

-: Module 2 :-

Introduction:-

- The goal in classification is to take an ^{Input} vector $x \in \mathbb{R}^D$ and assign it to one of K discrete classes C_k .
- The input space is divided into decision regions whose boundaries are called as decision boundaries.
- As we are discussing linear classification, it means that the decision surfaces are linear functions of Input vector $x \in \mathbb{R}^D$ are defined by $(D-1)$ dimensional hyperplanes within the D -dimensional Input space.
- There are 3 distinct approaches to the classification problems & the simplest involves constructing a discriminant function which directly assigns each vector ' x ' to a specific class.

Discriminant functions:-

- A discriminant function takes an \mathbb{R}^p vector ' x ' and assigns it to one of K classes denoted as C_k .

→ Here we assume linear discriminations in which decision surfaces are hyperplanes.

Two classes :-

Simplest representation of a linear discriminant function is

$$y(x) = w^T x + w_0$$

w → weight vector

w₀ → Bias

→ If vector x is assigned to class C₁ if y(x) ≥ 0
C₂ if ~~y(x)~~ otherwise.

→ The decision boundary is corresponding to a (D-1) dimensional hyperplane within the D-dimensional input space.

→ Consider two points x_A & x_B both lie on the decision surfaces. Because y(x_A) = y(x_B) = 0 we have w^T(x_A - x_B) = 0 & hence 'w' vector is orthogonal to every vector lying on the decision surface, ∴ 'w' determines the orientation of the decision surface.

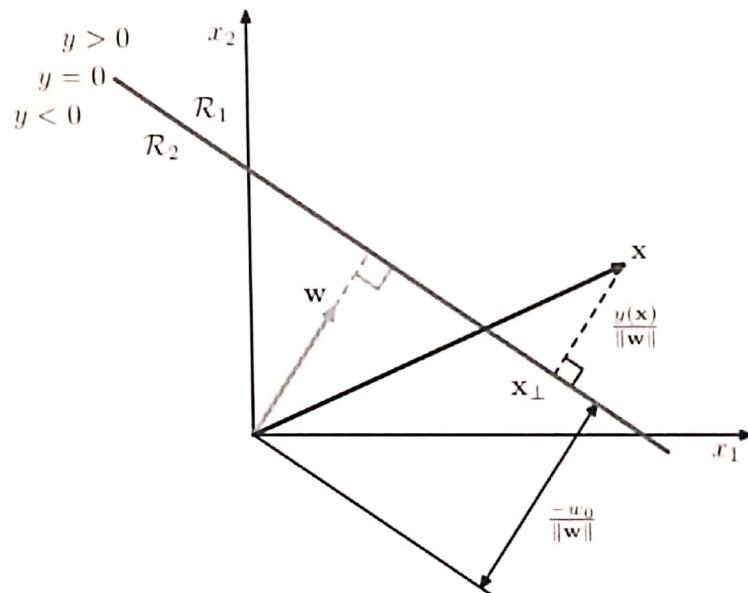
→ The normal distance from the Origin to
the decision surface is given by ③

$$\frac{\mathbf{w}^T \mathbf{x}}{\|\mathbf{w}\|} = -\frac{w_0}{\|\mathbf{w}\|}$$

→ Consider the diagram shown below.

182 4. LINEAR MODELS FOR CLASSIFICATION

Figure 4.1 Illustration of the geometry of a linear discriminant function in two dimensions. The decision surface, shown in red, is perpendicular to \mathbf{w} , and its displacement from the origin is controlled by the bias parameter w_0 . Also, the signed orthogonal distance of a general point x from the decision surface is given by $y(x)/\|\mathbf{w}\|$.



→ Here $y(x)$ gives a signed measure of the perpendicular distance r of the point x from the decision surface.

→ For Example let x be the arbitrary point if x_\perp be its orthogonal projection onto the decision surface.

$$x = x_\perp + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

→ multiplying both sides with w^T & adding (4)
 w_0 we get -

$$x \cdot w^T + w_0 = \left[x \perp + r \frac{w}{\|w\|} \right] w^T + w_0$$

* using $y(x) = w^T x + w_0$ ~~is zero~~

$$y(x \perp) = w^T x \perp + w_0 = 0$$

we get

$$\cancel{y(x)} = r \boxed{r = \frac{y(x)}{\|w\|}}$$

$$\boxed{y(x) = \tilde{w}^T \tilde{x}}$$

ii) Multiple classes :-

→ Consider Extension of linear discriminations
to $k > 2$ classes.

→ Consider the use of $k-1$ classifiers, each
of which solves a two class problem of
separating points into a particular class
 C_k from points not in that class.

→ This is known as a one-versus-the
rest classifier. (5)

→ Fig below shows an example of 3 classes with ambiguous classification.

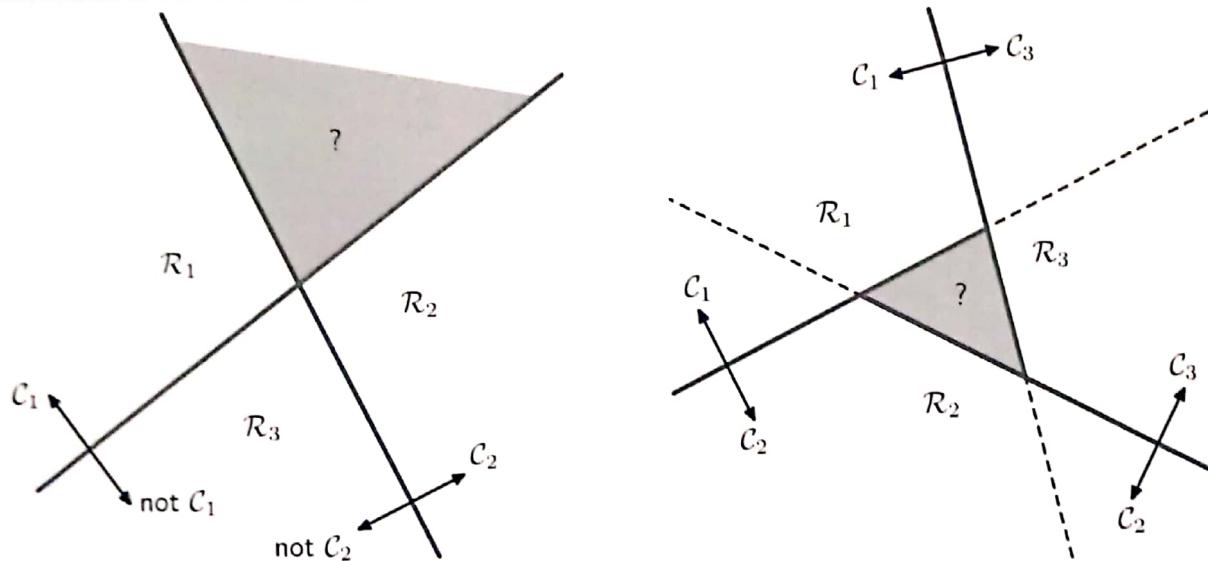


Figure 4.2 Attempting to construct a K class discriminant from a set of two class discriminants leads to ambiguous regions, shown in green. On the left is an example involving the use of two discriminants designed to distinguish points in class C_k from points not in class C_k . On the right is an example involving three discriminant functions each of which is used to separate a pair of classes C_k and C_j .

→ An alternate is to introduce $\frac{K(K-1)}{2}$ binary discriminant functions, one for every possible pair of classes. This is known as a one-versus-one classifier.

→ Further illustrating in Equations.

$$y_K(x) = w_K^T x + w_K^0$$

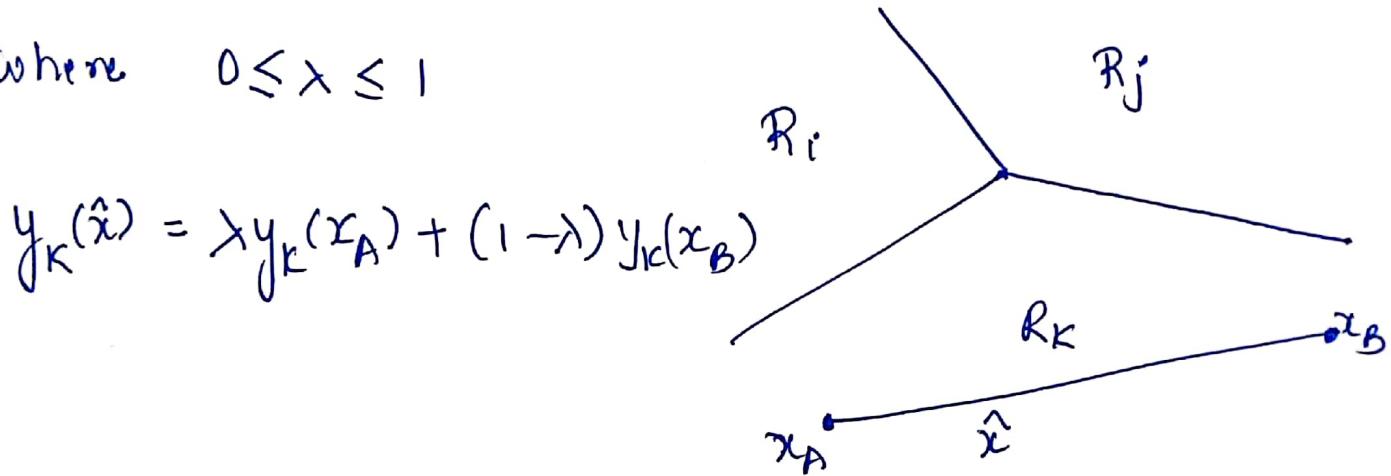
$$\{ (w_K - w_j)^T x + (w_{K0} - w_{j0}) = 0$$

where k & j are two classes C_k & C_j . (6)

→ Consider two points x_A & x_B both of which lie inside decision region R_k . Any point \hat{x} that lies on the line connecting x_A & x_B can be expressed in the form,

$$\hat{x} = \lambda x_A + (1-\lambda)x_B.$$

where $0 \leq \lambda \leq 1$



Note :- Both x_A & x_B lie inside R_k . R_k is singly connected & convex.

Least Squares for classification :-

→ We know that minimization of sum-of-squares error function led to a simple closed-form solution for the parameter values in regression, so that can be applied for classification.

→ Each class C_k described by its linear model

$$y_k(x) = \tilde{w}_k^T x + w_{k0}$$

where $k = 1, \dots, K$, grouping these together using vector notation.

$$y(x) = \tilde{w}^T x$$

→ we now determine parameter matrix \tilde{w} by minimizing sum-of-squares error function.

Error function is defined as

$$E_D(\tilde{w}) = \frac{1}{2} \text{Tr} \{ (\tilde{x}\tilde{w} - T)^T (\tilde{x}\tilde{w} - T) \}$$

* Setting the derivative with respect to \tilde{w} to zero. & rearranging we obtain

$$\tilde{w} = (\tilde{x}^T \tilde{x})^{-1} \tilde{x}^T T = \tilde{x}^+ T$$

\downarrow
Pseudo-Inverse of
the matrix \tilde{x} .

Note :- Sum-of-squares error function penalizes predictions that are 'too correct' & which lie along the decision boundary.

Principle Component Analysis :-

- It is a way of Identifying pattern in data and Expressing the data in such a way to highlight their similarities & differences
- It is also way to reduce the dimension [Dimensionality Reduction] to reduce the complexity of the Analysis .

Example problem :-

- Given the following data use PCA to reduce the dimension from 2 to 1 .

Feature	1	2	3	4
x	4	8	13	7
y	11	4	5	14

Step 1 :- From the data set find n & N

$$n = \text{no. of features} = 2$$

$$N = \text{no. of samples} = 4$$

Step 2 :- Computation of mean of variables

$$\bar{x} = \frac{4+8+13+7}{4} = 8$$

$$\bar{y} = \frac{11+4+5+14}{4} = 8.5$$

Step 3 :- Computation of covariance matrix

i) find ordered pairs $\rightarrow x, y$

$(x, x) (x, y) (y, x) (y, y)$

ii) find covariance of all ordered pairs

$$\text{Cov}(x, x) = \frac{1}{N-1} \sum_{k=1}^N (x_{ik} - \bar{x}_i)(x_{jk} - \bar{x}_j)$$

* or

$$\Rightarrow \text{Cov}(x, x) = \frac{1}{N-1} \sum_{k=1}^N (x_i - \bar{x})^2$$

$$= \frac{1}{4-1} \left[(4-8)^2 + (8-8)^2 + (13-8)^2 + (7-8)^2 \right]$$

$$= 14$$

$$\Rightarrow \text{Cov}(x, y) = \frac{1}{N-1} \left[(4-8)(11-8.5) + (8-8)(8-8.5) + (13-8)(5-8.5) + (7-8)(14-8.5) \right]$$

$$= -11$$

$$\text{cov}(y, x) = \text{cov}(x, y) = -11$$

$$\begin{aligned}\text{cov}(y, y) &= \frac{1}{4-1} \left[(11-8.5)^2 + (4-8.5)^2 + (5-8.5)^2 + (14-8.5)^2 \right] \\ &= 23\end{aligned}$$

Covariance matrix :-

$$\begin{aligned}S &= \begin{bmatrix} \text{cov}(x, x) & \text{cov}(x, y) \\ \text{cov}(y, x) & \text{cov}(y, y) \end{bmatrix} \\ &= \begin{bmatrix} 14 & -11 \\ -11 & 23 \end{bmatrix}\end{aligned}$$

Step 4 :- Eigen value , Eigen vector , normalized eigen vector .

i) Eigen value :-

$$\det(S - \lambda I) = 0$$

$$\text{Identity matrix } I = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \times \lambda = \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix}$$

$$S = \begin{bmatrix} 14 & -11 \\ -11 & 23 \end{bmatrix}$$

$$\therefore \det(S - \lambda I) = 0$$

$$\det \begin{bmatrix} 14-\lambda & -11 \\ -11 & 23-\lambda \end{bmatrix} = 0$$

$$(14-\lambda)(23-\lambda) - (-11 \times -11) = 0$$

$$\lambda^2 - 37\lambda + 201 = 0$$

$$\therefore \lambda_1 = 30.3849 \quad \left[\lambda_1 > \lambda_2 \right]$$

$$\lambda_2 = 6.6151$$

* pick largest value which becomes first principal component.

ii) Eigen vector of λ_1

$$(S - \lambda_1 I) U_1 = 0$$

\downarrow \hookrightarrow Eigen vector of λ_1

$$\begin{bmatrix} 14-\lambda_1 & -11 \\ -11 & 23-\lambda_1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = 0$$

$$\begin{bmatrix} (14-\lambda_1)u_1 - 11u_2 \\ -11u_1 + (23-\lambda_1)u_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$(14-\lambda_1)u_1 - 11u_2 = 0 \rightarrow (1)$$

$$-11u_1 + (23-\lambda_1)u_2 = 0 \rightarrow (2)$$

* use any one linear equation to find values of eigen vector

(12)

∴ Using Eqn (2)

$$(14 - \lambda_1) u_1 - 11 u_2 = 0$$

$$\frac{u_1}{11} = \frac{u_2}{14 - \lambda_1} = t$$

when $t = 1$

$$u_1 = 11, \quad u_2 = 14 - \lambda_1$$

Eigen vector u_1 of $\lambda_1 = \begin{bmatrix} 11 \\ 14 - \lambda_1 \end{bmatrix}$

$$= \begin{bmatrix} 11 \\ 14 - 30.3849 \end{bmatrix} = \begin{bmatrix} 11 \\ -16.3849 \end{bmatrix}$$

* Normalizing Eigen vector u_1

$$e_1 = \begin{bmatrix} \frac{11}{\sqrt{11^2 + (-16.3849)^2}} \\ \frac{-16.3849}{\sqrt{11^2 + (-16.3849)^2}} \end{bmatrix} = \begin{bmatrix} 0.5574 \\ -0.8303 \end{bmatrix}$$

$$11^{th} e_2 = \begin{bmatrix} 0.8303 \\ 0.5574 \end{bmatrix}$$

V

Step 5 :- Define new dataset-

Frist principle component	1	2	3	4
	P_{11}	P_{12}	P_{13}	P_{14}

$$P_{11} = e_1^T \begin{bmatrix} x_i - \bar{x} \\ y_i - \bar{y} \end{bmatrix} \quad \text{and} \quad e_1 = e_1^T \begin{bmatrix} 4 - 8 \\ 11 - 8.5 \end{bmatrix}$$

$$= \begin{bmatrix} 0.5574 & -0.8303 \end{bmatrix} \begin{bmatrix} -4 \\ 2.5 \end{bmatrix} = -4.3052$$

$$P_{12} = \begin{bmatrix} 0.5574 & -0.8303 \end{bmatrix} \begin{bmatrix} 8 & -8 \\ 4 & -8.5 \end{bmatrix} = 3.73$$

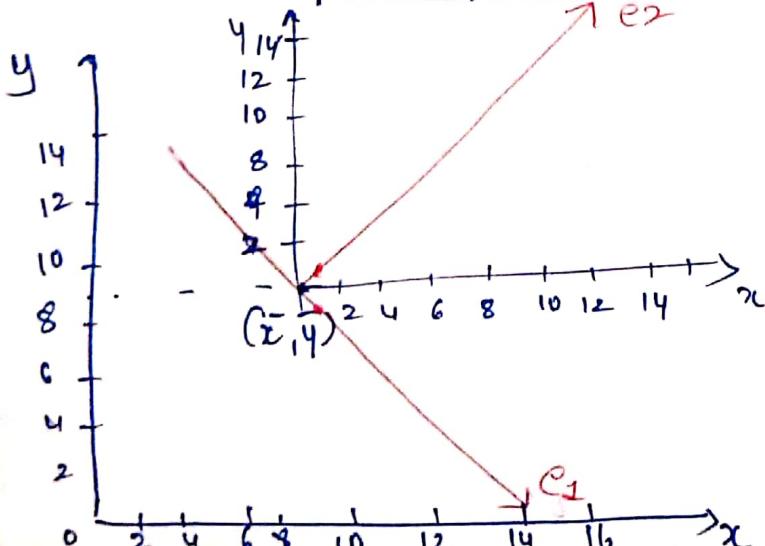
111th

$$P_{13} = 5.69$$

$$P_{14} = -5.12$$

PC1	1	2	3	4
	-4.3	3.7	5.6	-5.12

\Rightarrow new data set
with dimension 1



Linear discriminant analysis :-

(14)

- LDA is a dimensionality reduction technique used as a pre-processing step for pattern-classification & Machine learning applications.
- LDA is similar to PCA but LDA in addition finds the axes that maximizes the separation between multiple classes.

Goal :- To project a feature space (n -dimensional data) onto a smaller subspace K ($K \leq n-1$) while maintaining the class discriminating information

Example :- Let's take a 2-D dataset -

$$C_1 \Rightarrow x_1 = (x_1, x_2) = \{(4,1), (2,4), (2,3), (3,6), (4,4)\}$$

$$C_2 \Rightarrow x_2 = (x_1, x_2) = \{(9,10), (6,8), (9,5), (8,7), (10,8)\}$$

Step 1 :- Compute within-class scatter matrix (S_w)

$$S_w = S_1 + S_2$$

$S_1 \Rightarrow$ Covariance matrix for class C_1

$S_2 \Rightarrow$ Covariance matrix for class C_2 .

i) finding S_1

$$S_1 = \sum_{x \in C_1} (x - \mu_1)(x - \mu_1)^T$$

$\mu_1 \rightarrow$ mean of class C_1

$$\mu_1 = \left\{ \frac{4+2+2+3+4}{5}, \frac{1+4+3+6+4}{5} \right\}$$

$$\mu_1 = [3.00 \quad 3.60]$$

$$\text{Similarly } \mu_2 = [8.4 \quad 7.60]$$

$$\therefore (x - \mu_1) = \begin{bmatrix} 1 & -1 & -1 & 0 & 1 \\ -2.6 & 0.4 & -0.6 & 2.4 & 0.4 \end{bmatrix}$$

Now, for each x we calculate $(x - \mu_1)(x - \mu_1)^T$.

$$\therefore \begin{bmatrix} 1 \\ -2.6 \end{bmatrix} \begin{bmatrix} 1 & -2.6 \end{bmatrix} = \begin{bmatrix} 1 & -2.6 \\ -2.6 & 6.76 \end{bmatrix} \rightarrow (1)$$

$$\begin{bmatrix} -1 \\ 0.4 \end{bmatrix} \begin{bmatrix} -1 & 0.4 \end{bmatrix} = \begin{bmatrix} 1 & -0.4 \\ -0.4 & 0.16 \end{bmatrix} \rightarrow (2)$$

$$\begin{bmatrix} -1 \\ -0.6 \end{bmatrix} \begin{bmatrix} -1 & -0.6 \end{bmatrix} = \begin{bmatrix} 1 & 0.6 \\ 0.6 & 0.36 \end{bmatrix} \rightarrow (3)$$

$$\begin{bmatrix} 0 \\ 2.4 \end{bmatrix} \begin{bmatrix} 0 & 2.4 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0 & 5.76 \end{bmatrix} \rightarrow (4)$$

$$\begin{bmatrix} 1 \\ 0.4 \end{bmatrix} \begin{bmatrix} 1 & 0.4 \end{bmatrix} = \begin{bmatrix} 1 & 0.4 \\ 0.4 & 0.16 \end{bmatrix} \rightarrow (5)$$

(14)

Adding (2) + (2) + (3) + (4) + (5) & taking average
we get covariance matrix S_1

$$S_1 = \begin{bmatrix} \frac{4}{5} & -\frac{2}{5} \\ -\frac{2}{5} & \left(\frac{13}{2}\right) \end{bmatrix} = \begin{bmatrix} 0.8 & -0.4 \\ -0.4 & 6.5 \end{bmatrix}$$

Similarly for the class 2, the co-variance matrix is given by.

$$S_2 = \begin{bmatrix} 1.84 & -0.04 \\ -0.04 & 2.64 \end{bmatrix} \text{ & } \mu_2 = \begin{bmatrix} 8.4 & 7.6 \end{bmatrix}$$

$$S_w = S_1 + S_2$$

$$S_w = \begin{bmatrix} 2.64 & -0.44 \\ -0.44 & 5.28 \end{bmatrix}$$

Step 2:- Compute Between class scatter matrix (S_B)

$$S_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T$$

$$\mu_1 = \begin{pmatrix} 3 & 3.6 \end{pmatrix} \quad \mu_2 = \begin{pmatrix} 8.4 & 7.6 \end{pmatrix}$$

$$S_B = \begin{bmatrix} -5.4 \\ -4 \end{bmatrix} \begin{bmatrix} -5.4 & -4 \end{bmatrix} = \begin{bmatrix} 29.16 & 21.6 \\ 21.6 & 16.00 \end{bmatrix}$$

Step 3:- Find the best LDA projection vector

* Similar to principal component analysis we find this using eigen vectors using largest Eigen value.

$$S_w^{-1} S_B \cdot w = \lambda \cdot w \quad \xrightarrow{\text{Projection vector}} \quad (1)$$

i) Solving the eigen value

$$|S_w^{-1} S_B - \lambda I| = 0$$

$$S_w^{-1} = \begin{bmatrix} 2.64 & -0.44 \\ -0.44 & 5.28 \end{bmatrix}^{-1} = \begin{bmatrix} 0.38 & 0.03 \\ 0.03 & 0.19 \end{bmatrix}$$

$$S_B = \begin{bmatrix} 29.16 & 21.6 \\ 21.6 & 16.00 \end{bmatrix}$$

$$\lambda I = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \times \lambda = \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix}$$

$$S_w^{-1} \cdot S_B = \begin{bmatrix} 11.72 & 8.6 \\ 4.9 & 3.6 \end{bmatrix}$$

$$\therefore |S_w^{-1} S_B - \lambda I| = \begin{vmatrix} 11.72 - \lambda & 8.6 \\ 4.9 & 3.6 - \lambda \end{vmatrix} = 0$$

Solving for λ

$$(11.72 - \lambda)(3.6 - \lambda) - (4.9 \times 8.6) = 0$$

$$42.19 - 11.72\lambda - 3.6\lambda + \lambda^2 - 42.14 = 0$$

$$\lambda^2 - 15.32\lambda + 0.05 = 0.$$

$$\therefore \lambda_1 = 15.31, \lambda_2 = 3.26 \times 10^{-3}$$

Substituting highest Eigen value in Eqn (1)

$$\begin{bmatrix} 11.72 & 8.6 \\ 4.9 & 3.6 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = 15.31 \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

~~$$11.72w_1 + 8.6w_2 = 15.31w_1 \rightarrow (2)$$~~

$$4.9w_1 + 3.6w_2 = 15.31w_2 \rightarrow (3)$$

$$11.72w_1 - 15.31w_1 + 8.6w_2 = 0$$

$$-3.59w_1 + 8.6w_2 = 0 \rightarrow (4)$$

$$4.9w_1 + 3.6w_2 - 15.31w_2 = 0$$

$$4.9w_1 - 11.7w_2 = 0$$

$$\begin{bmatrix} -3.59 & 8.6 \\ 4.9 & -11.7 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$-3.59w_1 + 8.6w_2 = 0$$

$$w_1 = \frac{8.6}{3.59}w_2 = 2.39w_2$$

$$4.9w_1 + 11.7w_2 = 0$$

$$\text{Let } w_2 = 1$$

$$\therefore w_1 = \frac{11.7}{4.9} = 2.38$$

$$\begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 2.38 \\ 1 \end{bmatrix}$$

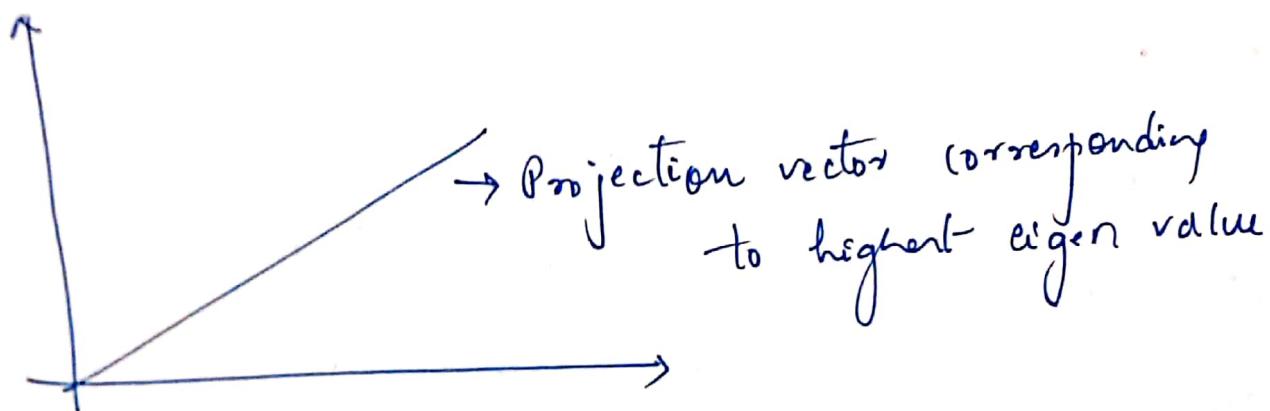
normalizing

$$\begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 2.38 \\ \sqrt{1^2 + 2.38^2} \end{bmatrix}$$

$$\therefore \text{we get } \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 0.91 \\ 0.39 \end{bmatrix}$$

Step 4 :- Dimension Reduction

$y = w^T x \rightarrow$ Input data sample.
 \downarrow
 Projection vector



The perceptron algorithm :-

→ Here let us consider two-class model in which input vector x is first transformed using a fixed nonlinear transformation to give a feature vector $\phi(x)$ & then construct a generalized linear model

$$y(x) = f(w^T \phi(x)).$$

where the non-linear activation function $f(\cdot)$ is a step function of form

$$f(a) = \begin{cases} +1 & a > 0 \\ -1 & a \leq 0 \end{cases}$$

- In Perception algorithm it convenient to use $t = +1$ for c_1 & $t = -1$ for c_2 as target values matching the choice of activation function.
- To determine the parameters w correctly we can minimize error function but alternatively we can also use perception criterion.

- To derive perception criterion based Error we 1st consider

pattern $x_n \rightarrow c_1 \Rightarrow$ when $w^T \phi(x_n) > 0$
 $x_n \rightarrow c_2 \Rightarrow$ when $w^T \phi(x_n) < 0$

$$t \in \{-1, +1\}$$

note :- v satisfy condition $w^T \phi(x_n) t_n > 0$
 All patterns must

→ perception criterion is given by

$$\mathcal{E}_p(w) = - \sum_{n \in M} w^T \phi_n t_n$$

$M \rightarrow$ denotes the set of misclassified patterns

→ we apply gradient-descent algorithm to this error function. The change in weight is given by

$$\begin{cases} w^{(t+1)} \\ = w^{(t)} - \eta \nabla \mathcal{E}_p(w) \\ = w^{(t)} + \eta \phi_n t_n. \end{cases}$$

where $\eta \rightarrow$ learning rate

$T \rightarrow$ integer values of steps of algorithm.

Note :- As the weight are change during training the misclassified pattern will change.
set of

Probabilistic Discriminative Models :-

→ There are 2 Models

- i) Generative
- ii) discriminative .

Analogy :-

Task \rightarrow Determine the language that someone is speaking

- i) Generative approach :- Learn each language & determine as to which language the speech belongs to.
- ii) Discriminative approach :- Determine the linguistic differences without learning any language

- \rightarrow We wish to learn $f: x \rightarrow y$ eg., $P(y|x)$
- i) In generative we model the joint distribution of all data
 - ii) In discriminative we model only points at the boundary.

Note :- Discriminative models make predictions on the unseen data based on conditional probability

- \rightarrow Discriminative classification is also called as "informative".

Probability Basics :-

(23)

* Prior, conditional & joint probability for random variables.

- 1) Prior Probability - $P(x)$
- 2) Conditional probability - $P(x_1|x_2), P(x_2|x_1)$
- 3) Joint probability :- $x = (x_1, x_2) \quad P(x) = P(x_1, x_2)$
- 4) Relationship $P(x_1, x_2) = P(x_2|x_1) P(x_1)$
 $= P(x_1|x_2) P(x_2)$

5) Independence

$$P(x_2|x_1) = P(x_2), \quad P(x_1|x_2) = P(x_1), \quad P(x_1, x_2) = P(x_1)P(x_2)$$

Bayes rule :-

Generative .

$$P(c|x) = \frac{P(x|c) P(c)}{P(x)}$$

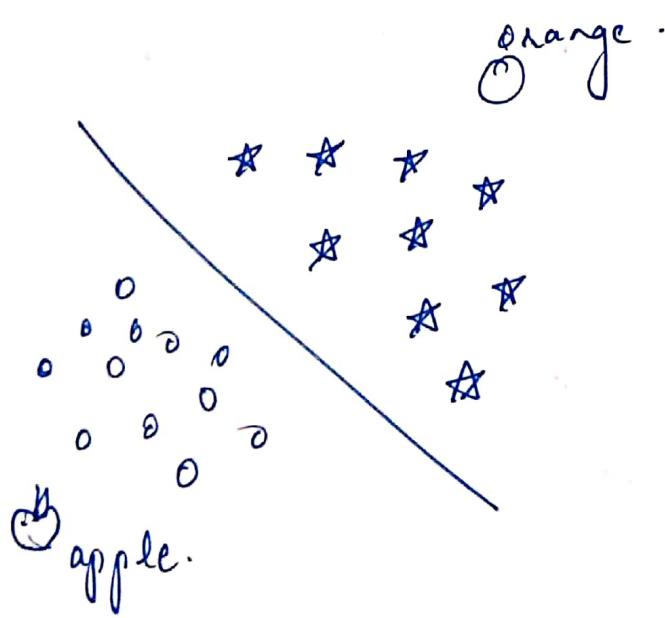
discriminative

$$\text{Posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{Evidence.}}$$

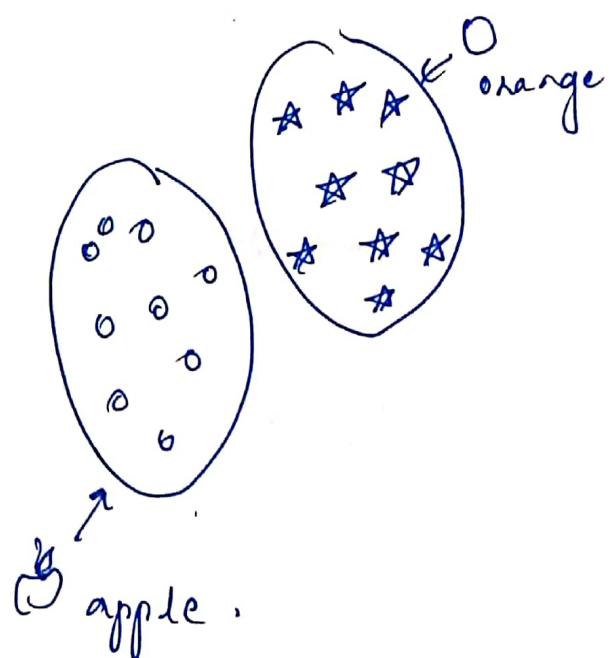
→ Given some data points, the discriminative model (24) learns to classify the data points into their respective classes by learning the decision boundary that separates the classes.

→ The generative models can also classify the given data points, but instead of learning the decision boundary, they learn the characteristics of the classes.

discriminative model



generative Model



* Above shows example of image classification, whether an input image is an apple/orange.

i) Discriminative Model learns the optimal

boundary that separates the apple & orange classes

- ii) Generative model learn their distribution by learning the characteristics of apple & orange classes.
- while both types of models are used to predict $P(y|x)$. Calculation methods for both Models are different.

i) Generative Model → first find joint probability distribution $P(x,y)$, then use Bayes' theorem to calculate $P(y|x)$

ii) discriminative Model → will learn the conditional probability $P(y|x)$ directly.

Ex :- Data set $(0,0), (0,1), (1,0), (1,1)$

$$\begin{array}{ccc}
 P(x,y) & & \\
 \begin{array}{c} x=0 \\ x=1 \end{array} &
 \begin{array}{c} y=0 \\ \hline y_1 \\ \hline y_2 \end{array} &
 \begin{array}{c} y=1 \\ \hline y_3 \\ \hline y_4 \end{array} \\
 & &
 \begin{array}{c} y_1 \Rightarrow P(x_1) \\ \hline y_2 \Rightarrow P(x_2) \end{array} \\
 & &
 \begin{array}{c} p(y)=4/8 \\ p(y)=1/8 \end{array}
 \end{array}$$

26

$$\text{and } P(y|x) = \frac{P(x,y)}{P(x)} \text{ or } \frac{P(x,y)}{P(y)}$$

$$P(x_1) = \frac{1}{2} \quad P(x_2) = \frac{1}{2}$$

$$\therefore P(y|x)$$

$x=0$	$y=0$	$y=1$
$x=1$	$\frac{1}{2}$	$\frac{1}{2}$

$$\frac{\text{perm}}{\text{nonperm}} = \frac{y_2}{y_1}$$

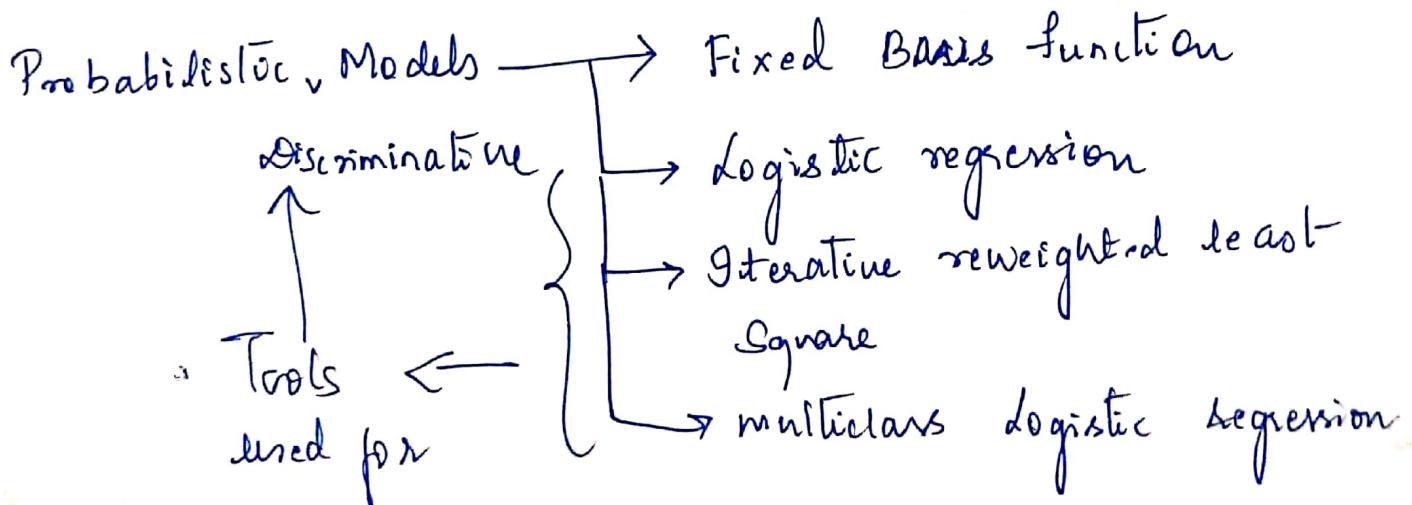
Summary :-

Generative Model

→ Works good for small data & delivers high accuracy.

Discriminative Model

- Accuracy is high
- Resource saving
- Fewer parameters to be determined.



Logistic Regression :-

(62)

→ Here we model the posterior probabilities directly assuming that they have a sigmoid-shaped distribution [without Modeling class prior & class conditional densities]

→ The sigmoid-shaped function (σ) is a model function logistic regression.

→ 1st non-linear transformation of inputs using a vector of basis function $\phi(x)$

$$P(c_1|\phi) = y(\phi) = \sigma(w^T\phi)$$

$$P(c_2|\phi) = \cancel{y(\phi)} \quad 1 - P(c_1|\phi)$$

∴ If we have M -dimensional feature space ϕ , we use $M+1$ adjustable parameters to represent

w directly.

→ Here $\sigma(\cdot) \rightarrow$ logistic sigmoid function or logistic regression.

→ We can use maximum likelihood to determine the parameters of the logistic regression model.

$$\frac{d\sigma}{da} = \sigma(1-\sigma)$$

→ For a data set $\{\phi_n, t_n\}$, $t_n \in \{0, 1\}$ & $\phi_n = \phi(x_n)$ with $n = 1 \dots N$.

Likelihood function can be written as

$$P(t|w) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n}$$

where $t = [t_1, \dots, t_N]^T$

$$y_n = P(c_1 | \phi_n) \quad \& \text{Error Equation}$$

$$\nabla E(w) = \sum_{n=1}^N (y_n - t_n) \phi_n$$

Iterative reweighted least squares :-

→ Logistic regression has no closed-form solution of $\nabla E(w) = 0$. due to non-linearity of logistic sigmoid function.

→ ∴ Error function can be minimized by an efficient iterative technique.

$$\text{ie., } w^{(\text{new})} = w^{(\text{old})} - H^{-1} \nabla E(w)$$

(20) 9

where $H \rightarrow$ Hessian matrix $\rightarrow 2^{\text{nd}}$ derivatives of $E(w)$
wrt 'w'

$$H = \nabla^2 E(w)$$

$$\nabla E(w) = \phi^T \phi w - \phi^T t.$$

$$\therefore H = \nabla \nabla E(w) = \sum_{n=1}^N \phi_n \phi_n^T = \underbrace{\phi^T \phi}_{\text{Diagonal matrix}}$$