</p> <table border="1"> <thead> <tr> <th>Age</th> <th>Weight</th> </tr> </thead> <tbody> <tr> <td>20</td> <td>65</td> </tr> <tr> <td>25</td> <td>60</td> </tr> <tr> <td>5</td> <td>15</td> </tr> <tr> <td>10</td> <td>20</td> </tr> </tbody> </table>		Classifier's Confidence	Prediction	Classifier 1	0.61	+1	Classifier 2	0.53	-1	Classifier 3	0.88	-1	Classifier 4	0.34	+1	Age	Weight	20	65	25	60	5	15	10	20	[2]
	Classifier's Confidence	Prediction																									
Classifier 1	0.61	+1																									
Classifier 2	0.53	-1																									
Classifier 3	0.88	-1																									
Classifier 4	0.34	+1																									
Age	Weight																										
20	65																										
25	60																										
5	15																										
10	20																										
3	<p>a. Consider the following neural network, comprising four layers: layer 1, layer 2, layer 3, and layer 4. Layers 1, 2, and 3 are hidden layers, and layer 4 is the output layer.</p> <p>Determine the following:</p> <ol style="list-style-type: none"> <li>dimensions of <math>W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, W^{[3]}, b^{[3]}, W^{[4]}, b^{[4]}</math></li> </ol>	[6]																									

- ii. number of trainable parameters for Layer 1.
- b. Given the Bayesian belief network (BBN) and the corresponding dataset shown below:



Mileage	Engine	Air Conditioner	Number of Records with Car Value = High	Number of Records with Car Value = Low
High	Good	Working	3	4
High	Good	Broken	1	2
High	Bad	Working	1	5
High	Bad	Broken	0	4
Low	Good	Working	10	0
Low	Good	Broken	3	1
Low	Bad	Working	2	2
Low	Bad	Broken	0	2

- i. Draw the probability table for each node in the BBN.  
 ii. Use the BBN to compute  $P(Engine = Bad, Air Conditioner = Broken)$ .

- 4
- a. For a classifier that distinguishes between cat vs non-cat, which of the four activations should be employed at the output layer: Sigmoid, Leaky ReLU, and Tanh? Justify your answer. [2]
- b. For a neural network that outputs the likelihood of occurrence of a disease out of six, which activation functions should be employed: Sigmoid, ReLU, SoftMax, and Tanh? Justify your answer. [2]
- c. Consider a dataset comprising 1 billion instances. Should a 70:30 train-test split be a good choice? Justify your choice. [2]
- d. For the following classification tasks, identify which metric, precision, or recall would be most suitable. Justify your answer.  
 i) COVID-19 identification using chest X-rays.  
 ii) Spam e-mail classification. [3]
- e. Write kernel K-means algorithm. [5]

5	<p>a. Consider the following diagram.</p> <p>Show that when conditioned on <b>C</b>, <b>A</b> and <b>B</b> are orthogonal to each other.</p> <p>b. Assume we have two users and three movies. The <math>2 \times 3</math> matrix <math>Y</math> is given below:</p> $Y = \begin{bmatrix} 1 & 2 & ? \\ ? & 10 & 15 \end{bmatrix}$ <p>Our goal is to find the matrices <math>U</math> and <math>V</math> such that <math>X = UV^T</math> closely approximates the observed ratings in <math>Y</math>. Assume we start by fixing <math>V</math> to initial values of <math>[2 \ 4 \ 6]^T</math>. Find the optimal <math>2 \times 1</math> vector <math>U</math> in this case. (Express your answer in terms of <math>\lambda</math>, where <math>\lambda</math> is a regularization parameter).</p> <p>c. Consider a problem with three binary random variables relating to the fuel system on a car. The variables are called <b>B</b>, representing the state of a battery that is either <b>charged</b> (<math>B = 1</math>) or <b>flat</b> (<math>B = 0</math>), <b>F</b> representing the state of the fuel tank that is either <b>full of fuel</b> (<math>F = 1</math>) or <b>empty</b> (<math>F = 0</math>), and <b>G</b>, which is the state of an electric fuel gauge and which indicates either <b>full</b> (<math>G = 1</math>) or <b>empty</b> (<math>G = 0</math>). The battery is either charged or flat, and independently, the fuel tank is either full or empty, with prior probabilities given below:</p> <p><math>p(B = 1) = 0.9</math>  <math>p(F = 1) = 0.9</math></p> <p>Also, consider the probability table given below:</p> <table style="margin-left: auto; margin-right: auto;"> <tr> <td><math>p(G = 1   B = 1, F = 1) = 0.8</math></td> </tr> <tr> <td><math>p(G = 1   B = 1, F = 0) = 0.2</math></td> </tr> <tr> <td><math>p(G = 1   B = 0, F = 1) = 0.2</math></td> </tr> <tr> <td><math>p(G = 1   B = 0, F = 0) = 0.1</math></td> </tr> </table> <p>Compute the following:</p> <ul style="list-style-type: none"> <li>i. <math>p(G = 0)</math></li> <li>ii. <math>p(F = 0   G = 0)</math></li> <li>iii. <math>p(F = 0   G = 0, B = 0)</math></li> </ul>	$p(G = 1   B = 1, F = 1) = 0.8$	$p(G = 1   B = 1, F = 0) = 0.2$	$p(G = 1   B = 0, F = 1) = 0.2$	$p(G = 1   B = 0, F = 0) = 0.1$	<p>[3]</p> <p>[4]</p> <p>[2+2 +3]</p>
$p(G = 1   B = 1, F = 1) = 0.8$						
$p(G = 1   B = 1, F = 0) = 0.2$						
$p(G = 1   B = 0, F = 1) = 0.2$						
$p(G = 1   B = 0, F = 0) = 0.1$						