

## Supervised Learning

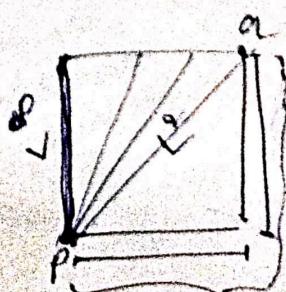
- A set of patterns with class [label] information. (i.e training pattern)
- is available
- Make a classification model with the help of training pattern.
- Assign the class label if to new pattern to a predefined classes by classifier through collection of similarity/dissimilarity measures.
- Probabilistic Model.
  - ↪ Baye's classifier.
- Similarity Based
  - ↪ MDC, KNN, LDC, QDC
- Neural network.
  - MLP, RBF, Hopfield
- Graph based decision tree
- SVM
- Random Forest

$L_p$  norm,  $p=1$ , City Block  
 $p=2$ , Euclidean distance

$p=\infty$ ,  $L^\infty \rightarrow$  supnorm

$\max_{\text{over all } i} |x_i - y_i|$

$$L_1 > L_2 > L^2 > L^\infty$$



$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

with all the features/dimensions

$n \rightarrow$  no. of features/dimensions

scoring function  $Z = \sum x_i y_i$  with others taken into account in next slide

Euclidean distance =  $\sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_n - y_n)^2}$  sum of the total units with respect to

dimensions  $\sqrt{\sum (x_i - y_i)^2}$  weighted Agarwal, explained by

if there are  $n$  features  
average distance =  $\frac{1}{n} \sum_{i=1}^n (x_i - y_i)^2$

if we consider distance between two points

$$\text{distance} = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_n - y_n)^2}$$

$$= Z^T \cdot Z = Z^T \cdot I_{n \times n} \cdot Z$$

eg.  $d^2((150, 4), (130, 3)) = \underbrace{400}_{\leq 401} + 1$

contributes more but what if first second feature was

more important give diff. weights to diff. features

$$\therefore \mathbf{d}_{\text{ew}}' = \begin{bmatrix} x_1 - y_1 & x_2 - y_2 & \dots & x_n - y_n \end{bmatrix}^T \begin{bmatrix} \lambda_1 & 0 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & \dots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \lambda_n \end{bmatrix} \begin{bmatrix} x_1 - y_1 \\ x_2 - y_2 \\ x_3 - y_3 \\ \vdots \\ x_n - y_n \end{bmatrix}$$

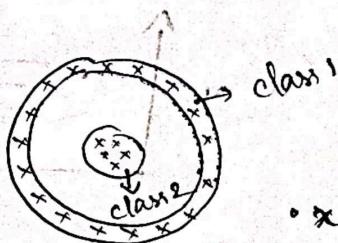
$$\mathbf{d}_{\text{ew}}' = \mathbf{z}' \mathbf{A}_{n \times n} \mathbf{z}$$

$\downarrow$   
A is symmetric positive definite Matrix.

Advantage of A

- all eigen-values  $> 0$
- $|A| > 0$
- $A^{-1}$  also be positive definite Matrix.

~~Lecture 2  
25/10/23~~  
•  $\Sigma$   $\rightarrow$  variance / covariance matrix.  
or dispersion matrix.



$x_0$  has almost equal distance from means of both class 1 and class 2

∴ we need to use second order statistics.

↳ ie variances etc.

$$\text{Var}(x) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)^2$$

calculates spread

want one variable/feature

$\Rightarrow f(x)$

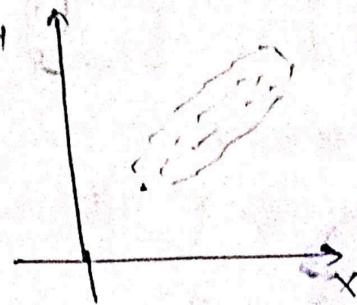
$$\text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)$$

two features

$$(f(x), f(y)) = (g(x), g(y))$$

$$(f(x), f(y)) = (g(x), g(y))$$

$$\text{Cov}(x, y) > 0$$

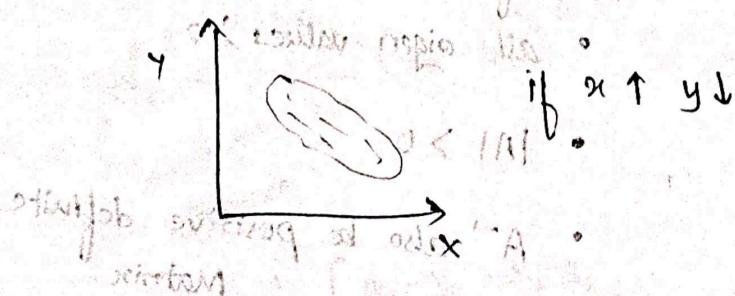


if  $x \uparrow y \uparrow$

$$\text{Span } x = b$$

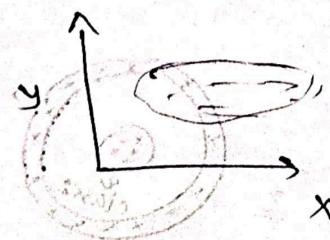
$$\text{Cov}(x, y) < 0$$

$\rightarrow$  A for negative



if  $x \uparrow y \downarrow$

stiff job surviving at odd x rotation



$$\text{Cov}(x, y) = 0$$

linear regression is zero in Z.

without regression

means not in

parallel lines

• Dispersion / Var-Cov Matrix

$$f_1(f_2, \dots, f_d) \in \mathbb{R}^{(d \times d)}$$

$$\sum = f_1 \begin{bmatrix} \text{Cov}(f_i, f_i) & \dots & \text{Cov}(f_i, f_d) \\ \vdots & \ddots & \vdots \\ \text{Cov}(f_d, f_i) & \dots & \text{Cov}(f_d, f_d) \end{bmatrix}$$

$$(f_1, f_2, \dots, f_d) \in \mathbb{R}^{(d \times 1)}$$

$(1, D \times D)$

$$\text{Cov}(f_i, f_i) = \text{Var}(f_i)$$

$$\text{Cov}(f_i, f_j) = \text{Cov}(f_j, f_i)$$

dispersion matrix is symmetric & ve definite matrix

$\Sigma_1 \rightarrow$  dispersion matrix for class 1.

How to calculate dispersion matrix easily?

$$X = \begin{bmatrix} f_1 & f_2 \\ 150 & 55 \\ 160 & 60 \\ 170 & 65 \\ \hline 160 & 60 \end{bmatrix}$$

$\Sigma = \frac{1}{3} \begin{bmatrix} f_1 & f_2 \\ 200/3 & 100/3 \\ 100/3 & 50/3 \end{bmatrix}$

(i) Mean  $\bar{x}$

$\begin{bmatrix} 20 & -7 & 0 & 2 \\ -7 & 0 & 2 & 0 \\ 0 & 2 & 0 & 2 \\ 2 & 0 & 2 & 0 \end{bmatrix}$

$\begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$

$\begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$

$\begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$

$\gamma =$  Mean subtraction of  $X$

$$= \begin{bmatrix} -10 & -5 \\ 0 & 10 \\ 10 & 5 \end{bmatrix}$$

$$\Sigma = \gamma^T \cdot \gamma$$

$$\{(4-x)^2 \leq (x-a)^2\} \Rightarrow$$
$$= \begin{bmatrix} -10 & 0 & 10 \\ -5 & 0 & 5 \end{bmatrix} \begin{bmatrix} -10 & -5 \\ 0 & 0 \\ 10 & 5 \end{bmatrix}$$

$2 \times 3 \quad 3 \times 2$

similarly  $\Sigma_2 = 1 \times 1 \times 1$

$$\Sigma_2 = \begin{bmatrix} 200 & 100 \\ 100 & 50 \end{bmatrix}$$

(NOTE:  $\Sigma_2$  is not so important)

## Classification

↳ supervised

- A training set

(con)

- The prior prob. and probability

density fn. of classes are given.

- Prior probability  $\Rightarrow$  to be in a class.

Suppose training data

Class 1 : 100	$P_1 = \frac{1}{6}$
2 : 200	$P_2 = \frac{1}{3}$
3 : 300	$P_3 = \frac{1}{2}$
<hr/>	
	600

$$\begin{bmatrix} 1 & 1 & 1 \\ 2 & 2 & 2 \\ 3 & 3 & 3 \\ 2 & 2 & 2 \end{bmatrix} = X$$

$$P_i | P(w_i)$$

$X \in \mathbb{R}^{n \times m}$ ; depending on size of class

$$\sum_i P_i = 1$$

$$1 > P_i > 0$$

- Probability density function

Gaussian

$$G = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} \left\{ (\alpha - \mu)^T \Sigma^{-1} (\alpha - \mu) \right\} \right\}$$

Focuses  
68.27%, 95.45%, 99.7%

$\alpha_{D \times 1}$ ,  $\mu_{D \times 1}$ ,  $\Sigma_{D \times D}$  Matrix.

In one dimensional space

$$G_1 = \frac{1}{\sigma \sqrt{2\pi}} \exp \left\{ -\frac{1}{2} \left\{ \frac{(\alpha - \mu)^2}{\sigma^2} \right\} \right\}$$

• pdf denoted by  $p_i$ ,  $P(x/w_i)$  for class  $i$   
 1. each assigned to

→ Baye's classification rule (Baye's classifier)

Let there be  $M$  classes ( $M \geq 2$ )

The prior probability be  $P_1, P_2, \dots, P_M$  where

$$P_i > 0 \text{ and } \sum_{i=1}^M P_i = 1 > 0$$

The prob. density functions.

$P_1(x), P_2(x), \dots, P_M(x)$  are given.

Assign new test pattern  $x_0$  to class  $i$

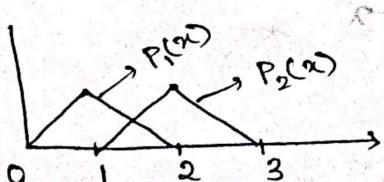
if  $P_i P_i(x_0) \geq P_j P_j(x_0) \quad \forall i, j \neq j$

Resolve ties arbitrarily

Example

$$M=2$$

$$P_1 = p \quad P_2 = 1 - p$$



$$P_1(x) = \begin{cases} x & 0 \leq x \leq 1 \\ 2-x & 1 \leq x \leq 2 \\ 0 & \text{otherwise} \end{cases}$$

$$P_2(x) = \begin{cases} \alpha-1 & x < 2 \\ 3-\alpha & 2 \leq x \leq 3 \\ 0 & \text{otherwise} \end{cases}$$

$$2p - p\alpha + \alpha - 1 = p\alpha + p - 1$$

$$P(2-x) = (1-p)(x-1). \quad P_{\alpha} = \frac{\alpha-1}{p+1}$$

Lecture 8  
27/10/93

Threshold value (of  $x$ ): a value such that if  $x < \text{th}$   
it belongs to class 1

and otherwise belongs to  
(if  $x = \text{th}$ ) class 2.

( $C \leq M$ ) consider  $M$  ad  $m$

Case I:  $0 \leq x \leq 1$  and  $P_1$  and  $P_2$  are

$$0 \leq x \leq 1$$

$$P_1 = p \quad P_2 = 1 - p$$

$$P_1(x) = \frac{x}{1-p}, \quad P_2(x) = \frac{1-x}{p}$$

$$P_1 P_1(x) > P_2 P_2(x) \quad \text{always here}$$

$$(x \in [0, 1] \setminus \{0.5\})$$

$\therefore x$  be in class 1.

Case II

$$2 \leq x \leq 3$$

$$1 - 1 = 0$$

$$9 - 9 = 0$$

$$P_1(x) = 0, \quad P_2(x) = 3 - x$$



$$P_2 P_2(x) > P_1 P_1(x)$$

$\therefore x$  is in class 2.

Case III

$$1 \leq x \leq 2.$$

$x$  be in class 1 }  $\{$   $\begin{cases} 0, q, 1 \\ 1, 2, 3 \end{cases}$   $\}$   $\{$   $\begin{cases} 1, 2, 3 \\ 4, 5, 6 \end{cases}$   $\}$

$$P_1 P_1(x) \geq P_2 P_2(x)$$

$$\Rightarrow (1-p)(2-x) \geq (1-p)(x-1)$$

$$\Rightarrow 2p - px \geq x - 1 - p + p$$

$$\Rightarrow p \geq x \Rightarrow x \leq 1+p$$

$\therefore x$  be in class 1 if  $x \leq 1+p$   
class 2 if  $x \geq 1+p$

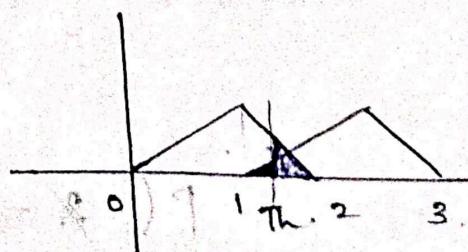
$$((e_{q+1}) - 1) + (q+1 - e) \cdot q$$

$$\omega = 0 \text{ to } 3.$$

↪ whole domain space

$$(0, q+1) \rightarrow (1 - (q+1)) (q+1) =$$

$$\omega_1 = 0 \text{ to } 1+p$$



$$\int_{\omega_1} P_1 P_1(x) + \int_{\omega_2} P_2 P_2(x) \rightarrow \underbrace{\int_{1-p}^q (q+x)p dx + \int_q^{1+p} (q+1)(1-p) dx}_{\text{error made by this classifier.}}$$

generally

$$\text{Error} = \sum_{i=1}^M \int_{\Omega_i^c} P_i P_i(x) dx$$

$$(x_1, q_1) \leq (x_2, q_2)$$

here  $E = (q - p) \int_{\Omega} P_1(x) dx + (1-p) \int_{\Omega} P_2(x) dx$

$$q - p - 1 + p \leq q - p$$

$$E = p \int_{\Omega} P_1(x) dx + (1-p) \int_{\Omega} P_2(x) dx.$$

$$= p \int_{\Omega} (q - x) dx + (1-p) \int_{\Omega} (x - p) dx$$

$$= p \left( 2(q - 1 - p) - \frac{1}{2}(4 - (1+p)^2) \right)$$

$$+ (1-p) \left[ \frac{1}{2}(1+p)^2 - 1 \right] + (1+p - 1)$$

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

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

$$E_B = \frac{P(1-P)}{2}$$

Classifier 2.

$$\Omega_1 = 0 \text{ to } 1.25$$

$$\Omega_2 = \{1.25 \text{ to } 3\}$$

2. coefficients of  $x^2$  &  $x^3$

$$E_{21} = P \int_{1.25}^2 (2-x) dx + (1-P) \int_2^3 (x-1) dx$$

\* Bayes classifier gives minimum misclassification probability

- Limitation of Bayes classification, pdf should be known.

④ Implementing Bayes classifier on datasets  
Uci - repository

- IRIS  
⇒ 150 patterns  
⇒ 4 attributes  
⇒ 3 classes.

after taking say 50% of data as training data

$$\text{calculate} \quad \left. \begin{array}{l} P_1 \Rightarrow \mu_1, \Sigma_1 \\ P_2 \Rightarrow \mu_2, \Sigma_2 \\ P_3 \Rightarrow \mu_3, \Sigma_3 \end{array} \right\} \text{assume Gaussian}$$

take  $x_0$  from test set

$$x_0 = \{ -, -, -, - \}$$

pass this to classifier

$x_0 \in P_1 \cup P_2 \cup P_3$   
if it give highest value of  $P_1, P_2, P_3(x_0)$

then 3 is the estimated class label.

To verify if it is correct compare it with actual class later.

- To measure goodness of classifier we have to follow some quantitative measure called evolution measure.

Lecture 4  
30/10/23

### Minimum Distance Classifier

Condition

1)  $P_1 = P_2 = \dots = P_M = 1/M$

All classes are equally probable.

2) Prob. density function follows Gaussian distribution

$$P_i = \frac{1}{(2\pi)^{D/2} |\Sigma_i|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mu_i)' \Sigma_i^{-1} (\mathbf{x} - \mu_i) \right\}$$

M of class means  $\mu_1, \mu_2, \dots, \mu_M$  with covariances  $\Sigma_1, \Sigma_2, \dots, \Sigma_M$  are precomputed.

3)  $\Sigma_1 = \Sigma_2 = \dots = \Sigma_M = I$  Identity Matrix.

$\mathbf{x}$  will be class  $i$

$$P_i P_i(\mathbf{x}) \geq P_j P_j(\mathbf{x}) \quad \forall i, j$$

$$\therefore P_1 = P_2 = \dots = P_M$$

$$P_i(\mathbf{x}) \geq P_j(\mathbf{x})$$

$$\frac{1}{(2\pi)^{D/2} |\Sigma_i|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mu_i)' I^{-1} (\mathbf{x} - \mu_i) \right\} \geq \frac{1}{(2\pi)^{D/2} |\Sigma_j|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mu_j)' I^{-1} (\mathbf{x} - \mu_j) \right\}$$

$\log_e$  in both sides in equation between  $\alpha$  &  $\mu_i$

$$\frac{1}{2} (\alpha - \mu_i)' I^{-1} (\alpha - \mu_i) \geq \frac{1}{2} (\alpha - \mu_j)' I^{-1} (\alpha - \mu_j)$$

$$(\alpha - \mu_i)' I^{-1} (\alpha - \mu_i) \leq (\alpha - \mu_j)' I^{-1} (\alpha - \mu_j)$$

$$(I^{-1} = I)$$

$$\underbrace{(\alpha - \mu_i)' (\alpha - \mu_i)} \leq (\alpha - \mu_j)' (\alpha - \mu_j)$$

$$d^2(\alpha, \mu_i)$$

$d(\alpha, \mu_i)$  = Euclidean distance  
shortest distance of  $\alpha$  from  $\mu_i$

assign  $\alpha$  in class  $i$  because

distance  $d(\alpha, \mu_i)$  is minimum.

MDC says that

1) calculate the  $\mu_i$  of each class  $i=1$  to  $M$   
based on training data

2) calculate the distance of test pattern  $\alpha$  to each  
of the class mean (i.e.  $d(\alpha, \mu_i) \forall i$ )

3) Assign  $\alpha$  in that class where distance  $d(\alpha, \mu_i)$   
 $i$  is minimum.

### Distance

(1) Euclidean distance

(2) weighted Euclidean distance

(3)  $L_p$  norm

(4) Mahalanobis distance

$$d_M(\alpha, \mu_i) = (\alpha - \mu_i)' \Sigma^{-1} (\alpha - \mu_i)$$

→ Mean of classes may be physical or virtual Pattern.  
 ↳ may not be a data point

→ Simple and less computational complexity method.

→ If above three conditions satisfy then MDC is equivalent to Bayes classifier as well as best one to apply.

### Confusion Matrix

it is a matrix which represents actual and predicted class label with their count.

		Predicted	
		+ve	-ve
Actual	+ve	True Positive = 600	False Negative = 50
	-ve	False Positive = 100	True Negative = 250

Overall Accuracy =  $\frac{\text{No. of Pattern correctly classified}}{\text{Total no. of Patterns}}$

as with  
all test  
patterns

$$= \frac{TP + TN}{TP + TN + FP + FN} \times 100 = \frac{850}{1000} \times 100 = 85\%$$

Class wise accuracy : for +ve class =  $\frac{TP}{TP + FN} = \frac{600}{650} \times 100\%$

for -ve class =  $\frac{TN}{TN + FP} = \frac{250}{350} \times 100\%$

For more than 2 classes

Actual	$C_1$	$C_2$	$\dots$	$C_M$
$C_1$	True Pos.	False Pos.	$\dots$	False Neg.
$C_2$	False Neg.	True Neg.	$\dots$	False Pos.
	Estimated to belong to $C_1$			

$$\text{overall Accuracy} = \frac{\text{sum of diagonal elements}}{\text{sum of all elements}} \times 100\%$$

$$\text{Class wise accuracy } C_i = \frac{\text{value of cell } i,i}{\text{sum of elements in row } i} \times 100\%$$

$$A = \frac{TP + TN}{TP + TN + FP + FN}$$

$$C_1 = \frac{TP}{TP + FN}$$

$$C_1 = \frac{TP}{TP + FN} = \frac{100}{100 + 90} = 0.555$$

## K - Nearest Neighbour.

1. Training Set:

$$S = \{x_1, x_2, \dots, x_n\} \in \mathbb{R}^D$$

and  $k$  +ve integer given to you.

Task  $\Rightarrow$  Assign the test pattern  $x_0$ .

2. Calculate:

$$d(x_0, x_i), \forall i=1 \text{ to } n$$

$d(\cdot) \Rightarrow$  Euclidean Distance.

3. Sort  $d(x_0, x_i)$  in non-decreasing order.

choose  $k$  patterns corresponding to these distances.

and treat them as  $k$  nearest neighbours.

4. let  $k_i$  = be the no. of nearest neighbours belongs to class  $i$

$$\sum_{i=1}^C k_i = K$$

$x_0$  be in class  $i$

if  $k_i \geq k_j$  &  $i, j \neq i$

Resolve ties arbitrarily.

## Observations

- 1.  $K$  is +ve integer

$$k = 1, 2, 3, 4, 5, \dots, 10, \dots, 2$$

Generally  $k$  is chosen as odd number.

$$k = 3, 5, 7, \dots, 11, \dots, 2$$

(to avoid ties)

- 2. if  $K=1$ , call it as Nearest Neighbour.

- 3. KNN follows weak learning. (there is no learning, only assigning is happening)

- 4. Different  $K$  values give different decisions.

- 5. If is computationally expensive

- 6. If a tie situation arise, resolve it arbitrarily. But if it happens frequently increase value of  $K$ .

## Bayes' Decision to KNN Decision

$x_0$  be in class  $i$

$$P_i P_i(x) > P_j P_j(x)$$

Let  $i \in \{1, 2, \dots, C\}$

$$\text{where } P_i = \text{prior prob.} = \frac{n_i}{n}$$

$$S = \{x_1, x_2, \dots, x_n\}, \quad \sum_{j=1}^C n_j = n$$

$p_i(x)$  = prob. density function of class  $i$

Mixen density function. 
$$p(x) = \sum_{i=1}^c p_i p_i(x)$$

Probability density =  $\frac{\text{Mass}}{\text{volume}}$

Let  $r$  be the distance between  $x_0$  and its  $k^{th}$  neighbour.

let  $A_n$  be the volume with radius  $r$  in  $D$ -dimensional domain

$$p_i(x) = \frac{k_i/n_i}{A_n} \rightarrow \begin{array}{l} k_i \Rightarrow \text{no. of patterns} \\ \text{in } i \text{ near to } x \\ n_i \Rightarrow \text{no. of patterns in } i \end{array}$$

Baip's Rule

$$p_i p_i(x) \geq p_j p_j(x)$$

$$\frac{n_i}{n} \cdot \frac{k_i}{A_n} \geq \frac{n_j}{n} \cdot \frac{k_j}{A_n}$$

$$k_i \geq k_j \quad x \text{ be in class } i$$

$\therefore KNN$  is a very good classifier.

Lecture 8  
8/11/23

## Linear Discriminant Analysis

$$M=2 \quad P_1 P_1(x) \geq P_2 P_2(x).$$

(1)  $p(x)$  follows Gaussian Distribution.

$$(2) \Sigma_1 = \Sigma_2 = \Sigma.$$

$$P_1 \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu_1)' \Sigma^{-1} (x - \mu_1) \right\} \geq P_2 \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}}$$

$$\exp \left\{ -\frac{1}{2} (x - \mu_2)' \Sigma^{-1} (x - \mu_2) \right\}$$

log both side

$$\log P_1 - \frac{1}{2} \{ (x - \mu_1)' \Sigma^{-1} (x - \mu_1) \} \geq \log P_2 - \frac{1}{2} \{ (x - \mu_2)' \Sigma^{-1} (x - \mu_2) \}$$

~~$$\log P_1 - \frac{1}{2} [x' \cancel{\Sigma} x - x' \Sigma \mu_1 - \mu_1' \Sigma x + \mu_1' \Sigma \mu_1]$$~~

$$\geq \log P_2 - \frac{1}{2} [\cancel{x' \Sigma x} - x' \Sigma \mu_2 - \mu_2' \Sigma x]$$

~~$$\text{neglect log base } e \text{ in max} \therefore + \mu_2' \Sigma \mu_2]$$~~

$$\begin{aligned} x' \mu &= \mu' x \\ x' I \mu &= \mu' I x \\ x' \Sigma \mu &= \mu' \Sigma x \end{aligned}$$

$$\Rightarrow \log \frac{P_1}{P_2} \geq \frac{1}{2} \left[ \alpha' \Sigma^{-1} (\mu_2 - \mu_1) + (\mu_1' \Sigma^{-1} \mu_1 - \mu_2' \Sigma^{-1} \mu_2) \right]$$

11/5

$$\Rightarrow \alpha' \underbrace{\Sigma^{-1} (\mu_2 - \mu_1)}_{D \times D \quad D \times 1} + \frac{1}{2} (\mu_1' \Sigma^{-1} \mu_1 - \mu_2' \Sigma^{-1} \mu_2) + \log \frac{P_2}{P_1} \leq 0$$

$$\Rightarrow \alpha' B + C \leq 0$$

• sogenannte Restriktionen auf  $\alpha$  fest:  $\alpha \in (-\infty)$

$$\alpha = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_D \end{bmatrix} \in \mathbb{R}^D \quad B = \begin{bmatrix} b_1 \\ \vdots \\ b_D \end{bmatrix} \in \mathbb{R}^{D \times 1}$$

$$b_1 \alpha_1 + b_2 \alpha_2 + \dots + b_D \alpha_D + \gamma \leq 0$$

$$\gamma = 0 + \text{rej. } \mathbb{R}$$

-> Bedingung für Realisierbarkeit der Parameter

Wahrsch. mehrerer Variable gleich Null

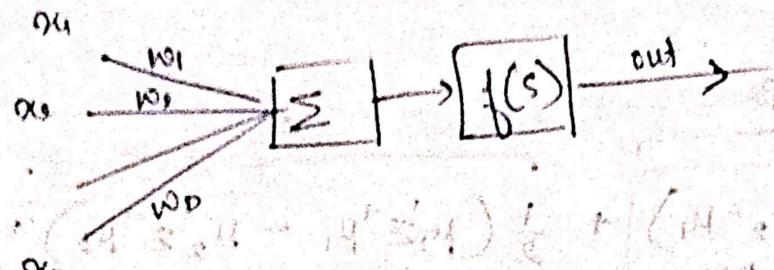
100% gewiss



Z -> T Gesamtwert

Lecture  
10/11/23

## Perception



$$s = \sum_{i=1}^D w_i x_i$$

$f(\cdot) \Rightarrow$  let it be threshold function

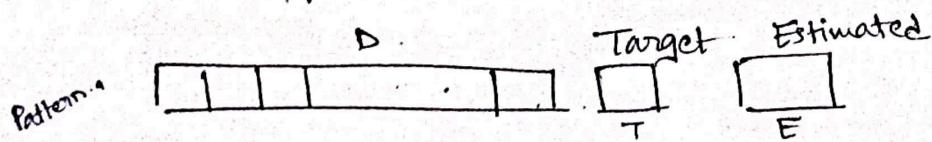
$$f\left(\sum w_i x_i\right) \begin{cases} > 0 & \text{output } 1 \\ \leq 0 & \text{output } 0. \end{cases}$$

$$\sum w_i x_i + b = 0.$$

Learning  $\Rightarrow$  Modification of weight

initially take random weight

Training set



$$\text{Error } \delta = T - E$$

if  $\delta \neq 0$ , we need to modify weight

$$w_i(t+1) = w_i(t) + \Delta w_i(t+1)$$

$$\Delta w_i(t+1) \propto \delta$$

$$\propto x_i$$

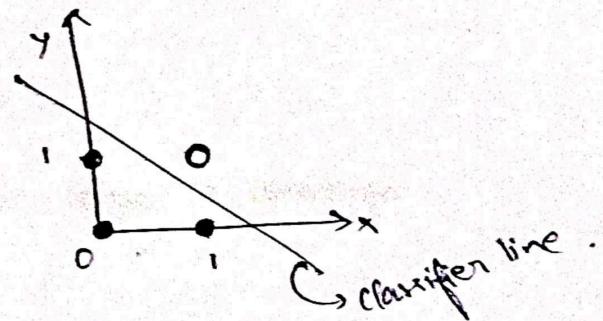
$$\Delta w_i(t+1) = \eta \delta x_i$$

$\eta$  learning rate constant

Example: Let  $x_1 = 1, x_2 = 2, \theta = 1.5, \eta = 0.5$

	Input	Output
$x_1 = 0, x_2 = 0$	0	0
$x_1 = 0, x_2 = 1$	1	0
$x_1 = 1, x_2 = 0$	0	0
$x_1 = 1, x_2 = 1$	1	1

$$\text{Eqn} : z = 1x_1 + 2x_2 - 1.5 = 1$$



Classifier line.

$$\text{let } \Theta = 1.5 \quad \eta = 0.5$$

Randomly initialize

$$w_1 = 0, w_2 = 1$$

$$z = 0 + 2 - 1.5 = 0.5$$

time/ t = 1

$$x = 0, y = 0$$

$$I_1 : s = 0 \times 0 + 1 \times 0 = 0 < 1.5 \therefore \text{class} = 0$$

(no change)

NO CHANGE

$$E = B_{out} T - E = 0$$

$$I_2: x=0 \quad y=1 \quad S = 0 \times 0 + 1 \times 1 = 1 \leq 1.5 \quad \text{Class 0}$$

Class 0.

$(0+1) \times 0 + (0+0) \times 1 = (0+1) \times 1$

$$\text{NO CHANGE, } \delta = T - E = 0 - 0 = 0$$

so

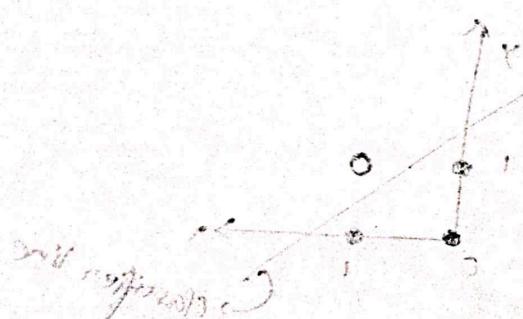
$$I_3: x=1 \quad y=0 \quad S = 0 \times 1 + 1 \times 0 = 0 \leq 1.5$$

Class 0

functions often produce 0

$$\text{NO CHANGE, } \delta = 0$$

$$I_4: \quad x=1 \quad y=1 \quad S = 0 \times 1 + 1 \times 1 = 1 \leq 1.5$$



function  $\delta = T - E$  Class 0

0	0	$\delta = T - E$
0	1	$= 1 - 0 = 1$
0	0	1
1	1	1
1	1	1

$$w_1 = 0 + 0.5 \times 1 \times 1 = 0.5$$

$$w_2 = 1 + 0.5 \times 1 \times 1 = 1.5$$

$$z_1 = f \quad z_2 = 0$$

$t=2$

$$w_1 = 0.5 \quad w_2 = 1.5$$

$$\& \eta = 0.5 \quad \theta = 1.5$$

$I_1$ : — no change —

$$I_2: \quad x=0 \quad y=1 \quad S = 1.5 \leq 1.5$$

so  $I_2 = 1.5$

Class = 0.

No change  $\delta = 0$ .

T<sub>3</sub>: No change.

T<sub>4</sub>:  $x=1 \quad y=1$

$$S = 0.5x + 1.5y = 2 > 1.5$$

Class 1

No change     $S = 1 - 1$   
 $= 0$

Stop. (As no change for all iterations)

∴ 
$$0.5x + 1.5y = 1.5$$
 → can be the classification line.

### • Learning Rate $\eta$

if  $\eta$  is small, time will be large to converge

if  $\eta$  is large, oscillates between original values and may not converge.

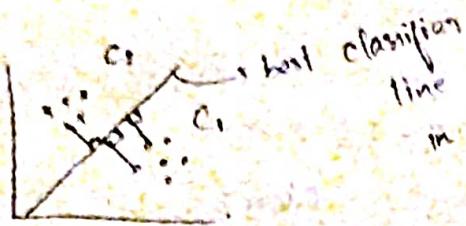
### • For multiclass problem ( $M > 2$ ) need multiple perceptrons.

↳ (not multi-layer. perception)

(H.W)  
• logical OR  
and logical NAND

Lecture  
15/11/20

Linear Classification:  $\omega^T x + b = 0$



line as  
in middle of  $c_1$  and  $c_2$

so that least prob. of  
miss classification

by distance of  
each pattern from classifier line.

can be  
perpendicular  
to line

Support vectors: those patterns which are on boundary  
of classes, where far distance  
from classifier line; helps deciding line

### Support vector Machines

Margin between two classes: distance of one boundary point with that



Maximize  
margin.

of other.

$[x_i, y_i]$  dataset  
 ↓  
 Pattern class label  
 $(-1, 1)$

$x_i : w^T x_i + b > 0$  then  $x_i$  be in  $+1$ , in class 1  
 $w^T x_i + b < 0$  then  $x_i$  be in  $-1$ , in class 2  
 $= 0$   $x_i$  is on the line.

Should be with appropriate  $\gamma_i$

$$\boxed{\gamma_i (w^T x_i + b) > 0}$$

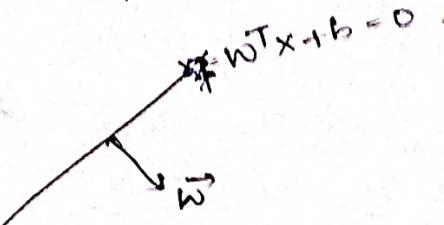
$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix}$$

∴ can represent them as vectors.

direction of  $\vec{w}$  will be perpendicular to the line

$\|w\| = \sqrt{w_1^2 + w_2^2 + \dots + w_D^2}$

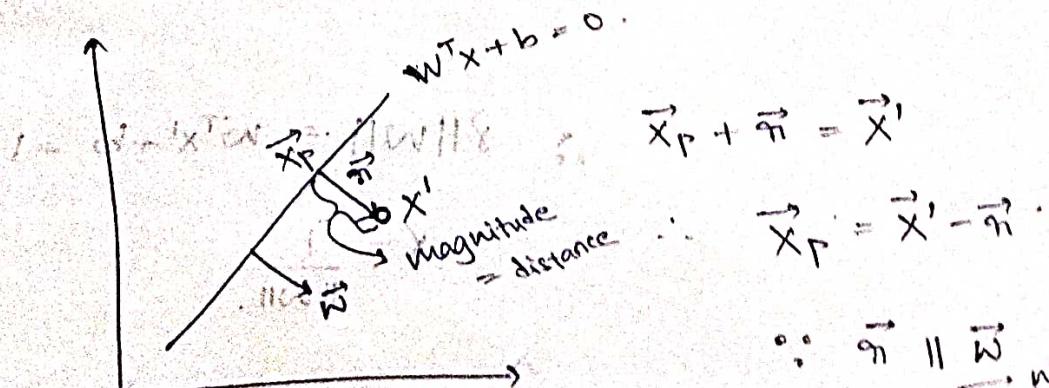
$\|w\| \vec{w} = w^T \vec{x}$



$\|w\| \vec{w}$   
 also  $\hat{w} = \|w\| \vec{w}$

After normalizing it we get

$$\hat{w} = \frac{\vec{w}}{\|w\|}$$



$\therefore \vec{r} \parallel \vec{w}$ , magnitude  
 $\therefore \vec{r} = \gamma \hat{w}$

$$\vec{w} \cdot \vec{w} = \frac{\vec{w} \cdot \vec{w}}{\|\vec{w}\|}$$

$\gamma = \text{Margin} = \text{distance of } x_p \text{ and } x'$

$$[0 < (d + x^T w) / \gamma]$$

$$w^T x_p + b = 0$$

$$\Rightarrow w^T [\vec{x}' - \vec{w}] + b = 0$$

$$\Rightarrow w^T \left[ x' - \frac{\gamma \vec{w}}{\|w\|} \right] + b = 0$$

$$\Rightarrow \left[ w^T x' - \gamma \frac{w^T \vec{w}}{\|w\|} \right] + b = 0$$

$$\because w^T w = \|w\|^2$$

$$\therefore w^T x' - \gamma \|w\| + b = 0$$

$$\Rightarrow w^T x' + b = \gamma \|w\|$$

time is independent with scaling operation.

$$w^T x' + b = 1.$$

$$\therefore \gamma \|w\| = w^T x' + b = 1$$

$$\gamma = \frac{1}{\|w\|}$$

similarly for  $x'$  belonging to  $-1$  i.e class 2.

$$\therefore \text{stat. } w^T x' + b = -1 \quad (1)$$

$$\therefore \gamma = \frac{-1}{\|w\|} \quad (\text{margin})$$

$$0 \leq 1 - (d + b)\gamma \quad (2)$$

$$\therefore \text{Total Margin} = 2\gamma = \frac{2}{\|w\|}$$

(1) P. condition: our task is to Maximize Total Margin  $= \frac{2}{\|w\|}$

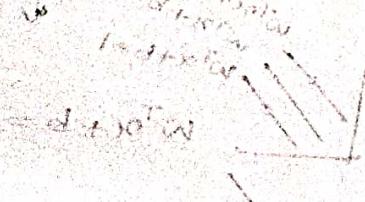
(2) Q.P. condition: Minimize  $\|w\|^2$

$$\therefore \text{Minimize } f(w, b) = \frac{\|w\|^2}{2}$$

with:

$$y_i (w^T x_i + b) > 0$$

[i.e. (1) & (2) constraint] in optimization problem



of  
vector  
space

is  $(w, b)$  if

max of

within set  $S$

and zero otherwise

Lecture  
17/11/23

## Minimization.

$$\phi(w) = \frac{1}{2} w \cdot w$$

constraint

$$y_i (w^T x_i + b) - 1 \geq 0$$

$\forall i=1 \text{ to } n$

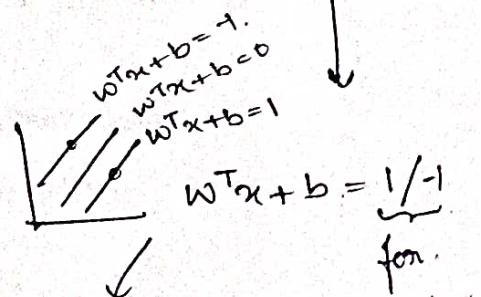
\* optimise  $f(x)$  constraint  $g_i(x)$ .

↳ unconstraint  $h(x) = f(x) - \sum \alpha_i g_i(x)$ .

$\alpha_i$  = Lagrange Multiplier

$$L(w, b) = \frac{1}{2} w \cdot w - \sum \alpha_i [y_i (w^T x_i + b) - 1]$$

$\alpha$  = Lagrange Multiplier.



$y_i (w^T x_i + b) > 1$

↳ for  $x_i$  to not be within the margin.

$$L(w, b) = \frac{1}{2} w \cdot w - \sum \alpha_i y_i w x_i - \sum \alpha_i y_i b + \sum \alpha_i$$

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

$$\frac{\partial L}{\partial b} = \sum \alpha_i y_i - 0 \quad \sum_{i=1}^n \alpha_i y_i = 0$$

$$\frac{\partial L}{\partial w} = w - \sum \alpha_i y_i x_i = 0$$

$$L = \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j (x_i \cdot x_j)$$

$$= \sum_i \sum_j \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) + \sum \alpha_i$$

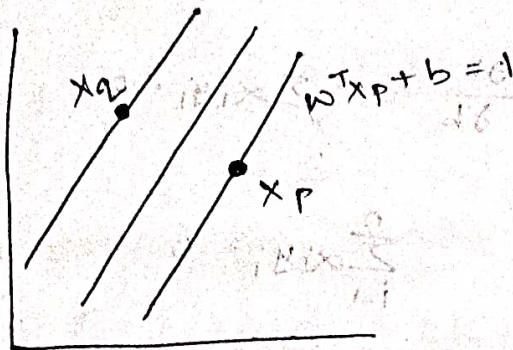
$$L = \sum \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j (x_i \cdot x_j)$$

$\alpha_i, y_i, x_i, x_j$  are known from training data  
 $\therefore L$  is optimization function of  $\alpha_i$

Solving it we get  $\alpha_i \ i=1 \text{ to } n$ .

$$\alpha_i \geq 0$$

$$w = \sum_{i=1}^n \alpha_i y_i x_i \quad (1)$$



$$\frac{1}{2} (w^T x_p + w^T x_q) = b \quad (2)$$

$$b = \frac{1}{2} \left[ \max_j \left( \sum_{i | y_i = -1} \alpha_i y_i (x_i \cdot x_j) \right) + \min_j \left( \sum_{i | y_i = +1} \alpha_i y_i (x_i \cdot x_j) \right) \right]$$

For any test pattern  $x_p$ .

$$w^T x_p + b > 0 \Rightarrow \text{class } +1$$

$$w^T x_p + b < 0 \Rightarrow \text{class } -1.$$

In SVM we say only support vectors are required  
but in eqn. of  $w$  all  $\alpha_i$  are contributing.

$\alpha_i = 0$ ; corresponding patterns are non-support vectors.

note: if  $\alpha_i > 0$ ; " " is support vector.

## Extraordinary cases

$\alpha_i$  very very large value.

$\alpha_i \gg 0 \rightarrow$  (noise pattern)

∴ discard them. i.e. don't use them in calc. of  $w$  and  $b$ .

- Multiclass  $M > 2$ .

SVM

SVM 1 → separate class 1 and Rest

SVM 2 → " class 2 and Rest

.

.

SVM M → " class M and Rest

- in real world. rarely linear classification occurs.

but non-linear can be transformed to linear.

in some higher dimension.

→ using kernel function

Kernel SVM very famous

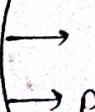
kernel function

↳ Transform the data set  
from lower dimensional

space to higher dimensional space.



$$f_3 = \sqrt{f_1^2 + f_2^2}$$



→ linear  
→ polynomial  
→ Radial basis function (RBF)  
Exponential ...

Lecture  
22/11/23

# Decision Tree (based Classification)

↳ follow tree structure.

↳ Each node contain attribute (Test of attribute)

↳ No. of branches depends on the no. of discrete values of attribute.

↳ The leaf node represents the decision

CGPA	Communication	Aptitude	Prog. skill	Job. opp.
H	G	H	G	
M	G	H	G	Yes
L	B	L	G	Yes
L	G	L	B	No
H	G	H	B	No
H	G	H	B	Yes
M	B	L	G	Yes
M	B	L	B	No
H	G	H	G	No
M	B	H	G	Yes
L	B	H	G	Yes
L	G	H	B	No
M	G	H	B	No
L	B	H	B	Yes
H	B	L	G	No
M	B	L	B	No
M	B	H	G	No
M	B	L	B	No
M	G	H	B	Yes

## Observations.

- works on discrete valued attribute set. (not continuous)
- works on numerical / non-numerical attribute set

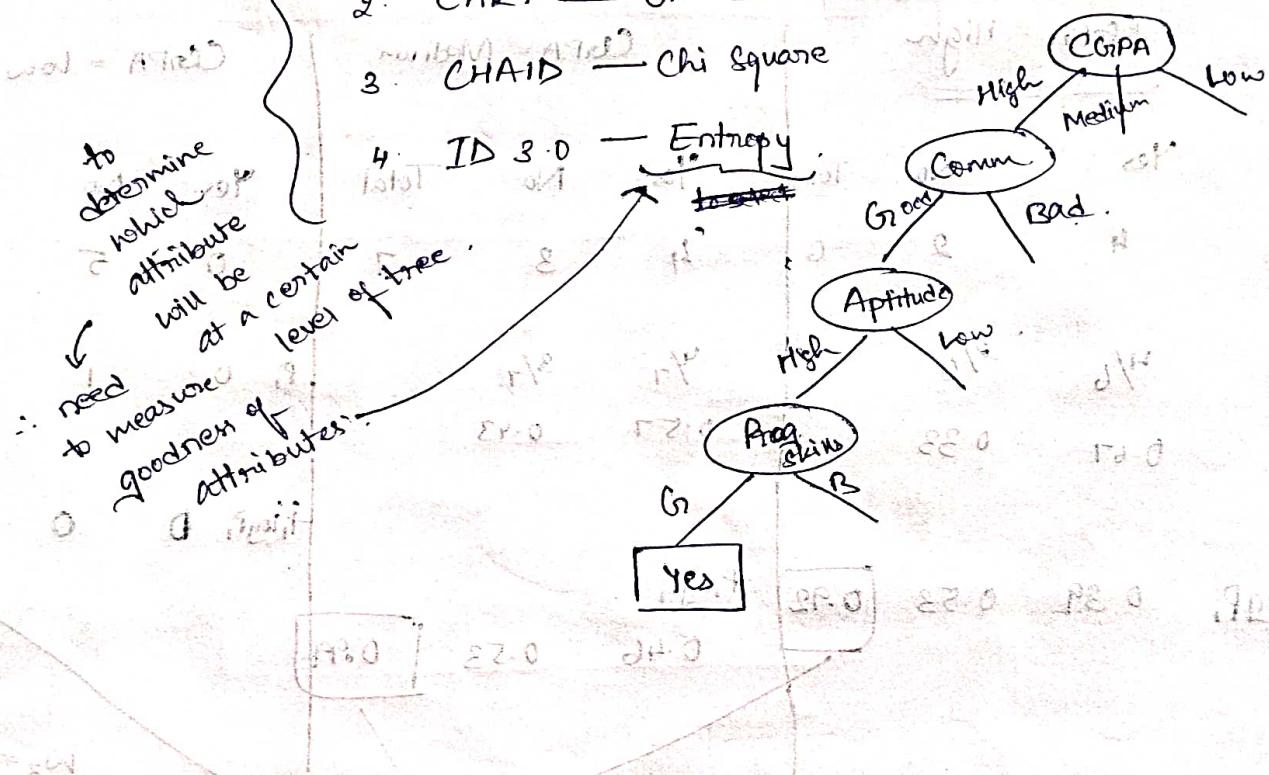
Examples of some classification techniques which follows decision tree based technique.

1. C 5.0 — Entropy

2. CART — Gini Index

3. CHAID — Chi square

4. ID 3.0 — Entropy



① Entropy based

↳ measures randomness present on the dataset

$$E_s = \sum_{i=1}^c -p_i \log p_i$$

$$p_i = \text{prior probability.} = \frac{n_i}{N}$$

$$E_{SA} = \sum w_i E_i$$

ie partition

w.r.t attribute A

$$\text{Information Gain. (A)} = E_s - E_{SA}$$

Eq.	S	Yes.	No.	Total.
count	8	10	18	

$$P_i = \frac{8}{18} \quad \frac{10}{18}$$

$$= 0.44 \quad 0.56$$

Yes      No      Total

$$-P_i \log P_i = -0.67 \log 0.44 = 0.47 \quad 0.47 + 0.52 = 0.99$$

$$\text{Indirect effect} = 0.52 \quad \text{Total effect} = 0.99$$

$$E_S = 0.99$$

$$A = \underline{\text{CGPA}}$$

CGPA = High.			CGPA = Medium			CGPA = Low.			
Yes	No	Total	Yes	No	Total	Yes	No	Total	
Count	4	2	6	4	3	7	0	5	5
P <sub>i</sub>	4/6	2/6		4/7	3/7		0	1	
-P <sub>i</sub> log P <sub>i</sub>	0.39	0.53	0.92	0.46	0.53	0.99	0	0	0

$$w_1 = 6/18$$

$$w_2 = 7/18$$

$$w_3 = 5/18$$

ie Prohibition ie High, Medium, Low

$$E_{SA} = w_1 E_1 + w_2 E_2 + w_3 E_3$$

$$(A) = \frac{6}{18} \times 0.92 + \frac{7}{18} \times 0.99 + \frac{5}{18} \times 0$$

$$\text{Total} = 0.69$$

$$IG_1(C_{GPA}) = 0.99 - 0.69 = 0.30.$$

Why  $IG_1(\text{Comm}) = ?$

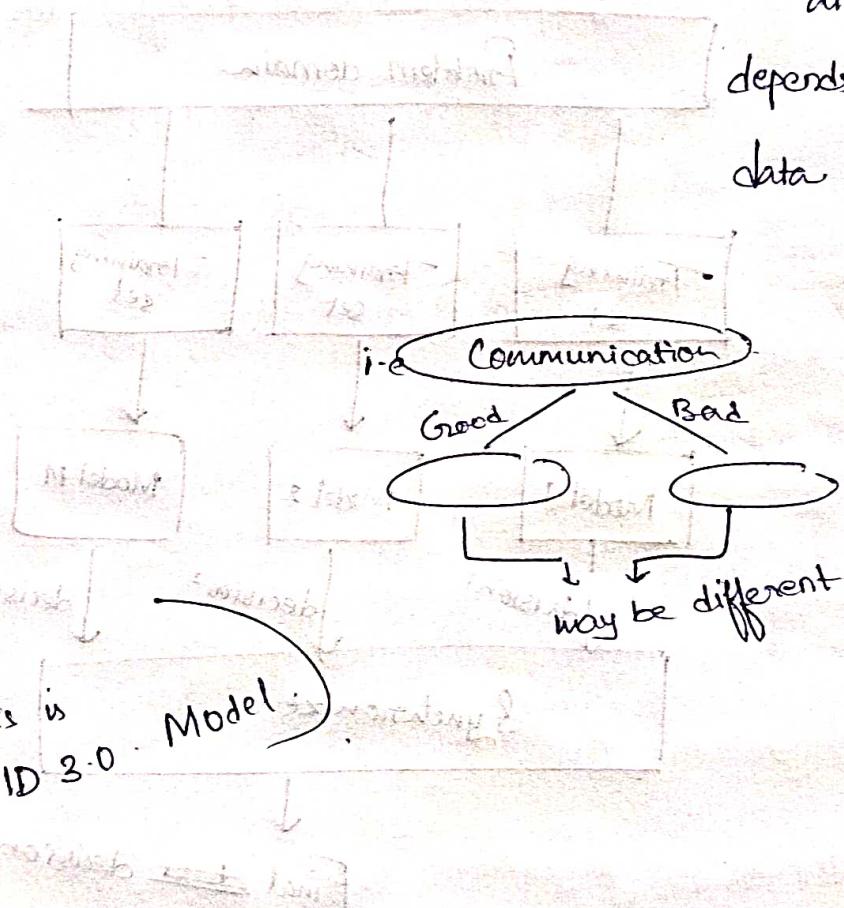
$$IG_1(\text{Aptitude}) = ?$$

$$IG_1(\text{Prog. skill}) = ?$$

select the attribute with Least  $IG_1$  as ~~root~~<sup>root</sup> attribute.

Now, for next layer/level of tree, again calculate for next attribute as it

depends on each partition data of previous level attribute.



Lecture  
24/11/23

## Ensemble of classifier.

$$0.2 \cdot 0 = 0.2 \cdot 0 - 0.0 = (ABCD) \cdot 0$$

- It allows us to incorporate more than one similar and different model and then synchronise their decision to take the final decision.

Problem domain.

Training set      Training set      Training set

Model 1

Model 2

Model M

Synchronize

decision<sup>1</sup>

decision<sup>2</sup>

decision<sup>M</sup>

Final ~~decide~~ decision

eg. Average decision /  
Majority Voting

similar → means working principle is same

different → may give diff result for. same pattern.

## Random Forest Classification

↳ based on Ensemble of classifier.

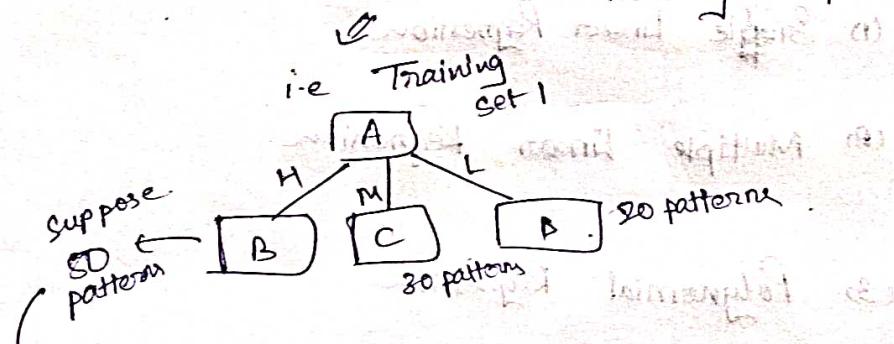
↳ It allows us to incorporate more than one decision tree model and then synchronise their decision to take the final decision.

(~~we get different decision tree~~, not because of diff. quantifiers used)

but because, & different training patterns are chosen randomly for each model hence each set may.

have diff. decision tree (subtrees)

↳ which are formed base on training set patterns.



↳ based on these 50 patterns next attribute is selected → subtree depends on diff trees for diff training set

## Regression

Supervised

Intelligent or intelligent

Estimating / Approximating.

a dependent variable (predicted ~~value~~<sup>variable</sup>) from one or more independent variables (predictor variable)

with the help of prior information.

$$y = f(x_1, x_2, x_3, \dots, x_n)$$

(1) Simple Linear Regression,

(2) Multiple Linear Regression.

(3) Polynomial Regression

(4) Non-linear Regression

(5) Logistic Regression.

for classification