

## LDA (Linear Discriminant Analysis)

Supervised machine learning technique

Classification (single class & multi-class) & Dimensionality Reduction

developed by R.A Fisher

### Assumptions-

- Normality → The data follows normal (or gaussian) distribution
- Homogeneity of variances → The covariance of every classes is same.
- Independence → The observations are assumed to be independently sampled

Ques → ① Suppose if we have any point  $x_i$ , dimension D

$$\begin{bmatrix} x_i^1 \\ x_i^2 \\ x_i^3 \\ \vdots \\ x_i^D \end{bmatrix}$$

$$\begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_D \end{bmatrix}$$

② And any plane →

$$w_1x^1 + w_2x^2 + w_3x^3 + \dots + w_Dx^D = 0$$

passing through origin.

③ So projection of any point

$$w^T x_i$$

projection of  $\vec{B}$  on  $\vec{A}$  is dot-product

$$(w_1, w_2, \dots, w_D) \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_D \end{bmatrix}$$

→ Objective: LDA aims to find a linear combination of features that best separates two classes of objects or events.

It projects the data onto a lower dimensional space with good class-separability.

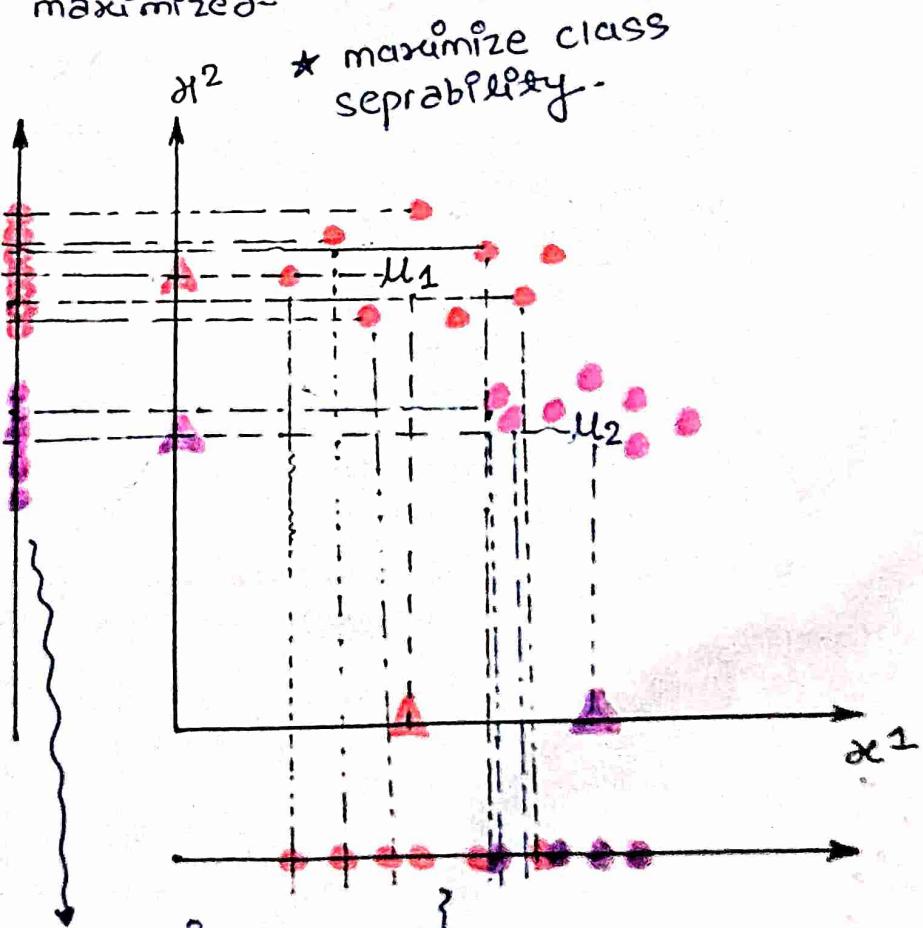
reduce complexity of eqn.

$$y = w_1x_1 + w_2x_2 + \dots + w_nx_n + w_0$$

associated weights

input features

The main task in LDA is to project the data onto a line or plane such that the distance between the points is maximized.



more precise for separation (if nearly separable).

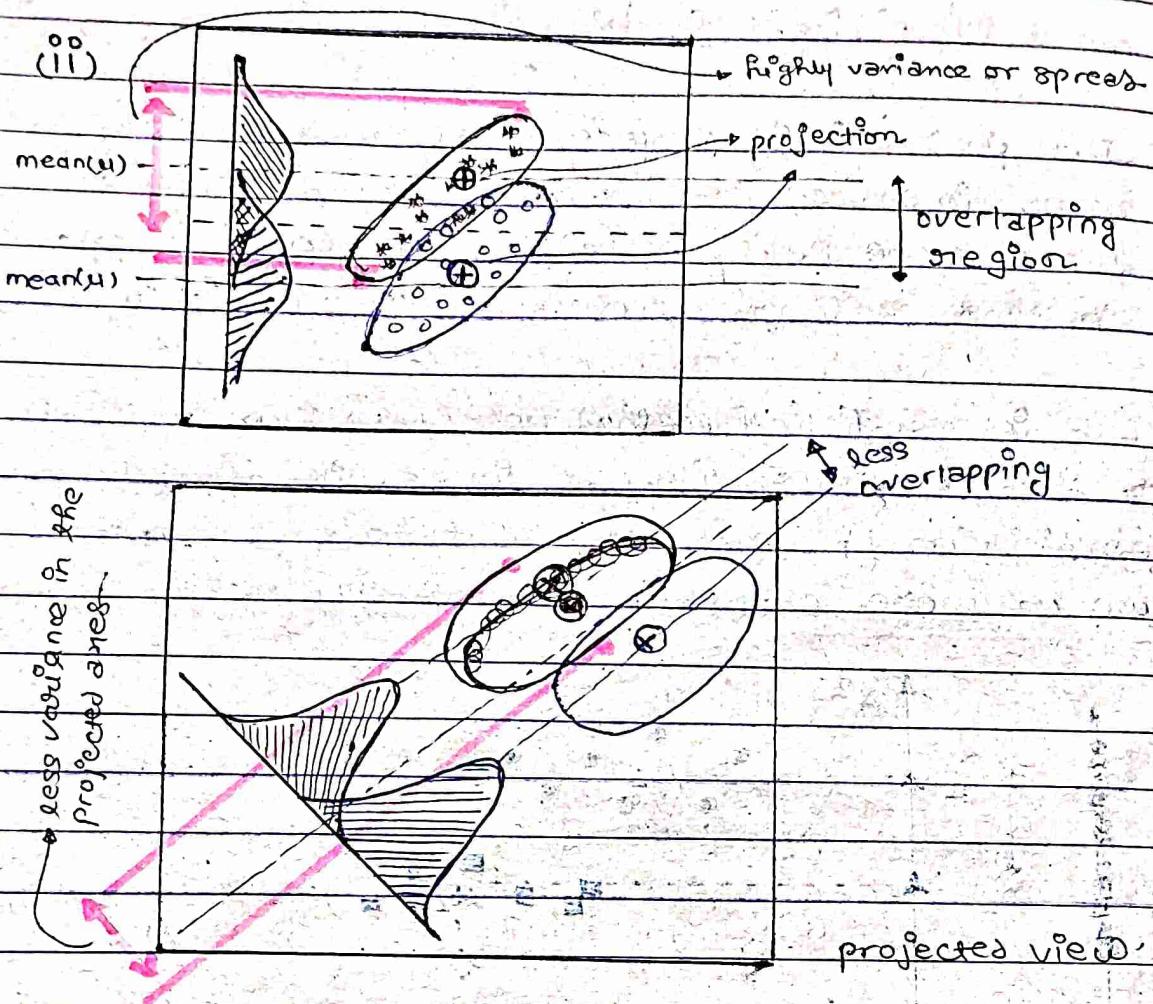
overlapping

$\mu$  → original data class mean values

$\tilde{\mu}$  → projected data class mean values

Statement 1: Modify  
be done by Incorporating  
the square of the coeff.  
Statement 2: P

Thus we want that projection mean location of different classes should be far away from each other.



So, for better projection we need  $\Rightarrow$

- The variance of the projected points of the different classes should be minimum

Outcome

$$1) \max [ \tilde{\mu}_1 - \tilde{\mu}_0 ]^2 \quad \tilde{\mu}_1 \Rightarrow \text{projected mean of class 1 points}$$

$\tilde{\mu}_0 \Rightarrow$  projected mean of class 2 points

$$2) \min \tilde{s}_1^2 \Rightarrow \text{variance of projected points on class 1}$$

$$\tilde{s}_0 \Rightarrow \text{variance of projected points on class 0}$$

location of  
from

So, we need to find  $w'$  vector  $w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}$

$$\text{Objective} \Rightarrow \max_{w \in S} \frac{(\bar{\mu}_1 - \bar{\mu}_0)^2}{S_0^2 + S_1^2}$$

or spread

Hence, LDA  
(Fisher Linear Discriminant)

To create this new axis & reduce

Dimensionality —

maximize  
reduce the distance between projected mean of two classes

→ How to find the projection of any point on a

→ Algorithm:

$x_i^{(1)}$   
Ex matrix  
here  $(6 \times 2)$

$\mu_1 \Rightarrow$

$\frac{1}{N_1} \sum_{i=1}^N x_i^{(1)2}$
$\frac{1}{N_1} \sum_{i=1}^N x_i^{(2)2}$
$\vdots$
$\frac{1}{N_1} \sum_{i=1}^N x_i^{(B)2}$

$\Omega \times 1$  matrix

$\mu_1 = w^T \mu_1$

$1 \times B$

dim 1	dim 2	class
3 $x_1^1$	4 $x_1^2$	1
2 $x_2^1$	3 $x_2^2$	1
1 $x_3^1$	1 $x_3^2$	1
5	5	0
6	6	0
7	7	0
$\bar{x}_N^{(1)}$	$\bar{x}_N^{(2)}$	

No points

6  
3  
 $\frac{8}{3}$   
 $\frac{18}{3}$

$$\mu_1 \Rightarrow \begin{bmatrix} 2 \\ 8/3 \end{bmatrix}_{2 \times 1} ; \mu_0 \Rightarrow \begin{bmatrix} 6 \\ 6 \end{bmatrix} ; \tilde{\mu}_1 - \tilde{\mu}_0 = \underbrace{w^T(\mu_1 - \mu_0)}_{\text{gap by } w \text{ the mean location in original feature space}}$$

$$(\tilde{\mu}_1 - \tilde{\mu}_0)^2 \Rightarrow \underbrace{\sum w^T w \rightarrow \frac{1}{2} \|w\|^2 \rightarrow \sqrt{w_1^2 + w_2^2 + \dots + w_B^2}}$$

$$(w^T(\mu_1 - \mu_0))(w^T(\mu_1 - \mu_0))^T \quad \downarrow (AB)^T = B^T A^T$$

$$(w^T)_{1 \times D} (\mu_1 - \mu_0)_{D \times 1} \quad (\mu_1 - \mu_0)_{D \times 1}^T w_{D \times 1}$$

original data class mean values available if training data is known.

$$(\tilde{\mu}_1 - \tilde{\mu}_0)^2 = w^T S_B w \rightarrow \text{scatter-matrix between classes}$$

between-class scatter matrix (S<sub>B</sub>) - scatter b/w the means of different classes

$$\Rightarrow (\mu_1 - \mu_0)(\mu_1 - \mu_0)^T$$

$$\begin{bmatrix} -4 \\ -10/3 \end{bmatrix} \begin{bmatrix} -4 & -10/3 \end{bmatrix} = \begin{bmatrix} 16 & 40/3 \\ 40/3 & 100/9 \end{bmatrix}$$

Class 1 :

$$\text{dim 1 : } \begin{bmatrix} 3 & 2 & 1 \end{bmatrix} = x_1$$

$$\text{dim 2 : } \begin{bmatrix} 4 & 3 & 1 \end{bmatrix} = x_0$$

separate  $N_1$  &  $N_0$  points

$S_1^2 \rightarrow$  variance of projected points of class 1.

$$\frac{1}{N_1} \sum_{i=1}^{N_1} (w^T x_i^0 - \tilde{\mu}_1)^2$$

↑  
only known parameter

$$\Rightarrow (w^T x_i^0 - w^T \mu_1) \rightarrow \text{known}$$

$$\Rightarrow w^T (x_i^0 - \mu_1)$$

$$w^T (x_1 - \mu_1) \rightarrow \partial x N_1$$

↓  
 $1 \times B$

$$\begin{bmatrix} w^T (x_1 - \mu_1) & w^T (x_2 - \mu_2) & \dots & w^T (x_n - \mu_n) \end{bmatrix}$$

$$S_1^2 \Rightarrow [w^T (x_1 - \mu_1)] [w^T (x_1 - \mu_1)]^T$$

similar

$$S_0^2 \Rightarrow [w^T (x_0 - \mu_0)] [w^T (x_0 - \mu_0)]^T$$

$$S_1^2 + S_0^2 \Rightarrow w^T (x_1 - \mu_1) [x_1 - \mu_1]^T w + w^T (x_0 - \mu_0) [x_0 - \mu_0]^T w$$

$$\Rightarrow w^T [(x_1 - \mu_1) (x_1 - \mu_1)^T + (x_0 - \mu_0) (x_0 - \mu_0)^T] w$$

↓

$S_w$

(Scatter matrix within-class)

measure the scatter (variance) within each class

→ from eqn - ①,

$$\frac{\partial A}{\partial \omega} = \frac{A}{B} \frac{\partial B}{\partial \omega}$$

$$S_B \omega = \left( \frac{A}{B} \right) S_w \omega$$

constant

$$S_w^{-1} S_B \omega = \left( \frac{A}{B} \right) \omega$$

const

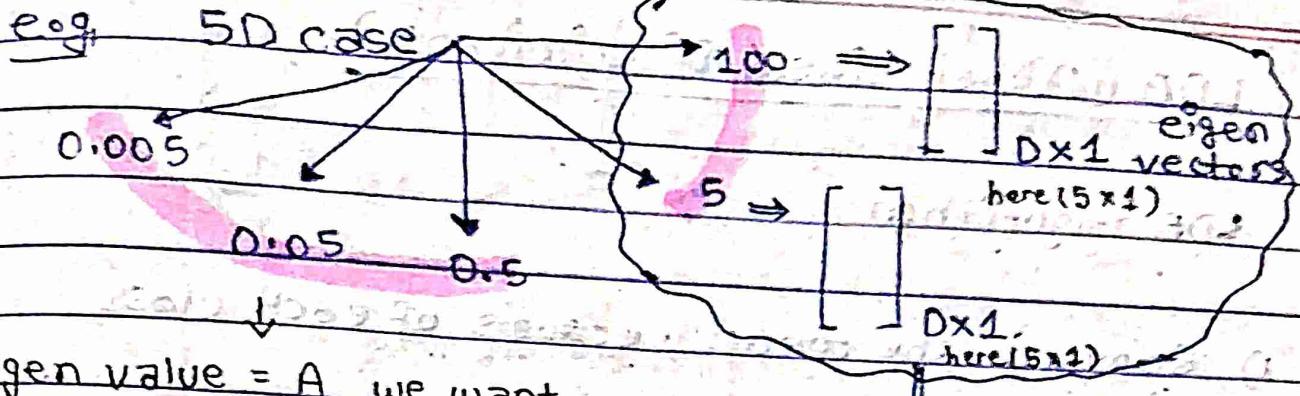
$Ax = \lambda x$  where  $x$  is eigen vector  
&  $\lambda$  is eigen value

$$\max \frac{\omega^T S_B \omega}{\omega^T S_w \omega} \Rightarrow S_w^{-1} S_B \omega = \frac{A}{B} \omega$$

known      eigen value  
eigen vector

- we have data of  $D$  dimension
- $S_w, S_B$  can be found  $\Rightarrow D \times D$  (size)
- So  $\underbrace{S_w^{-1} S_B}_{D \times D} \omega = \frac{A}{B} \omega \rightarrow$  So, we can get  $B$  nos of eigen values

we will find eigen vectors for each & every corresponding eigen value.



eigen value =  $\frac{A}{B}$  we want

to maximize, hence, doesn't take lower values.

significant or important

matrix containing eigen vector from both

$$\begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \quad D \times 2$$

original datapoint  $(x_i^o) \Rightarrow$

$$\begin{bmatrix} \cdot & \cdot & \cdot & \cdot & \cdot \end{bmatrix} \quad D \times 1$$

here  $(5 \times 1)$

$$\xrightarrow{\text{map}} w^T x_i^o$$

$$\begin{bmatrix} \cdot & \cdot & \cdot & \cdot & \cdot \end{bmatrix} \quad 2 \times D$$

here  $(2 \times 5)$

$D \times 1$

corresponding

$$\begin{bmatrix} \cdot \end{bmatrix} \quad 2 \times 1$$

dimensionality reduction

$$S_1^2 + S_2^2 \Rightarrow [(x_1 - \mu_1)(x_1 - \mu_1)^T + (x_0 - \mu_0)(x_0 - \mu_0)^T]$$

$$\begin{bmatrix} 1 & 0 & -1 \\ 4/3 & 1/3 & -5/3 \end{bmatrix} \begin{bmatrix} 1 & 4/3 \\ 0 & 4/3 \\ -1 & -5/3 \end{bmatrix} \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix} \begin{bmatrix} -1 & 1 \\ 0 & 0 \\ 1 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 2 & 3 \\ 3 & 14/3 \end{bmatrix} \quad \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix} \rightarrow \begin{bmatrix} 4 & 5 \\ 5 & 20/3 \end{bmatrix}$$

↓  
SW.

$$J(r) = \max \left( \frac{\omega^T S_B \omega}{\omega^T S_W \omega} \right) \rightarrow \max [r_1 - r_0]^2 \rightarrow \text{maximize the distance between projected mean of two classes}$$

$$\max \left( \frac{\omega^T S_B \omega}{\omega^T S_W \omega} \right) \rightarrow \min [S_1^2 + S_0^2] \rightarrow \text{reduce the variance by the projected mean of two classes}$$

$$\frac{d}{dw} \left( \frac{\omega^T S_B \omega}{\omega^T S_W \omega} \right) \rightarrow \frac{d}{dw} (\omega^T S_B \omega) \Rightarrow 2 S_B \omega$$

$$\rightarrow \underline{\frac{d}{dw} (\omega^T S_W \omega)} \Rightarrow 2 S_W \omega$$

$$\text{So, } S\omega^T S_B \Rightarrow \begin{bmatrix} 16 & 4013 \\ 4013 & 10019 \end{bmatrix} \xrightarrow{\begin{array}{l} \text{Row 1} - 4 \times \text{Row 2} \\ \text{Row 2} - 5 \times \text{Row 1} \end{array}} \begin{bmatrix} 1 & 0 \\ 0 & 40 \times 20 - 5 \times 5 \end{bmatrix} \quad \boxed{P6 20/3 - 5}$$

$$(S\omega^T S_B) \omega = \lambda \omega$$

$$\begin{pmatrix} a & b \\ c & d \end{pmatrix} \xrightarrow{\begin{array}{l} \text{Row 1} - \text{Row 2} \\ \text{Row 2} - \frac{c}{a} \text{ Row 1} \end{array}} \begin{pmatrix} a & b \\ 0 & ad - bc \end{pmatrix}$$

$$\frac{3}{5} \begin{bmatrix} 40 & 100/3 \\ -\frac{80}{3} & -200/9 \end{bmatrix} \omega = \lambda \omega$$

eigen value &  
eigen vector

$$AX = \lambda X$$

$$AX - \lambda X = 0$$

$$(A - \lambda I) X = 0$$

eigen vector  $X$  are calculated by solving this eqn

for finding eigen value of  $A$

$$\Rightarrow |A - \lambda I| = 0$$

## LDA (Fisher Linear Discriminant)

### LDA Algorithm

- 1) Compute the mean vectors of each class
- 2) Compute the scatter matrices  $S_w$  &  $S_B$ .
- 3) Compute the eigen vectors & eigen values for the matrix  $S_w^{-1} S_B$ .
- 4) Sort the eigen vectors by decreasing eigen values and choose the top k eigen vectors to form a matrix  $W$ .
- 5) Transform the samples to new space using  $W$ .

### Testing

So first transform  
the point  $x_t$

$$\rightarrow (W^T x_t)$$

Now, find distance  
of transformed  
point from the  
 $\sim \sim$   
 $M_1, M_0$ .

So, If dist from  $M_1 <$  dist  
from  $M_0$

→ Class 1

else class 0.

## questions

Q) What is the primary goal of LDA?

(A) To maximize the variance of data → PCA

(B) To minimize the covariance b/w classes.

(C) To maximize the separability b/w different classes. ✓

(D) To find non-linear boundary b/w classes.

Q) Assumption of LDA?

Q) In LDA, the Within-Class Scatter Matrix ( $S_W$ ) measures?

(A) The variance within each class →  $S_W = S_1^2 + S_0^2$

(B) The variance b/w the classes.  $(x_1 - \mu_1)(x_1 - \mu_1)^T + (x_0 - \mu_0)(x_0 - \mu_0)^T$

Q) Key diff b/w LDA & PCA?

LDA maximizes class separability  
PCA maximizes variance.

Q) In contrast to LDA, primary goal of eigen value & eigen vector.

Transform the data into new space with reduced dimensions.

Q) LDA is particularly effective when -

(A) The classes are linearly separable ✓

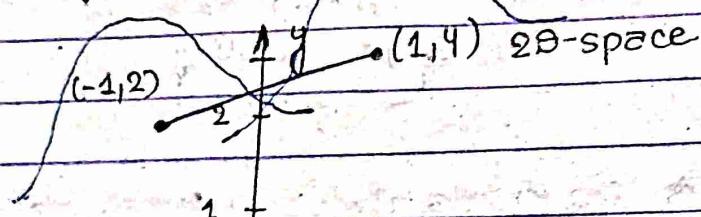
(B) The classes are not linearly separable )

(C) There are more features than samples )

(D) The data follow a non-gaussian distribution )

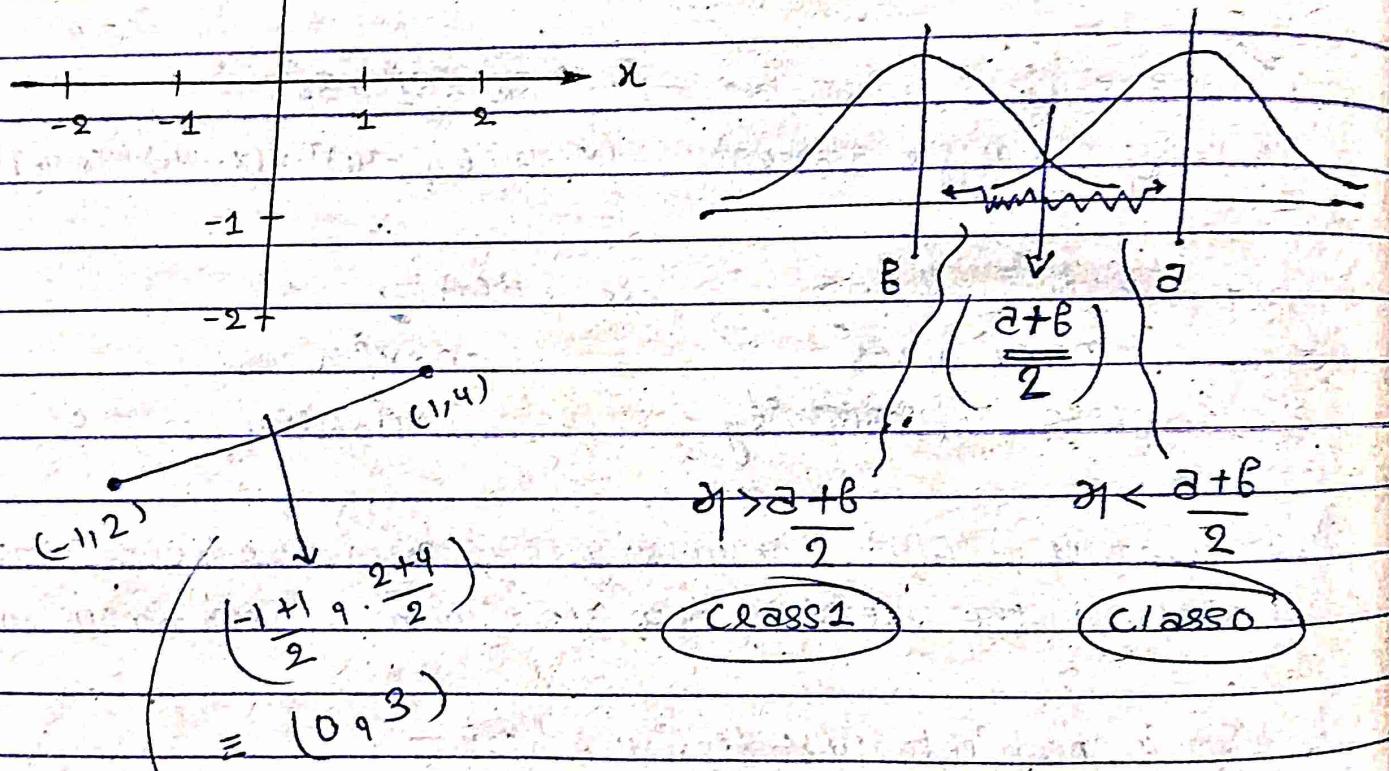
Q) Consider the case where two classes follows gaussian distribution which are centered at  $(-1, 2)$  &  $(1, 4)$  and have identity covariance matrix. Which of the following is the separating decision boundary?

- (A)  $y - x = 3$
- (B)  $x + y = 3$
- (C)  $x + y = 6$
- (D) (B) & (C) are possible



A. Not on gaussian distribution

C(kewness) same  $\Rightarrow$



classifier should be  
perpendicular to this  
line passing through  
this point  $(0, 3)$

$$\therefore \text{so slope} = -1 \quad y - y_1 = m(x - x_1)$$

$$\therefore \text{so line } \therefore y = -x + 3$$

Q. In LDA, what is the effect of class priors on decision boundary?

- (a) class priors have no effect on the decision boundary
- (b) Higher prior shift decision boundary towards that class
- (c) Lower prior shift the decision boundary towards that class
- (d) class prior make the decision boundary curved

