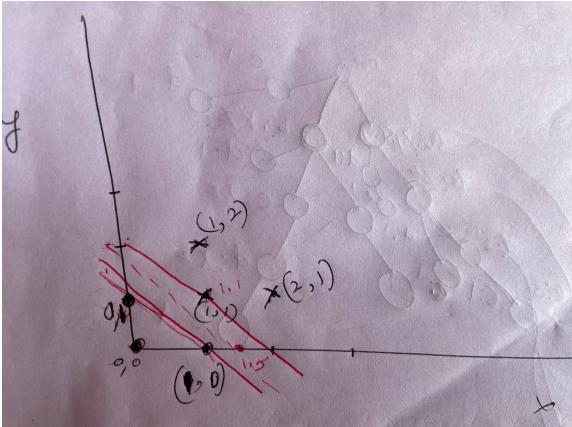


## DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

### MID SEMESTER EXAMINATION-III – Scheme & Solution

<b>Course Title with code</b>	Introduction to Machine Learning, 18CSE751	<b>Maximum Marks</b>	<b>30 Marks</b>			
<b>Date and Time</b>	02-02-22, 2 pm to 3 pm	<b>No. of Hours</b>	1.0			
<b>Course Instructor(s)</b>	Dr. Vani V					
<b>Instructions to Students</b>						
1. Answer any <b>two full questions.</b> 2. Any missing data may assume suitably.						

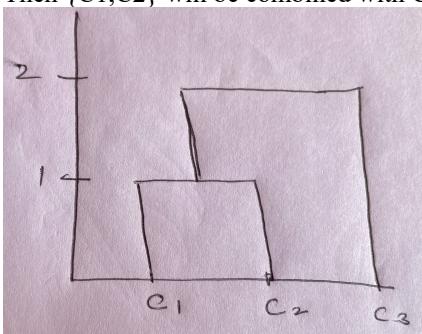
Q. No	Question	MAX MARKS	CO	BL	PO/P SO
1. a	<p>Compare parametric approach with non-parametric approach with necessary examples</p> <p><b>Solution:</b></p> <p><b>Parametric approach :</b></p> <p>The sample comes from a known distribution.</p> <ul style="list-style-type: none"> <li>- The advantage of any parametric approach is that given a model, the problem reduces to the estimation of a small number of parameters, which, in the case of density estimation, are the sufficient statistics of the density.</li> <li>- For example: the mean and covariance in the case of Gaussian densities.</li> <li>- used quite frequently</li> <li>- Assume rigid parametric model may be a source of bias in many applications where this assumption does not hold.</li> </ul> <p><b>Examples:</b></p> <p>Logistic Regression, Linear Discriminant Analysis, Perceptron, Naive Bayes, Simple Neural Networks</p> <p><b>Non-Parametric Approach:</b></p> <p>Algorithms that do not make strong assumptions about the form of the mapping function are called nonparametric machine learning algorithms.</p> <p>By not making assumptions, they are free to learn any functional form from the training data.</p> <p>Non-parametric methods seek to best fit the training data in constructing the mapping function, whilst maintaining some ability to generalize to unseen data.</p> <p>As such, they are able to fit a large number of functional forms. An easy to understand nonparametric model is the k-nearest neighbours algorithm that makes predictions based on the k most similar training patterns for a new data instance.</p> <p>The method does not assume anything about the form of the mapping function other than patterns that are close are likely to have a similar output variable.</p> <p><b>Examples:</b></p> <p>k-Nearest Neighbors, Decision Trees like CART and C4.5, Support Vector Machines</p>	4	4	2	1,2,3/2

1. b	<p>Consider the following training data and answer the questions (a) to (c).</p> <p>class 1: (1,1), (1,2), (2,1)      class 2: (0,0), (1,0), (0,1)</p> <p><b>Solution:</b></p> <p>(a) Plot these six training points. Are the classes {class1, class2} linearly separable?</p> 	6	4	3	1,2,3/2
1.c	<p>The maximum margin hyperplane should have a slope of <math>-1</math> and should satisfy <math>x_1 = 3/2</math>, <math>x_2 = 0</math>. Therefore its equation is <math>x_1 + x_2 = 3/2</math>, and the weight vector is <math>(1, 1)^T</math></p> <p>(c) If you remove one of the support vectors does the size of the optimal margin decrease, stay the same, or increase?      In this specific dataset the optimal margin increases when we remove the support vectors</p> <p><b>Scheme:</b>      Each answer carries 2 marks  <math>3 \times 2 = 6</math> Marks</p> <p>Discuss the need for linear, polynomial and RBF kernel functions in Support Vector Machine (SVM) and show how to use SVM for multi-class classification.</p> <p>Kernel Function generally transforms the training set of data so that a non-linear decision surface can be transformed to a linear equation in a higher number of dimension spaces.</p> <p>Linear Kernel: used when data is linearly separable.</p> <p>Polynomial Kernel: It represents the similarity of vectors in training set of data in a feature space over polynomials of the original variables used in kernel.</p> <p>RBF: It is used to perform transformation when there is no prior knowledge about data</p> <ul style="list-style-type: none"> <li>polynomials up to some degree <math>s</math> in the elements <math>x_k</math> of the input vector (e.g., <math>x_3^3</math> or <math>x_1 \times x_4</math>) with kernel:</li> </ul> $K(x, y) = (1 + x^T y)^s \quad (8.19)$ <p>For <math>s = 1</math> this gives a linear kernel</p> <ul style="list-style-type: none"> <li>radial basis function expansions of the <math>x_k</math>s with parameter <math>\sigma</math> and kernel:</li> </ul> $K(x, y) = \exp(-(x - y)^2 / 2\sigma^2)$ <ul style="list-style-type: none"> <li>The SVM only works for two classes.</li> <li>For the problem of N-class classification, <ul style="list-style-type: none"> <li>train an SVM that learns to classify class one from all other classes,</li> <li>then another that classifies class two from all the others.</li> <li>So, for N-classes, we have N SVMs.</li> </ul> </li> <li>This still leaves one problem: How do we decide which of these SVMs</li> </ul>	5	4	2	1,2,3/2

	<p>is the one that recognises the input?</p> <ul style="list-style-type: none"> <li>- <i>The answer is just to choose the one that makes the strongest prediction</i></li> <li>- <i>Classifier return either the class label (as the sign of <math>y</math>) or the value of <math>y</math></i></li> <li>- <i>The value of <math>y</math> is telling us how far away from the decision boundary it is, and clearly it will be negative if it is a misclassification.</i></li> </ul> <p><b>Therefore, use the maximum value of this soft boundary as the best classifier</b></p> <p>Different kernels produce different decision boundaries</p> <p>Scheme: 3 marks kernel, 2marks for Multiclass SVM</p>																																																																																																																																								
2. a	<p>For the given intial set of three clusters (<math>k=3</math>)</p> $C_1 = \{(1,3),(3,6),(3,5)\}$ $C_2 = \{(5,3),(6,7),(2,2)\}$ $C_3 = \{(6,5),(3,1),(2,3)\}$ <p>Apply k-means algorithm and show when the clusters converge.</p> <p>Note: Use Manhattan Distance</p> <p><b>Solution :</b></p> <p>Let points be considered as 1 to 9</p> <table border="1"> <thead> <tr> <th colspan="3">Initial</th> <th>M1</th> <th>M2</th> <th>M3</th> </tr> <tr> <th></th> <th>Data Points</th> <th></th> <th>2.33</th> <th>4.6667</th> <th>4.333</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>1</td> <td>3</td> <td>3</td> <td></td> <td>4.333</td> </tr> <tr> <td>2</td> <td>3</td> <td>6</td> <td>2</td> <td></td> <td>3.333</td> </tr> <tr> <td>3</td> <td>3</td> <td>5</td> <td>1</td> <td></td> <td>2.333</td> </tr> <tr> <td>4</td> <td>5</td> <td>3</td> <td>4.33</td> <td></td> <td>1.667</td> </tr> <tr> <td>5</td> <td>6</td> <td>7</td> <td>6</td> <td></td> <td>4.667</td> </tr> <tr> <td>6</td> <td>2</td> <td>2</td> <td>3</td> <td></td> <td>4.333</td> </tr> <tr> <td>7</td> <td>6</td> <td>5</td> <td>4</td> <td></td> <td>2.667</td> </tr> <tr> <td>8</td> <td>3</td> <td>1</td> <td>4.33</td> <td></td> <td>4.333</td> </tr> <tr> <td>9</td> <td>2</td> <td>3</td> <td>2</td> <td></td> <td>3.333</td> </tr> </tbody> </table> <table border="1"> <thead> <tr> <th colspan="3">Iteration 1:</th> <th>M1</th> <th>M2</th> <th>M3</th> </tr> <tr> <th></th> <th>Data Points</th> <th></th> <th>3</th> <th>5.5</th> <th>6</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>1</td> <td>3</td> <td>4.5</td> <td></td> <td>8</td> </tr> <tr> <td>2</td> <td>3</td> <td>6</td> <td>0.5</td> <td></td> <td>3</td> </tr> <tr> <td>3</td> <td>3</td> <td>5</td> <td>0.5</td> <td></td> <td>4</td> </tr> <tr> <td>4</td> <td>5</td> <td>3</td> <td>4.5</td> <td></td> <td>4</td> </tr> <tr> <td>5</td> <td>6</td> <td>7</td> <td>4.5</td> <td></td> <td>1</td> </tr> <tr> <td>6</td> <td>2</td> <td>2</td> <td>4.5</td> <td></td> <td>8</td> </tr> <tr> <td>7</td> <td>6</td> <td>5</td> <td>3.5</td> <td></td> <td>1</td> </tr> <tr> <td>8</td> <td>3</td> <td>1</td> <td>4.5</td> <td></td> <td>8</td> </tr> <tr> <td>9</td> <td>2</td> <td>3</td> <td>3.5</td> <td></td> <td>7</td> </tr> </tbody> </table> <p>Initial and Iteration 1 are same the algorithm converges</p> <p>C1(2,3) C2(5,7) C3(1,4,6,8,9)</p> <p><b>Scheme :</b> Each iteration : 5 Marks , <math>5 \times 2 = 10</math> Marks</p>	Initial			M1	M2	M3		Data Points		2.33	4.6667	4.333	1	1	3	3		4.333	2	3	6	2		3.333	3	3	5	1		2.333	4	5	3	4.33		1.667	5	6	7	6		4.667	6	2	2	3		4.333	7	6	5	4		2.667	8	3	1	4.33		4.333	9	2	3	2		3.333	Iteration 1:			M1	M2	M3		Data Points		3	5.5	6	1	1	3	4.5		8	2	3	6	0.5		3	3	3	5	0.5		4	4	5	3	4.5		4	5	6	7	4.5		1	6	2	2	4.5		8	7	6	5	3.5		1	8	3	1	4.5		8	9	2	3	3.5		7	10	5	3	1,2,3/2
Initial			M1	M2	M3																																																																																																																																				
	Data Points		2.33	4.6667	4.333																																																																																																																																				
1	1	3	3		4.333																																																																																																																																				
2	3	6	2		3.333																																																																																																																																				
3	3	5	1		2.333																																																																																																																																				
4	5	3	4.33		1.667																																																																																																																																				
5	6	7	6		4.667																																																																																																																																				
6	2	2	3		4.333																																																																																																																																				
7	6	5	4		2.667																																																																																																																																				
8	3	1	4.33		4.333																																																																																																																																				
9	2	3	2		3.333																																																																																																																																				
Iteration 1:			M1	M2	M3																																																																																																																																				
	Data Points		3	5.5	6																																																																																																																																				
1	1	3	4.5		8																																																																																																																																				
2	3	6	0.5		3																																																																																																																																				
3	3	5	0.5		4																																																																																																																																				
4	5	3	4.5		4																																																																																																																																				
5	6	7	4.5		1																																																																																																																																				
6	2	2	4.5		8																																																																																																																																				
7	6	5	3.5		1																																																																																																																																				
8	3	1	4.5		8																																																																																																																																				
9	2	3	3.5		7																																																																																																																																				
2. b	<p>Discuss the use of clustering in supervised learning with suitable examples</p> <p><b>Solution:</b></p> <ul style="list-style-type: none"> <li>■ Dimensionality reduction methods find correlations between features and group features</li> <li>■ Clustering methods find similarities between instances and group instances</li> <li>■ Allows knowledge extraction through</li> </ul>	5	5	2	1,2,3/2																																																																																																																																				

	<p><i>number of clusters, prior probabilities, cluster parameters, i.e., center, range of features.</i></p> <p>Example: CRM, customer segmentation</p> <ul style="list-style-type: none"> <li>■ Estimated group labels <math>h_j</math> (soft) or <math>b_j</math> (hard) may be seen as the dimensions of a new <math>k</math> dimensional space, where we can then learn our discriminant or regressor.</li> <li>■ Local representation (only one <math>b_j</math> is 1, all others are 0; only few <math>h_j</math> are nonzero) vs Distributed representation (After PCA; all <math>z_j</math> are nonzero)</li> </ul> <p><b>Scheme:</b> 5 x 1 = 5 Marks</p>			
3. a	<p>Identify a suitable algorithm to estimate the gender of the students studying in a college given their height information. Discuss the identified algorithm with its iterative steps.</p> <p><b>Solution:</b></p> <ul style="list-style-type: none"> <li>■ The Expectation-Maximization (EM) algorithm is used in maximum likelihood estimation where the problem involves two sets of random variables of which one, <b>X, is observable and the other, Z, is hidden.</b></li> <li>■ The goal of the algorithm is to find the parameter vector <math>\Phi</math> that maximizes the likelihood of the observed values of X, <math>L(\Phi X)</math>. But in cases where this is not feasible, we associate the extra hidden variables Z and express the underlying model using both, to maximize the likelihood of the joint distribution of X and Z, the complete likelihood <math>L_c(\Phi X,Z)</math>.</li> <li>■ Since the Z values are not observed, we cannot work directly with the complete data likelihood <math>L_c</math>;</li> <li>■ Instead, we work with its expectation, <math>Q</math> given X and the current parameter values <math>\Phi^l</math>, where l indexes iteration. This is the expectation (E) step of the algorithm.</li> <li>■ Then in the maximization (M) step, we look for the new parameter values, <math>\Phi^{l+1}</math>, that maximize this.</li> <li>■ Iterate the two steps <ol style="list-style-type: none"> <li>1. E-step: Estimate <math>z</math> given X and current <math>\Phi</math></li> <li>2. M-step: Find new <math>\Phi'</math> given <math>z</math>, X, and old <math>\Phi</math>.</li> </ol> <math display="block">\text{E - step : } Q(\Phi   \Phi') = E[\mathcal{L}_c(\Phi   X, Z)   X, \Phi']</math> <math display="block">\text{M - step : } \Phi'^{l+1} = \underset{\Phi}{\operatorname{argmax}} \mathcal{Q}(\Phi   \Phi')</math> <p>An increase in Q increases incomplete likelihood</p> <math display="block">\mathcal{L}(\Phi'^{l+1}   X) \geq \mathcal{L}(\Phi'   X)</math> </li> </ul> <p><b>Scheme:</b> identifying and explaining EM algorithm : 2 Marks E-Step : 1.5 Mark M-Step: 1.5 Mark</p>	5	5	2
3. b	<p>For the given initial set of three clusters</p> $C_1 = \{(1,3),(3,6),(3,5)\}$ $C_2 = \{(5,3),(6,7),(2,2)\}$ $C_3 = \{(6,5),(3,1),(2,3)\}$ <p>Apply the agglomerative algorithm with complete cluster distance measures to produce a dendrogram tree</p> <p>Note : Use Manhattan distance to compute distance matrix</p> <p><b>Solution:</b></p> <p>C1 -&gt; C2 : max distance : <math>(5+4) = 9</math>  C3-&gt;C1 : max distance <math>(5+2) = 7</math>  C3-&gt; C1 : max distance <math>(4+4)= 8</math></p>	10	5	3

C1 and C2 will be combined after calculating max distances  
 Then {C1,C2} will be combined with C3



Calculation of max distances between clusters :  $3 \times 3 = 9$  Marks  
 Dendrogram : 1 Mark

Faculty Signature	Course Co-Ordinator/Mentor Signature	HoD Signature