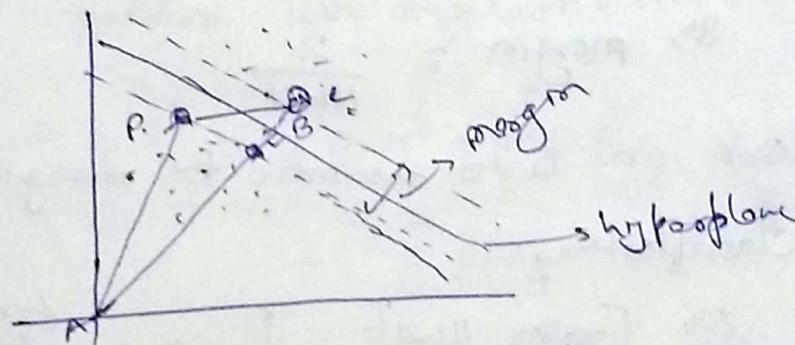


Q2) (a) Ans

Vladimir Vapnik is credited for developing Support Vector Machine (SVM). His main contribution was to use statistical approach for classifying rather than probabilistic approach, he has been awarded by Bayesian Postulate model.



When we project vector w i.e. normal to hyperplane, we get it lies on either side of the hyperplane

$$\vec{w} \cdot \vec{x} + b \geq c \quad \text{if } c = +1$$

Decision Rule $\begin{cases} \vec{w} \cdot \vec{x} + b \geq 0 & \text{then } \text{same} \\ \vec{w} \cdot \vec{x} + b \leq 0 & \text{then } \text{not} \end{cases}$

or $\vec{w} \cdot \vec{x} + b \geq 1 \quad \text{if } y=+1$

$\text{and } \vec{w} \cdot \vec{x} + b \leq -1 \quad \text{if } y=-1$

combining these two,

$$y(1)(\vec{w} \cdot \vec{x} + b) - 1 \geq 0 \quad \text{---(3)}$$

AB when projection on unit vector

$$(\vec{x}_+ - \vec{x}_-) \cdot \frac{\vec{w}}{\|\vec{w}\|}$$

From ③, we get

$$\vec{x}_+ \cdot \vec{w} + b - 1 \geq 0 \Rightarrow \vec{x}_+ \cdot \vec{w} = \frac{1-b}{\|\vec{w}\|}$$

For negative vec

$$\vec{x}_- = \frac{b-1}{\|\vec{w}\|}$$

$$\text{so, margin} = \frac{2}{\|\vec{w}\|}$$

Our aim is to maximize the margin for better classification.

$$\text{so, } \left[\text{maximize } \frac{\|\vec{w}\|^2}{2} \right] \xrightarrow{\text{for mathematically convenience}}$$

$$\text{subject to } y_i (\vec{w} \cdot \vec{x}_i + b) - 1 \geq 0$$

Lagrangian ~~problem~~

$$L = \frac{1}{2} \|\vec{w}\|^2 - \sum_i [y_i (\vec{w} \cdot \vec{x}_i + b) - 1]$$

\downarrow
Lagrange's multiplier

$$\frac{\partial L}{\partial \vec{w}} = 0 \quad \& \quad \frac{\partial L}{\partial b} = 0$$

$$\text{so, } \frac{\partial L}{\partial \vec{w}} = 0 \Rightarrow \vec{w} - \sum_i y_i \vec{x}_i = 0$$
$$\Rightarrow \vec{w} = \sum_i y_i \vec{x}_i$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow -\text{Edit}_i=0$$

$$\Rightarrow \boxed{-\text{Edit}_i=0}$$

For dual problem we will put these 2 values in Θ or Lagrange).

$$L(w, b, \alpha) = \sum \alpha_i - \frac{1}{2} \sum_m \sum_n \alpha_i \alpha_j y_i y_j x_i x_j$$

Dual problem with regularization term.

$$L(w, b, \alpha, \epsilon_i, \beta_i) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i (y_i(w \cdot x_i + b) - 1 + \epsilon_i) - \sum_{i=1}^n \beta_i \epsilon_i$$

$$\text{s.t. } \sum_{i=1}^m y_i \alpha_i = 0$$

$$\alpha_i \geq 0, i=1, 2, 3, \dots, n$$

Q2) (b) Ans \rightarrow

Dual problem was important because in primal we took α i.e. Lagrangian multiplier on support vectors or,

i) complexity increased with increase in no. of support vector

ii) support vectors are to be mentioned before hand.

so, primal equation failed. That's why we need dual ~~equation~~ problem.

$$\text{Max } L_g(\omega) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{r=1}^m \sum_{j=1}^m \alpha_r g_j g_r < O(n) O(m)$$

-①

Subject to $\alpha_i \geq 0, i=1, 2, \dots, m$

$$\sum_{i=1}^n \alpha_i g_i = 0$$

min. $J(\omega)$

s.t. $h_i(\omega) \leq 0 ; i=1, 2, 3, \dots$

$g_i(\omega) \leq 0 ; i=1, 2, 3, \dots$

$$L(\omega, \gamma, \beta) = J(\omega) - (\sum_{i=1}^n \alpha_i g_i(\omega) + \sum_{i=1}^m \beta_i h_i(\omega)) \quad -②$$

$$L_p(\omega) = \max_{\gamma, \beta \geq 0} L(\omega, \gamma, \beta) \quad -③$$

$$P = \min_{\omega} \max_{\gamma, \beta \geq 0} L(\omega, \gamma, \beta)$$

$$= \min_{\omega} L_p(\omega) \quad -④$$

$L_p(\omega)$ if $g(\omega) > 0$ i.e. violating the constraints

If violated,

$$L_p(\omega) \rightarrow \infty$$

If $h_i(\omega) = 0$ (violating constraints)

$$L_p(\omega) \rightarrow \infty$$

8) $L_p(\omega) = f(\omega)$, If constraint satisfied
 \sim otherwise

$$D^* = \max_{\alpha, \beta \geq 0} \min_w L(w, \alpha, \beta)$$

$$= \max_{\alpha, \beta \geq 0} L(\alpha, \beta)$$

$$D^* = P^* \text{ for max margin,}$$

solution of dual problem,

$$D^* = P^* = U(w^*, \alpha^*, \beta^*)$$

To max over w^* , α^* , β^* , we must satisfy,

$$\frac{\partial}{\partial w} L(w^*, \alpha^*, \beta) = 0$$

$$\frac{\partial}{\partial \alpha_i} L(w^*, \alpha^*, \beta) = 0$$

$$\alpha^* g(w^*) = 0$$

$$g(w^*) \leq 0$$

Basic coordinate descent algorithm,

max $w(d_1, d_2, \dots, d_m)$ (without bias constant)

Repeat { for $i=1$ to m

$$d_i = w(d_1, d_2, \dots, d_m)$$

Sequential minimal optimization in SMO, 2d
are taken at a time.

$$d_i \gamma - \sum_j \gamma_j d_j \leq 0$$

selected

- Select d_1, d_2 heuristically
- Hold all d_i fixed except d_1, d_2
- Optimize w.r.t. w.r.t. d_1, d_2 (s.t. s.t.)

lets put $\Gamma_{21} \& \Gamma_{22}$
randomly

update $d_1, d_2 \leftarrow \sum_{\Gamma=1}^n d_i y_i z_i$

$$d_1 y^{(1)} + d_2 y^{(2)} = - \sum_{\Gamma=1}^n d_i y_i - \varepsilon;$$

$$0 \leq d_i \leq C$$

Repeat until convergence,

1) a select some part of d_1, d_2 do update
w.r.t respect to α .

2) reoptimize w.r.t. w.r.t. d_1, d_2 hold
all d_k ($k \neq i, j$)

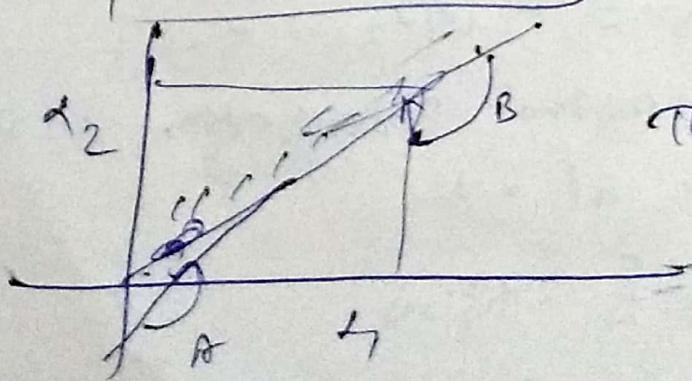
$$d_1 y^{(1)} + d_2 y^{(2)} = - \varepsilon - \varepsilon_i = \varepsilon_j$$

$$d_i = \frac{\varepsilon_j - d_2 y^{(1)}}{y^{(2)}}$$

$$w(x) = Ax^2 + Bx + c$$

$$L \leq d_i \leq H$$

Then we can minimize
quadratic opt.



$$d_2 \in \begin{cases} 1, 2, 3 & L \leq d_2 \leq H \\ L+1, \dots, H & \text{else} \end{cases}$$

05) (a) Ans

A generative model learns the joint probability distribution $p(x, y)$. It predicts the conditional probability with the help of Bayes theorem.

A discriminative model learns the conditional probability distribution $P(y|x)$. Both of these models are generally used in supervised learning problem.

Gaussian Discriminant Analysis (GDA) algorithm is a generative learning algorithm. As we know, generative learning algorithms make strong assumptions and can't predict the conditional probability directly. GDA assumes that $p(x|y)$ is distributed according to a multivariate normal distribution. ~~Assume~~ A multivariate gaussian distribution is fit for each class. This allows us to find $P(y)$ and $P(x|y)$. Using this two we can finally find $P(y|x)$, which is required for prediction.

GDA

- i.) Strong modelling Assumptions
- ii.) $d(d+1)/2 + 2d+1$ parameters
- iii.) A simple closed form solution.

Noise Bayes

- i.) Few stronger modelling assumptions.
- ii.) $3d+1$ parameters
- iii.) Super simple closed form solution.

UNDA makes an assumption about the probability distribution of the $P(X|y \leq k)$ where k is one of the classes. And it can be easily proved that for UNDA, if the initial assumptions about distribution of $\{P(X|y=k)\}$ (Gaussian) and $P(y)$ (Bernoulli) are true, the $P(y/X)$ can be expressed as sigmoid but reverse is not true.

Q4) Ans → Genetic Algorithm invented by Prof. John Holland.

Problem → $f(x) = x - x^2$, $x \in [0, 1]$; $x \in \mathbb{Z}$

We will take x in bits i.e. $(0000)_2$ to $(1111)_2$

Initial chromosome population of

Selection:

Initial Population	x	$f(x)$			
1010	10	-90	0.36	1.45	1
1100	12	-132	0.53	2.13	2
0011	3	-6	0.02	0.09	0
0101					

Iteration 1:

Initial population	α	$f(\alpha)$	$\frac{f_i}{\sum f_i}$	$\frac{f_i}{F}$	Putable wheel
1010	10	-90	0.36	2.18	2
1100	12	-132	0.53	3.20	3
0011	3	-6	0.02	0.14	6
0101	5	-20	0.08	0.48	0
0000	0	0	0	0	0
0001	1	0	0	0	0

$$\sum = -248$$

$$f = -41.3$$

mapping :-

$$\begin{matrix} 1100 \\ 1010 \end{matrix} \quad \begin{matrix} 1110 \\ 1000 \end{matrix}$$

$$\begin{matrix} 1100 \\ 0011 \end{matrix} \quad \begin{matrix} 1111 \\ 0000 \end{matrix}$$

$$\begin{matrix} 1100 \\ 0101 \end{matrix} \quad \begin{matrix} 1101 \\ 0100 \end{matrix}$$

Iteration 2:

Population	α	$f(\alpha)$	$\frac{f_i}{\sum f_i}$	$\frac{f_i}{F}$	Putable wheel
1110	14	182	0.3	1.22	2
1000	8	56	0.09	0.54	1
1101	13	156	0.25	1.52	2
0100	4	12	0.019	0.12	0
1111	15	210	0.34	2.04	2
0000	0	0	0	0	0

$$\sum = 616$$

$$f = 102.4$$

Scheme theory :

$$N = \boxed{1 \mid x \mid + \mid 1 \mid A}$$

Instances of N are : $\rightarrow 10010, 10011, 10110,$
 $10111, 11010, 11011, 11110,$
 11111

Order of $N = 2$ (No. of fixed bits)

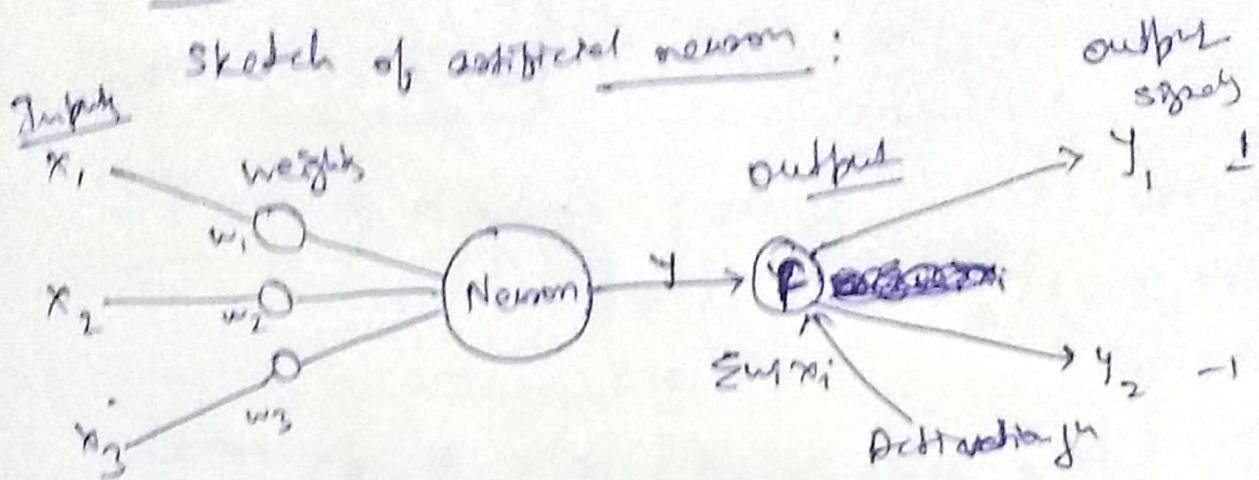
Defining length of following schemes:-

$$H_1 : x x 100 \Rightarrow 5-3 = 2$$

$$H_2 : x 0 x 1 x 0 x \Rightarrow 6-2 = 4$$

$$H_3 : 1 x x x x x 0 \Rightarrow 8-1 = 7$$

03) Any →



$$y = \begin{cases} +1, & \text{if } \sum x_i w_i \geq 0 \\ -1, & \text{if } \sum x_i w_i < 0 \end{cases}$$

Activation function used for generality signals can be any sigm function.

The middle layer hides the desired output that's why it is also called hidden layer. The neurons present in the hidden layer can be observed through the change in the behaviour of the network's input / output. The hidden layer present in between the input and output layers and from outside is present as abstract art.

The gradient of the error function with respect to the layer present is called as error gradient.

(a) Error gradient terms;

i) first term is the difference between output and target value.

ii) second term is the derivative of output layer activation function

iii) third term is the activation output of node j for the hidden layer.

For any weight w_{ij} , it can be written:

$$= \frac{\partial E}{\partial w_{ij}}$$

$$S_v(\theta) = y_v(\theta) / \log_v(\theta) e_v(\theta)$$

$$\text{where, } e_v(\theta) = y_{d,v}(\theta) - y_v(\theta)$$

The error gradient is required to calculate in order to find the effect of node v change on all the weights present in the neural network.

Any particular weight can be calculated as

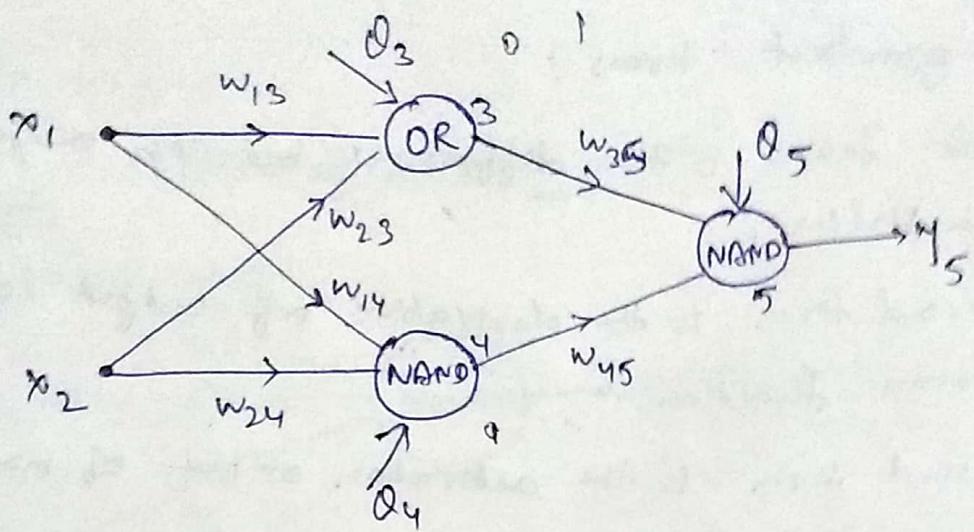
$$(w_{ij})_+ = w_{ij} + \eta \cdot \frac{\partial E}{\partial w_{ij}}$$

Working of NNML:

Neural network is a three-layer network, which consists of three neurons as,

i) 2 input neurons

ii) 1 output neuron



x_1	x_2	y_5
0	0	1
0	1	0
1	0	0
1	1	1

works let's bias terms θ has fixed initial values

i) Initialization of the weights:

$$w_{13} = 0.5, w_{14} = 0.9, w_{23} = 0.4, w_{24} = 1.0,$$

$$w_{35} = -1.2, w_{45} = 1.1, \theta_3 = 0.8, \theta_4 = 0.1, \theta_5 = 0.3$$

ii) Activation function:

Pass x_i forward through the network.

Activation $f^i \rightarrow \text{sigmoid}$.

$$y_3 = \text{sigmoid}(x_1 w_{13} + x_2 w_{23} + \theta_3)$$

$$= \frac{1}{1 + e^{-(0.5x_1 + 1 \times 0.4) - 0.8}} = 0.5265$$

here x_1 and $x_2 = 1$

$$y_4 = \text{sigmoid}(x_1 w_{14} + x_2 w_{24} - \theta_4)$$
$$= 0.8808$$

Now output of y_5 :

$$= \text{sigmoid}(y_3 w_{35} + y_4 w_{45} - \theta_5)$$
$$= 0.5092$$

Now, compute loss \rightarrow

$$e = y_{\text{des}} - y_5 = 1 - 0.5092 = 0.4908$$

→ Compute all derivatives)

$$\delta_5 = y_5(1-y_5) e = 0.5092(1-0.5092)(0.4908)$$
$$= 0.1225$$

$$\Delta w_{35} = \alpha \cdot y_3 \cdot \delta_5 \quad (\alpha=0.1)$$

$$= 0.1 \times 0.5260 \times 0.1225 = 0.0067$$

$$\Delta w_{45} = \alpha \cdot y_4 \cdot \delta_5$$
$$= 0.1 \times 0.8808 \times 0.1225 = 0.010$$

compute gradient for neuron 3 by,

$$\delta_3 = y_3(1-y_3) \delta_5 \cdot w_{23} = 0.8316$$

$$\delta_4 = y_4(1-y_4) \delta_5 \cdot w_{45} = 0.0141$$

The change in weight

$$\Delta w_{13} = d \alpha_1 s_3 = 0.1 \times 1 \times 0.0311 = 0.00311$$

$$\Delta w_{23} = d \alpha_2 s_3 = 0.1 \times 1 \times 0.0141 = 0.00141$$

$$\Delta Q_3 = d (-1) s_3 = -0.00311$$

$$\Delta w_{14} = d \alpha_1 s_4 = 0.00141$$

$$\Delta w_{24} = d \alpha_2 s_4 = 0.00141$$

$$\Delta Q_4 = d (-1) s_4 = -0.00141$$

\rightarrow weight update by:

$$w_{13} = w_{13} + \Delta w_{13} = 0.50311$$

$$w_{14} = w_{14} + \Delta w_{14} = 0.90141$$

$$w_{23} = w_{23} + \Delta w_{23} = 0.40141$$

$$w_{24} = w_{24} + \Delta w_{24} = 1.00141$$

$$w_{35} = w_{35} + \Delta w_{35} = 1.1916$$

$$w_{45} = w_{45} + \Delta w_{45} = 1.110$$

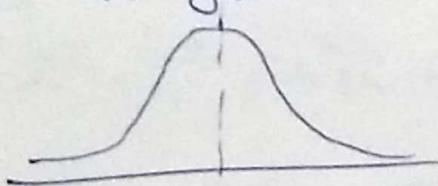
$$Q_3 = Q_3 + \Delta Q_3 = 0.2964$$

$$Q_4 = Q_4 + \Delta Q_4 = -0.14141$$

$$Q_5 = Q_5 + \Delta Q_5 = 0.28331$$

Q8.) (b) Ans

In the experiment we were told to classify the data in two classes using LDA with raw data itself and one by changing the raw data into normal distribution using Box-Muller Transformation. As we can see that after applying box-muller transformation the accuracy significantly increases. After applying the box-muller transformation, the data is distributed in gaussian bell shape as follows:



Data belongs to some class may be found near to mean or away from mean. This result in some form of clusters. The prob. of data point near to mean is high as compared to with points away from mean resulting in better probability model to classify the data.

Q5) (c) Ans

Let's assume we have training sample

$\{x^1, x^2, \dots, x^N\}$ which belongs to
 w_1 and N_2 belong to w_2 .

$$y = w^T x \quad \text{where } x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} \quad \text{and } w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_m \end{bmatrix}$$

So the mean of each class in \mathbb{R}^m feature space is.

$$\mu_i = \frac{1}{N} \sum_{x \in w_i} x$$

$$\begin{aligned} \mu_i^y &= \frac{1}{N} \sum_{y \in w_i} y = \frac{1}{N_i} \sum_{x \in w_i} w^T x \\ &= \frac{w^T}{N} \sum_{x \in w_i} x = w^T \mu_i \end{aligned}$$

So we can choose distance between the projected means as our objective function.

or cost.

$$\begin{aligned} J(w) &= (\mu_1^y - \mu_2^y)^2 = |w^T \mu_1 - w^T \mu_2|^2 \\ &= |w^T(\mu_1 - \mu_2)|^2 \end{aligned}$$

This is Fisher's solution to maximise a func.

$$\begin{aligned} S_i^y &= \sum_{y \in w_i} (y - \mu_i^y)^2 \\ &\text{variables} \\ &\text{with} \\ &\text{in class } i \end{aligned}$$

$$g) J(w) = \frac{(u_i - \hat{u}_i)^2}{s_1^2 + s_2^2} \rightarrow w \text{ has two class}$$

Now projection of y as a β of w in feature space x :

$$\begin{aligned}\hat{s}_i^2 &= \sum_{y \in \mathcal{Y}} (y \hat{u}_i)^2 = \sum_{x \in \mathcal{X}} (w^T x - w^T u_i)^2 \\ &= \sum_{x \in \mathcal{X}} w^T (x - u_i) (x - u_i)^T w \\ &= w^T s_i^2 w\end{aligned}$$

Now,

$$s_1^2 + s_2^2 = w^T w + w^T s_2^2 w = w^T (s_1^2 + s_2^2) w$$

now difference between this two.

$$\begin{aligned}(u_i - \hat{u}_i)^2 &= (w^T x_i - w^T u_i)^2 \\ &= w^T (u_i - x_i) (u_i - x_i)^T w \\ &= w^T s_B^2 w = \hat{s}_B^2\end{aligned}$$

Hence, Analog,

$$J(w) = \frac{(\hat{u}_i - u_i)^2}{s_1^2 s_2^2} = \frac{w^T s_B^2 w}{w^T s_B w}$$

\Rightarrow A relation of $J(w)$

- A dataset divided into train & test set, both PCF & LDA are applied to train data.
- Produce PCF & LDA model
- Original data in training set is transformed train set are produced which consist PCF instead of original feature & uses LDA.
- Then other transformed datasets are generated so that train finally have both

Q1) Ans: Given,

If $x \in A_1(0.5)$ AND

$y \in B_1(0.2)$ \Rightarrow Then $z \in C_1$

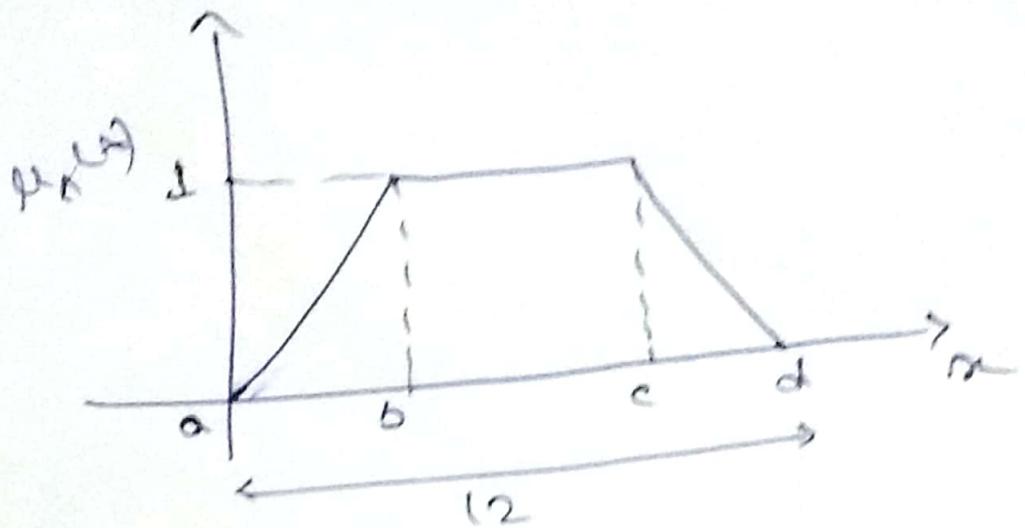
and,

If $x \in A_1(0.2)$ OR

$y \in B_1(0.5)$ then $z \in C_1$

where C_1 is a membership function defined by an isosceles trapezium.

Discrete $(0, 2)$, $[0, 1]$



C) $\mu_A(x) = 0$ if $x \leq a$

$$= \frac{x-a}{b-a} \quad \text{if } a \leq x \leq b$$

$$= 1 \quad \text{if } b \leq x \leq c$$

$$= \frac{d-x}{d-c} \quad \text{if } c \leq x \leq d$$

$$= 0 \quad \text{if } d \leq x$$

MI12019090

1:25 PM