

USN

## Internal Assessment Test 3 – March 2024

Sub:	<b>Artificial Intelligence and Machine Learning</b>				Sub Code:	<b>21CS54</b>	Branch:	<b>CSE</b>		
Date:	<b>14/03/2024</b>	Duration:	90 mins	Max Marks:	50	Sem/Sec:	V / A, B&C			OBE
<b><u>Answer any FIVE FULL Questions</u></b>								MAR KS	CO	RBT
1	Write the C4.5 algorithm. How is a C4.5 algorithm better than ID3? What are the attribute selection measures used in the algorithm?						(5+3+2)	CO2	L2	
2	Consider a dataset containing information about patients to predict whether a person has “Heart Disease” or not. Apply the CART algorithm to build a decision tree model						10	CO3	L3	
		<b>Chest Pain</b>	<b>Good Blood Circulation</b>	<b>Blocked Arteries</b>	<b>Heart Disease</b>					
		No	No	No	No					
		Yes	Yes	Yes	Yes					
		Yes	Yes	No	No					
		No	No	Yes	Yes					
		No	Yes	No	No					
		Yes	No	Yes	Yes					
3	State Bayes theorem. Write an algorithm of the Naïve Naves classifier. List out and write the working principle of variants of Naïve Bayes classifiers.						(2+5+3)	CO2	L2	

USN

## Internal Assessment Test 3 – March 2024

Sub:	<b>Artificial Intelligence and Machine Learning</b>				Sub Code:	<b>21CS54</b>	Branch:	<b>CSE</b>		
Date:	<b>14/03/2024</b>	Duration:	90 mins	Max Marks:	50	Sem/Sec:	V / A, B&C			OBE
<b><u>Answer any FIVE FULL Questions</u></b>								MA RKS	CO	RBT
1	Write the C4.5 algorithm. How is a C4.5 algorithm better than ID3? What are the attribute selection measures used in the algorithm?						(5+3+2)	CO2	L2	
2	Consider a dataset containing information about patients to predict whether a person has “Heart Disease” or not. Apply the CART algorithm to build a decision tree model						10	CO3	L3	
		<b>Chest Pain</b>	<b>Good Blood Circulation</b>	<b>Blocked Arteries</b>	<b>Heart Disease</b>					
		No	No	No	No					
		Yes	Yes	Yes	Yes					
		Yes	Yes	No	No					
		Yes	No	No	Yes					
		No	No	Yes	Yes					
		No	Yes	No	No					
		Yes	No	Yes	Yes					
3	State Bayes theorem. Write an algorithm of the Naïve Naves classifier. List out and write the working principle of variants of Naïve Bayes classifiers.						(2+5+3)	CO2	L2	

4	(a) Draw and explain the structure of an Artificial Neural Network. (b) Write an algorithm of the perceptron Learning.	5 5	CO2	L2															
5	Consider the network architecture with 4 input and 2 output units. Consider four training samples each vector of length 4. Training samples: I1: (1,1,1,0) I2: (0,0,1,1) I3: (1,0,0,1) I4: (0,0,1,0) Output units: Unit1 and Unit2, learning rate 0.6. Initial weight matrix: Unit 1: 0.3 0.8 0.7 0.2 Unit 2: 0.6 0.7 0.4 0.6 Identify an algorithm to learn without supervision. Show the weight updates of the SOFM during the first epoch of the learning algorithm.	10	CO4	L3															
6	Consider the following dataset. Cluster it using the K-Means algorithm with the initial value of objects 2 and 4 as initial seeds.	10	CO5	L3															
<table border="1"> <thead> <tr> <th>S.No</th> <th>X</th> <th>Y</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>3</td> <td>5</td> </tr> <tr> <td>2</td> <td>7</td> <td>8</td> </tr> <tr> <td>3</td> <td>12</td> <td>5</td> </tr> <tr> <td>4</td> <td>16</td> <td>9</td> </tr> </tbody> </table>		S.No	X	Y	1	3	5	2	7	8	3	12	5	4	16	9			
S.No	X	Y																	
1	3	5																	
2	7	8																	
3	12	5																	
4	16	9																	

CI

CCI

HoD

4	(a) Draw and explain the structure of an Artificial Neural Network. (b) Write an algorithm of the perceptron Learning.	5 5	CO2	L2															
5	Consider the network architecture with 4 input and 2 output units. Consider four training samples each vector of length 4. Training samples: I1: (1,1,1,0) I2: (0,0,1,1) I3: (1,0,0,1) I4: (0,0,1,0) Output units: Unit1 and Unit2, learning rate 0.6. Initial weight matrix: Unit 1: 0.3 0.8 0.7 0.2 Unit 2: 0.6 0.7 0.4 0.6 Identify an algorithm to learn without supervision. Show the weight updates of the SOFM during the first epoch of the learning algorithm.	10	CO4	L3															
6	Consider the following dataset. Cluster it using the K-Means algorithm with the initial value of objects 2 and 4 as initial seeds.	10	CO5	L3															
<table border="1"> <thead> <tr> <th>S.No</th> <th>X</th> <th>Y</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>3</td> <td>5</td> </tr> <tr> <td>2</td> <td>7</td> <td>8</td> </tr> <tr> <td>3</td> <td>12</td> <td>5</td> </tr> <tr> <td>4</td> <td>16</td> <td>9</td> </tr> </tbody> </table>		S.No	X	Y	1	3	5	2	7	8	3	12	5	4	16	9			
S.No	X	Y																	
1	3	5																	
2	7	8																	
3	12	5																	
4	16	9																	

CI

CCI

HoD

---

**(1.) Write the C4.5 algorithm. How is a C4.5 algorithm better than ID3? What are the attribute selection measures used in the algorithm?**

**C4.5 algorithm:**

(5)

1. Compute entropy for the whole training dataset based on the target attribute.

$$\text{Entropy\_Info}(T) = - \sum_{i=1}^m P_i \log_2 P_i$$

2. Compute Entropy, Information gain, Split information and Gain Ratio for each of the attribute in the training dataset.

$$\text{Entropy\_Info}(T, A) = \sum_{i=1}^v \frac{|A_i|}{|T|} \times \text{Entropy\_Info}(A_i)$$

$$\text{Information\_Gain}(A) = \text{Entropy\_Info}(T) - \text{Entropy\_Info}(T, A)$$

$$\text{Split\_Info}(T, A) = - \sum_{i=1}^v \frac{|A_i|}{|T|} \times \log_2 \frac{|A_i|}{|T|}$$

$$\text{Gain\_Ratio}(A) = \frac{\text{Info\_Gain}(A)}{\text{Split\_Info}(T, A)}$$

3. Choose the attribute for which the maximum Gain ratio is the best-split attribute.

4. The best-split attribute is placed as the root node.

5. The root node is branched into subtrees with each subtree as an outcome of the test condition of the root node attribute. Accordingly, the training dataset is also split into subsets.

6. Recursively apply the same operation for the subset of the training set with the remaining attributes until a leaf node is derived or no more training instances are available in the subset.

<i>Features</i>	<i>ID3</i>	<i>C4.5</i>
Type of data	Categorical	Continuous and Categorical
Speed	Low	Faster than ID3
Boosting	Not supported	Not supported
Pruning	No	Pre-pruning
Missing Values	Can't deal with	Can't deal with
Formula	Use information entropy and information Gain	Use split info and gain ratio

(3)

Metrics:

(2)

Entropy, Information gain, Split information and Gain Ratio

**(3) State Bayes theorem. Write an algorithm of the Naïve Bayes classifier. List out and write the working principle of variants of Naïve Bayes classifiers.**

**Bayes theorem:**

**(2)**

It can be written as:

$$P(\text{Hypothesis } h \mid \text{Evidence } E) = \frac{P(\text{Evidence } E \mid \text{Hypothesis } h) P(\text{Hypothesis } h)}{P(\text{Evidence } E)}$$

**Algorithm: Naïve Bayes Classifier:**

**(5)**

1. Compute the prior probability for the target class.
2. Compute the Frequency matrix and likelihood probability for each of the attributes.
3. Use Bayes theorem to calculate the probability of all hypotheses.

It can be written as:

$$P(\text{Hypothesis } h \mid \text{Evidence } E) = \frac{P(\text{Evidence } E \mid \text{Hypothesis } h) P(\text{Hypothesis } h)}{P(\text{Evidence } E)}$$

4. Use Maximum A Posteriori (MAP) hypothesis  $h_{MAP}$  to classify the test object to the hypothesis with the highest probability.

$$\begin{aligned} h_{MAP} &= \max_{h \in H} P(\text{Hypothesis } h \mid \text{Evidence } E) \\ &= \max_{h \in H} \frac{P(\text{Evidence } E \mid \text{Hypothesis } h) P(\text{Hypothesis } h)}{P(\text{Evidence } E)} \\ &= \max_{h \in H} P(\text{Evidence } E \mid \text{Hypothesis } h) P(\text{Hypothesis } h) \end{aligned}$$

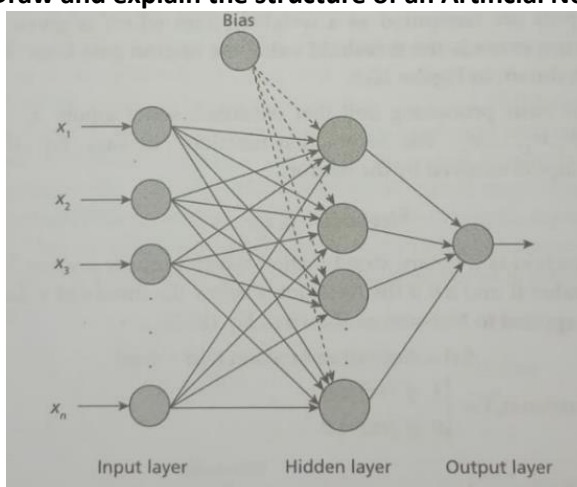
Variants of Naïve Bayes classifiers:

**(3)**

- 1) Bernoulli Naïve Bayes classifier
  - Discrete features
  - Boolean variables
- 2) Multinomial Naïve Bayes classifier
  - Generalization of Bernoulli Naïve Bayes classifier
- 3) Multi-class Naïve Bayes classifier
  - Multiclassification

**(4) (a) Draw and explain the structure of an Artificial Neural Network.**

**(5)**



**(b) Write an algorithm of perceptron Learning.**

**(5)**

### Algorithm 10.1: Perceptron Algorithm

Set initial weights  $w_1, w_2, \dots, w_n$  and bias  $\theta$  to a random value in the range  $[-0.5, 0.5]$ .

For each Epoch,

1. Compute the weighted sum by multiplying the inputs with the weights and add the products.

2. Apply the activation function on the weighted sum:

$$Y = \text{Step}((x_1 w_1 + x_2 w_2) - \theta)$$

3. If the sum is above the threshold value, output the value as positive else output the value as negative.

4. Calculate the error by subtracting the estimated output  $Y_{\text{estimated}}$  from the desired output  $Y_{\text{desired}}$ :

$$\text{error } e(t) = Y_{\text{desired}} - Y_{\text{estimated}}$$

[If error  $e(t)$  is positive, increase the perceptron output  $Y$  and if it is negative, decrease the perceptron output  $Y$ .]

5. Update the weights if there is an error:

$$\Delta w_i = \alpha \times e(t) \times x_i$$

$$w_i = w_i + \Delta w_i$$

where,  $x_i$  is the input value,  $e(t)$  is the error at step  $t$ ,  $\alpha$  is the learning rate and  $\Delta w_i$  is the difference in weight that has to be added to  $w_i$ .

Q:2

<u>Dataset:</u>		Good Blood Circulation	Blocked Arteries	Heart Disease
	Chest pain			
①	No	No	No	No
②	Yes	Yes	Yes	Yes
③	Yes	Yes	No	No
④	Yes	Yes	No	Yes
⑤	Yes	No	Yes	Yes
⑥	No	No	No	No
⑦	No	Yes	No	No
⑧	Yes	Yes	Yes	Yes
⑨	Yes	No	Yes	Yes

Solution:

Iteration 1:

Gini-Index for the whole dataset

$$\text{Gini-Index}(T) = 1 - \left(\frac{4}{7}\right)^2 - \left(\frac{3}{7}\right)^2$$

[formula:  $\text{Gini-Index}(T) = 1 - \sum_{i=1}^m p_i^2$ ]

$$= 1 - 0.327 - 0.184$$

$$= 0.489$$

First attribute: 'chest pain'

possible values: No, Yes, so subsets: 4

No. of Yes: 4

No. of No: 3

Chest Pain	Heart Disease = Yes	Heart Disease = No
Yes	3	1
No	1	2

subsets:

{Yes}, {No}, ~~{Yes, No}~~

$$\text{Gini-Index}(T, \text{chest pain} \in \{\text{Yes}\})$$

$$= 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2$$

$$= 1 - 0.563 - 0.063$$

$$= 0.374$$

$$\text{Gini-Index}(T, \text{chest pain} \in \{\text{No}\})$$

$$= 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2$$

$$= 1 - 0.11 - 0.44$$

$$= 0.45$$

$$\text{Gini-Index}(T, \text{chest pain} \in \{\text{Yes, No}\})$$

$$= \frac{|S_1|}{T} \text{Gini}(S_1) + \frac{|S_2|}{T} \text{Gini}(S_2)$$

$$= \frac{4}{7} \times 0.374 + \frac{3}{7} \times 0.45$$

$$= 0.214 + 0.19$$

$$= 0.404$$

$$\Delta \text{Gini}(\text{Chest pain}) = \text{Gini}(T) - \text{Gini}(\text{Chest pain})$$

$$= 0.489 - 0.404$$

$$= 0.085$$

Second attribute:

Yes: 3, No: 4

Good Blood Circulation

Good Blood Circulation

Heart Disease = Yes

Heart Disease = No

Yes

1

2

No

3

1

Subsets:

{Yes}, {No}

$$\text{Gini-Index}(T, \text{GBC} \in \{\text{Yes}\})$$

$$= 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2$$

$$= 1 - 0.11 - 0.44 = 0.45$$

$$\text{Gini-Index}(T, \text{GBC} \in \{\text{No}\})$$

$$= 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 = 1 - 0.563 - 0.063$$

$$= 0.374$$



$$\text{Gini-Index (Good B.C. } \in \{Yes, No\})$$

$$= \frac{3}{7} \times 0.45 + \frac{4}{7} \times 0.374$$

$$= ~~0.193~~ 0.404$$

$$\Delta \text{Gini (GBC)} = \text{Gini(T)} - \text{Gini (GBC)}$$

$$= 0.085$$

Third attribute: Blocked Asterres

Yes: 3, No: 4

Blocked Asterres

Hersot Disease = Yes

Hersot Disease = No

Yes

3

0

NO

1

3

Subsets:

{Yes}, {No}

$$\text{Gini-Index (BA } \in \{Yes\})$$

$$= 1 - \left(\frac{3}{3}\right)^2$$

$$= 0$$

$$\text{Gini-Index (BA } \in \{NO\})$$

$$= 1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2$$

$$= ~~0.374~~ 0.374$$

$$\text{Gini-Index (BA } \in \{Yes, No\})$$

$$= \frac{3}{7} \times 0 + \frac{4}{7} \times 0.374$$

$$= 0.214$$

$$\Delta \text{Gini (BA)} = \text{Gini(T)} - \text{Gini (BA)}$$

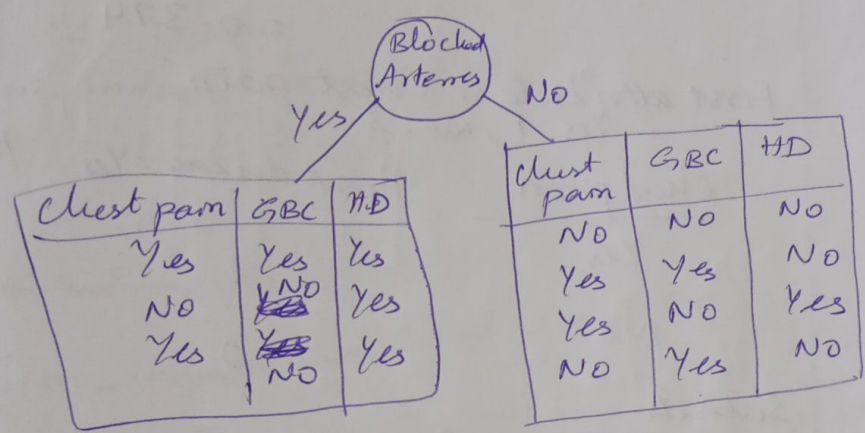
$$= 0.489 - 0.214$$

$$= 0.275$$

Gini-index &  $\Delta$ Gini for all the attributes

Attribute	Gini-index	$\Delta$ Gini
Chest pain	0.404	0.085
GBC	0.404	0.085
BA	0.214	0.275

Blocked Arteries has the highest  $\Delta$ Gini value. So 'Blocked Arteries' is the root node.

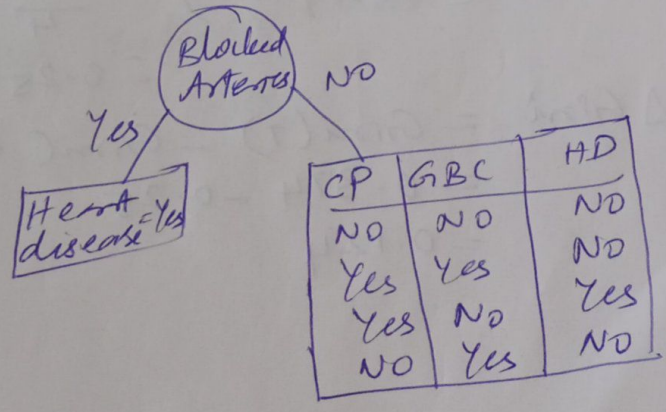


Iteration 2: Left subtree

Chest pain	GBC	HD
Yes	Yes	Yes
NO	NO	Yes
Yes	NO	Yes

same class

$$\text{Gini Index (T)} = 1 - \left(\frac{3}{3}\right)^2 - \left(\frac{0}{3}\right)^2 = 0$$



Iteration 2:

CP	GBC	HD
NO	NO	NO
Yes	Yes	NO
Yes	NO	Yes
NO	Yes	NO

$$\text{Gini-Index}(T) = 1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2$$
$$= 0.374$$

First attribute: 'Chest pain'  
Yes: 1, No: 3

Chest pain	Heart disease = Yes	HD = No
Yes	1	1
No	0	2

Subsets:

{ Yes }, { No }

$$\text{Gini-Index}(CP \in \{\text{yes}\})$$

$$= 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2$$

$$= 1 - 0.25 - 0.25$$

$$= 0.5$$

$$\text{Gini-Index}(CP \in \{\text{no}\})$$

$$= 1 - \left(\frac{0}{2}\right)^2 - \left(\frac{2}{2}\right)^2$$

$$= 0$$

$$\text{Gini-Index}(\text{Chest pain}) = \frac{2}{4} \times 0.5$$

$$= 0.25$$

$$\Delta \text{Gini} = \text{Gini}(T) - \text{Gini}(\text{Chest pain})$$

$$= 0.374 - 0.25$$

$$= 0.124$$

2<sup>nd</sup> attribute: Good Blood Circulation

Yes: 2, No: 2

GBC	HD = Yes	HD = NO
Yes	0	2
NO	1	1

~~Gini-index~~

subset:

{Yes}, {No}

Gini-index (GBC  $\in$  {Yes})

$$= 0$$

Gini-index (GBC  $\in$  {No})

$$= 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2$$

$$= 0.5$$

$$\Delta Gini = Gini(T) -$$

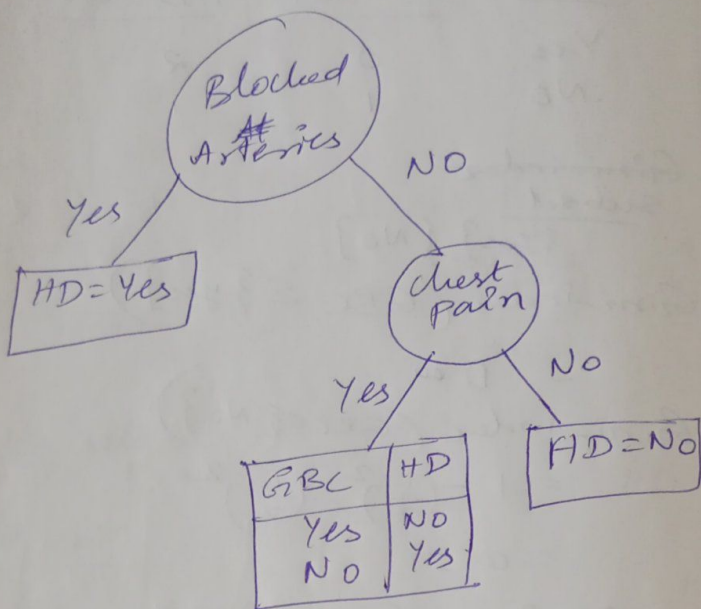
$$\begin{aligned} \text{Gini-index (Good Blood Circulation)} & \left| \begin{aligned} \Delta Gini(GBC) &= Gini(T) - \\ & Gini(GBC) \\ &= 0.375 - \\ & \quad 0.25 \\ &= 0.124 \end{aligned} \right. \\ &= \frac{2}{4} \times 0 + \frac{2}{4} \times 0.5 \\ &= 0.25 \end{aligned}$$

Gini-index &  $\Delta Gini$  values

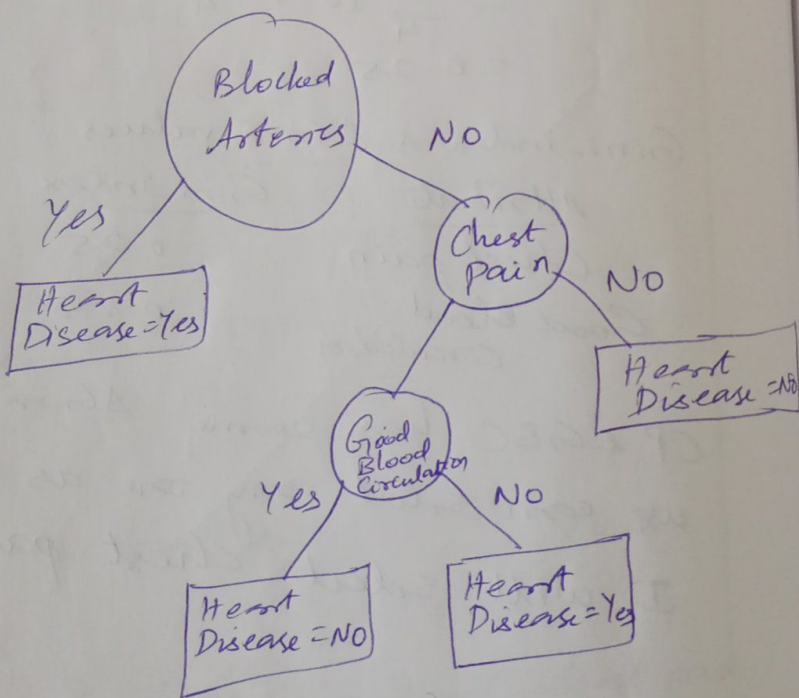
Attribute	Gini-index	$\Delta Gini$
Chest pain	0.25	0.124
Good Blood Circulation	0.25	0.124

CP & GBC has same  $\Delta Gini$  value,  
we can select any one as root node  
I will select 'Chest pain' as root node.

# The classification decision tree



## Final tree:



15/3/2024  
Friday

## Network Architecture

Q: 5  
4 input units  
2 output units (unit 1 & unit 2)

I/P: 4 training samples (each vector of length 4)

$$I_1: (1, 1, 1, 0)$$

$$I_2: (0, 0, 1, 1)$$

$$I_3: (1, 0, 0, 1)$$

$$I_4: (0, 0, 1, 0)$$

Learning rate = 0.6

Initial weight matrix:

$$\begin{array}{l} \text{unit 1:} \\ \text{unit 2:} \end{array} \begin{bmatrix} 0.3 & 0.8 & 0.7 & 0.2 \\ 0.6 & 0.7 & 0.4 & 0.6 \end{bmatrix}$$

Solution:

Iteration 1:

Training sample  $I_1: (1, 1, 1, 0)$   
weight matrix

$$\begin{array}{l} \text{unit 1:} \\ \text{unit 2:} \end{array} \begin{bmatrix} 0.3 & 0.8 & 0.7 & 0.2 \\ 0.6 & 0.7 & 0.4 & 0.6 \end{bmatrix}$$

Compute Euclidean distance between  $I_1$  and unit 1 weights

$$\begin{aligned} d^2 &= (0.3-1)^2 + (0.8-1)^2 + (0.7-1)^2 + (0.2-0)^2 \\ &= 0.49 + 0.04 + 0.09 + 0.04 \\ &= 0.66 \end{aligned}$$

Compute Euclidean distance between  $I_1$  and unit 2 weights

$$\begin{aligned} d^2 &= (0.6-1)^2 + (0.7-1)^2 + (0.4-1)^2 + (0.6-0)^2 \\ &= 0.16 + 0.09 + 0.36 + 0.36 \\ &= 0.97 \end{aligned}$$

- Distance between  $I_1$  and Unit 1 is shorter, so Unit 1 wins.

- Update the weights of Unit 1

Formula:

$$w(t+1) = w(t) + \eta (\text{Input} - w(t))$$

Unit 1's weight vector

$$= [0.3 \ 0.8 \ 0.7 \ 0.2] + 0.6 \left( [1 \ 1 \ 1 \ 0] - [0.3 \ 0.8 \ 0.7 \ 0.2] \right)$$

$$= [0.3 \ 0.8 \ 0.7 \ 0.2] + 0.6 (0.7 \ 0.2 \ 0.3 \ -0.2)$$

$$= [0.3 \ 0.8 \ 0.7 \ 0.2] + [0.42 \ 0.12 \ 0.18 \ -0.12]$$

$$= [0.72 \ 0.92 \ 0.88 \ 0.08]$$

$$\begin{array}{l} \text{Unit 1} : [0.72 \ 0.92 \ 0.88 \ 0.08] \\ \text{Unit 2} : [0.6 \ 0.7 \ 0.4 \ 0.6] \end{array}$$

Iteration 2:  $I_2 : (0, 0, 1, 1)$

Euclidean distance between  $I_2$  & Unit 1

$$d^2 = (0.72-0)^2 + (0.92-0)^2 + (0.88-1)^2 + (0.08-1)^2$$

$$= 0.52 + 0.85 + 0.014 + 0.85$$

$$= 2.234$$

Euclidean distance between  $I_2$  & Unit 2

$$d^2 = (0.6-0)^2 + (0.7-0)^2 + (0.4-1)^2 + (0.6-1)^2$$

$$= 0.36 + 0.49 + 0.36 + 0.16$$

$$= 1.37$$

Unit 2 wins

update weights of

$$\text{Unit 2} = [0.6 \ 0.7 \ 0.4 \ 0.6] + 0.6$$

$$\left( [0 \ 0 \ 1 \ 1] - [0.6 \ 0.7 \ 0.4 \ 0.6] \right)$$

$$= [0.6 \ 0.7 \ 0.4 \ 0.6] + 0.6 (-0.6 \ -0.7 \ 0.6 \ 0.4)$$

$$= [0.6 \ 0.7 \ 0.4 \ 0.6] + [-0.36 \ -0.42 \ 0.36 \ 0.24]$$

$$= [0.24 \ 0.28 \ 0.76 \ 0.84]$$

$$\text{Unit 1} \quad \begin{bmatrix} 0.72 & 0.92 & 0.88 & 0.08 \\ 0.24 & 0.28 & 0.04 & 0.36 \end{bmatrix}$$

Iteration 3:

$$I_3: (1, 0, 0, 1)$$

Euclidean distance between  $I_3$  & Unit 1

$$d^2 = (0.72 - 1)^2 + (0.92 - 0)^2 + (0.88 - 0)^2 + (0.08 - 1)^2$$

$$= 0.08 + 0.85 + 0.77 + 0.85$$

$$= 2.55$$

Euclidean distance between  $I_3$  & Unit 2

$$d^2 = (0.24 - 1)^2 + (0.28 - 0)^2 + (0.04 - 0)^2 + (0.36 - 1)^2$$

$$= 0.58 + 0.08 + 0.002 + 0.41$$

$$= 1.072 \quad 1.27$$

Unit 2 wins

update weights  
of unit 2

$$= [0.24 \quad 0.28 \quad 0.04 \quad 0.36] + 0.6 ([1, 0, 0, 1] - [0.24 \quad 0.28 \quad 0.04 \quad 0.36])$$

$$= [0.24 \quad 0.28 \quad 0.04 \quad 0.36] + [0.46 \quad -0.17 \quad -0.024 \quad 0.38]$$

$$= [0.7 \quad 0.11 \quad 0.02 \quad 0.74]$$

$$\text{Unit 1: } \begin{bmatrix} 0.72 & 0.92 & 0.88 & 0.08 \\ 0.7 & 0.11 & 0.02 & 0.74 \end{bmatrix}$$

Iteration 4:  $I_4: (0, 0, 1, 0)$

Euclidean distance b/w  $I_4$  and Unit 1

$$d^2 = (0.72 - 0)^2 + (0.92 - 0)^2 + (0.88 - 1)^2 + (0.08 - 0)^2$$

$$= 0.52 + 0.85 + 0.014 + 0.006$$

$$= 1.39$$



Euclidean distance between  $I_4$  & Unit 2

$$d^2 = (0.7 - 0)^2 + (0.11 - 0)^2 + (0.02 - 1)^2 + (0.74 - 0)^2$$

$$= 0.49 + 0.012 + 0.96 + 0.55$$

$$= 2.012$$

Euclidean distance between  $I_4$  & Unit 1

Unit 1 wins

update weights of unit 1

$$([0 \ 0 \ 1 \ 0] - [0.72 \ 0.92 \ 0.88 \ 0.08]) + 0.6$$

$$= [0.72 \ 0.92 \ 0.88 \ 0.08] + 0.6$$

$$= [0.72 \ 0.92 \ 0.88 \ 0.08] + [-0.43 \ -0.55 \ 0.07 \ -0.05]$$

$$= [0.29 \ 0.37 \ 0.95 \ 0.03]$$

$$\left. \begin{array}{l} \text{unit 1} \\ \text{unit 2} \end{array} \right\} = \begin{bmatrix} 0.29 & 0.37 & 0.95 & 0.03 \\ 0.7 & 0.11 & 0.02 & 0.74 \end{bmatrix}$$

Best mapping units for each of the inputs taken one

$I_1 : (1, 1, 1, 0) \rightarrow$  Unit 1

$I_2 : (0, 0, 1, 1) \rightarrow$  Unit 2

$I_3 : (1, 0, 0, 1) \rightarrow$  Unit 2

$I_4 : (0, 0, 1, 0) \rightarrow$  Unit 1

- This is the result after epoch 1.

12/10/2024  
Monday  
p. 6

Cluster the dataset using the k-Means algorithm with the initial value of objects 2 and 4 as initial seeds

S.No	x	y
1	3	5
2	7	8
3	12	5
4	16	9

Solution:

(7, 8) & (16, 9) are initial seeds.

- These data points are stored as two clusters. So  $k=2$ .

- Centroids: (7, 8) & (16, 9)

Iteration:

- Compare all the data points with the centroid and assign to the nearest cluster.

<sup>1st</sup> data point: (3, 5)

- calculate the distance between the data point and the centroids of cluster 1 and cluster 2.

$$\begin{aligned} \text{Dist}(1, \text{centroid}_1) &= \sqrt{(3-7)^2 + (5-8)^2} \\ &= \sqrt{4^2 + 3^2} = \sqrt{16+9} \\ &= \sqrt{25} = 5 \end{aligned}$$

$$\begin{aligned} \text{Dist}(1, \text{centroid}_2) &= \sqrt{(3-16)^2 + (5-9)^2} \\ &= \sqrt{13^2 + 4^2} \\ &= \sqrt{169+16} = \sqrt{185} \\ &= 13.60 \end{aligned}$$

$(3, 5)$  is closer to cluster 1.  
Object 3:  $(12, 5)$

$$\begin{aligned}\text{Dist}(3, \text{centroid}_1) &= \sqrt{(12-7)^2 + (5-8)^2} \\ &= \sqrt{5^2 + 3^2} = \sqrt{25+9} \\ &= \sqrt{34} = 5.831\end{aligned}$$

$$\begin{aligned}\text{Dist}(3, \text{centroid}_2) &= \sqrt{(12-16)^2 + (5-9)^2} \\ &= \sqrt{4^2 + 4^2} = \sqrt{16+16} \\ &= \sqrt{32} = 5.66\end{aligned}$$

Object 3 is closer to cluster 2.

	cluster 1	cluster 2
	$(7, 8)$	$(16, 9)$
	$(3, 5)$	$(12, 5)$
centroid:	$(\bar{x}, \bar{y})$	$(\bar{x}, \bar{y})$
	$(\frac{10}{2}, \frac{13}{2})$	$(\frac{28}{2}, \frac{14}{2})$
centroid:	$(5, 6.5)$	$(14, 7)$

Iteration 2:

Object 1:  $(7, 8)$       centroid 1:  $(5, 6.5)$   
centroid 2:  $(14, 7)$

$$\begin{aligned}\text{Dist}(1, \text{centroid}_1) &= \sqrt{(7-5)^2 + (8-6.5)^2} \\ &= \sqrt{2^2 + 1.5^2} \\ &= \sqrt{4 + 2.25} \\ &= \sqrt{6.25} = 2.5\end{aligned}$$

$$\begin{aligned} \text{Dist}(1, \text{centroid}_2) &= \sqrt{(7-14)^2 + (8-7)^2} \\ &= \sqrt{49+1} \\ &= 7.07 \end{aligned}$$

Object 1 belongs to the same cluster 2.

Object 3: (12, 5)      centroid 1: (5, 6.5)  
centroid 2: (14, 7)

$$\begin{aligned} \text{Dist}(3, \text{centroid}_1) &= \sqrt{(12-5)^2 + (5-6.5)^2} \\ &= \sqrt{49 + 2.25} = \sqrt{51.25} \\ &= 7.16 \end{aligned}$$

Object 3 is closer to

$$\begin{aligned} \text{Dist}(3, \text{centroid}_2) &= \sqrt{(12-14)^2 + (5-7)^2} \\ &= \sqrt{2^2 + 2^2} = \sqrt{8} \\ &= 2.83 \end{aligned}$$

Object 3 is closer to cluster 2

cluster 1	cluster 2
(7, 8)	(16, 9)
(3, 5)	(12, 5)
centroid: (5, 6.5)	(14, 7)

- There is no change in the clusters.
- Therefore k-means algorithm terminates with two clusters with the given data points.