

*. Normal density & Discriminant function

$$g_i^*(x) = \ln p(x|w_i) + \ln p(w_i)$$

$$p(x|w_i) = \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp \left[-\frac{1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) \right]$$

$$\Rightarrow g_i^*(x) = -\frac{1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) - \frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln |\Sigma_i| + \ln p(w_i)$$

Independent of
classes.

case II : $\Sigma_i = \Sigma$

(*) discriminant function :-

$$P(w_i|x) > P(w_j|x)$$

; w_1 and w_2 are classes.

(-) Bayes minimum error classifier:-

In this we compute $P(w_i|x)$ for all the classes. and choose the class which have highest value of $P(w_i|x)$.

(-) Bayes minimum risk classifier :-

In this we compute $R(\alpha_i|x)$ risk for all classes.

$$R(\alpha_i|x) = \sum_{\forall w_j} \lambda(\alpha_i|w_j) P(w_j|x)$$

$$\begin{aligned} \lambda(\alpha_i|w_j) &= 0 & i &= j \\ &= 1 & i &\neq j \end{aligned}$$

choose the class which have highest value of $-R(\alpha_i|x)$.

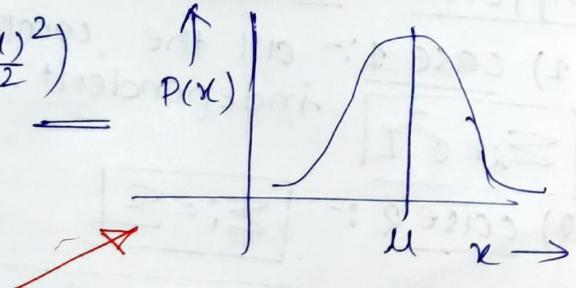
$$\Rightarrow g_i(x) = P(w_i|x)$$

$$g_i(x) = P(x|w_i) \cdot P(w_i)$$

$$g_i(x) = \ln P(x|w_i) + \ln P(w_i)$$

(*) Discriminant function under multivariate normal distribution :-

$$P(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{\frac{-(x-\mu)^2}{2\sigma^2}}$$



Normal Distribution.

Let us assume that our feature vector is in d dimension.

$x \rightarrow d$ dimension

then probability distribution,

$$P(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right]$$

$$P(x|w_i) = \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp \left[-\frac{1}{2} (x-\mu_i)^T \Sigma_i^{-1} (x-\mu_i) \right]$$

we know that,

$$g_i(x) = \ln P(x|w_i) + \ln(P(w_i))$$

$$\therefore g_i(x) = \ln \left(\frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp \left[-\frac{1}{2} (x-\mu_i)^T \Sigma_i^{-1} (x-\mu_i) \right] \right) + \ln P(w_i)$$

$$= -\frac{1}{2} (x-\mu_i)^T \Sigma_i^{-1} (x-\mu_i) - \frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln |\Sigma_i| + \ln P(w_i)$$

independent of class
so, it can be ignored.

$$\therefore g_i(x) = -\frac{1}{2} (x-\mu_i)^T \Sigma_i^{-1} (x-\mu_i) - \frac{1}{2} \ln |\Sigma_i| + \ln P(w_i)$$

Different cases :-

1) case 1 :- all the components are statistically independent. which means variance = 0.

$$\Sigma_i = \sigma^2 I$$

2) case 2 :- $\Sigma_i = \Sigma$

3) case 3 :- $\Sigma_i = \text{arbitrary}$

For case I:-

$$g_i(x) = \frac{1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) - \frac{1}{2} \ln |\Sigma_i| + \ln P(w_i)$$
$$\Sigma_i = \sigma^2 I.$$

$$\therefore g_i(x) = \frac{1}{2\sigma^2} (x - \mu_i)^T (x - \mu_i) + \ln P(w_i)$$

$$\therefore g_i(x) = \frac{1}{2\sigma^2} [x^T x - 2\mu_i^T x + \mu_i^T \mu_i] + \ln P(w_i)$$

$$\therefore g_i(x) = \frac{\mu_i^T x}{\sigma^2} - \frac{\mu_i^T \mu_i}{2\sigma^2} + \ln P(w_i)$$

class independent

$$g_i(x) \approx w_i^T x + w_{i0}$$

Linear

$$w_i = \frac{\mu_i^T}{\sigma^2} ; w_{i0} = \frac{1}{2\sigma^2} \mu_i^T \mu_i + \ln P(w_i)$$

For boundary,

$$g(x) = g_i(x) - g_j(x) = 0$$

$$g_i(x) = \frac{1}{2\sigma^2} (x - \mu_i)^T (x - \mu_i) + \ln P(w_i)$$

$$g_j(x) = \frac{1}{2\sigma^2} (x - \mu_j)^T (x - \mu_j) + \ln P(w_j)$$

$$g(x) = \frac{\mu_i^T x}{\sigma^2} - \frac{\mu_i^T \mu_i}{2\sigma^2} + \ln P(w_i) - \frac{\mu_j^T x}{\sigma^2} + \frac{\mu_j^T \mu_j}{2\sigma^2} - \ln P(w_j)$$

$$g(x) = \frac{(\mu_i - \mu_j)^T x}{\sigma^2} - \frac{\mu_i^T \mu_i}{2\sigma^2} + \frac{\mu_j^T \mu_j}{2\sigma^2} + \ln \frac{P(w_i)}{P(w_j)} = 0$$

$$g(x) = (\mu_i - \mu_j)^T x - \frac{1}{2} \left(\frac{\mu_i^T \mu_i - \mu_j^T \mu_j}{\sigma^2} + \frac{\ln \frac{P(w_i)}{P(w_j)}}{\sigma^2} \right) = 0$$

$$g(x) \approx w^T (x - x_0) = 0$$

$$w = (\mu_i - \mu_j) \quad \text{and} \quad x_0 = \frac{1}{2} (\mu_i + \mu_j) - \frac{\sigma^2}{\|\mu_i - \mu_j\|^2} \left(\ln \left(\frac{P(w_i)}{P(w_j)} \right) (\mu_i - \mu_j) \right)$$

Physical significance of $\Sigma_i = \sigma^2 I$.

- 1) Different component of feature vector are statistically independent.
and every component has same variance.

- 2) spread \rightarrow circle in 2D
 \rightarrow sphere in 3D
 \rightarrow Hypersphere in nD.

(-) Example :-

$$\omega_1 \Rightarrow \begin{pmatrix} 12 \\ 4 \end{pmatrix} \begin{pmatrix} 12 \\ 8 \end{pmatrix} \begin{pmatrix} 40 \\ 6 \end{pmatrix} \begin{pmatrix} 14 \\ 6 \end{pmatrix} \quad \omega_2 \Rightarrow \begin{pmatrix} 9 \\ 10 \end{pmatrix} \begin{pmatrix} 9 \\ 14 \end{pmatrix} \begin{pmatrix} 7 \\ 12 \end{pmatrix} \begin{pmatrix} 11 \\ 12 \end{pmatrix}$$

$$u_1 = \left[\left(\frac{12}{4} \right) + \left(\frac{12}{8} \right) + \left(\frac{40}{6} \right) + \left(\frac{14}{6} \right) \right] \div \left(\frac{1}{4} \right)$$

$$u_1 = \frac{\begin{pmatrix} 48 \\ 24 \end{pmatrix}}{\left(\frac{1}{4} \right)} = \boxed{\begin{pmatrix} 12 \\ 6 \end{pmatrix}} \quad u_2 = \boxed{\begin{pmatrix} 9 \\ 12 \end{pmatrix}}$$

$$(x_1 - u_1)(x_1 - u_1)^T = \begin{pmatrix} 0 \\ -2 \end{pmatrix} \begin{pmatrix} 0 & -2 \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ 0 & 4 \end{pmatrix} = M_1$$

$$(x_2 - u_1)(x_2 - u_1)^T = \begin{pmatrix} 0 \\ 2 \end{pmatrix} \begin{pmatrix} 0 & 2 \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ 0 & 4 \end{pmatrix} = M_2$$

$$(x_3 - u_1)(x_3 - u_1)^T = \begin{pmatrix} -2 \\ 0 \end{pmatrix} \begin{pmatrix} -2 & 0 \end{pmatrix} = \begin{pmatrix} 4 & 0 \\ 0 & 0 \end{pmatrix} = M_3$$

$$(x_4 - u_1)(x_4 - u_1)^T = \begin{pmatrix} 2 \\ 0 \end{pmatrix} \begin{pmatrix} 2 & 0 \end{pmatrix} = \begin{pmatrix} 4 & 0 \\ 0 & 0 \end{pmatrix} = M_4$$

$$\Sigma_1 = \frac{1}{4} [M_1 + M_2 + M_3 + M_4]$$

$$\Sigma_1 = \frac{1}{4} \begin{bmatrix} 8 & 0 \\ 0 & 8 \end{bmatrix}$$

$$\boxed{\Sigma_1 = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}}$$

$$(x_1 - u_2)(x_1 - u_2)^T = \begin{pmatrix} 0 \\ -2 \end{pmatrix}(0 - 2) = \begin{pmatrix} 0 & 0 \\ 0 & 4 \end{pmatrix} = m_1'$$

$$(x_2 - u_2)(x_2 - u_2)^T = \begin{pmatrix} 0 \\ 2 \end{pmatrix}(0 - 2) = \begin{pmatrix} 0 & 0 \\ 0 & 4 \end{pmatrix} = m_2'$$

$$(x_3 - u_2)(x_3 - u_2)^T = \begin{pmatrix} -2 \\ 0 \end{pmatrix}(2 - 0) = \begin{pmatrix} 4 & 0 \\ 0 & 0 \end{pmatrix} = m_3'$$

$$(x_4 - u_2)(x_4 - u_2)^T = \begin{pmatrix} 2 \\ 0 \end{pmatrix}(2 - 0) = \begin{pmatrix} 4 & 0 \\ 0 & 0 \end{pmatrix} = m_4'$$

$$\Sigma_2 = \frac{1}{4} [m_1' + m_2' + m_3' + m_4']$$

$$\Sigma_2 = \frac{1}{4} \begin{pmatrix} 8 & 0 \\ 0 & 8 \end{pmatrix} = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix} = \boxed{2I} \checkmark$$

$$\boxed{\Sigma_1 = \Sigma_2 = 2I} \approx \underline{\underline{\sigma^2 I}} \quad \therefore \boxed{\sigma = \sqrt{2}}$$

Here, $\boxed{P(w_1) = P(w_2)}$ $[\because \# w_1 = \# w_2]$

$$\therefore x_0 = \frac{1}{2}(u_1 + u_2) \quad \text{and} \quad w = u_1 - u_2.$$

$$\therefore x_0 = \frac{1}{2} \left(\begin{pmatrix} 1 & 2 \\ 6 & 1 \end{pmatrix} + \begin{pmatrix} 9 & 1 \\ 1 & 2 \end{pmatrix} \right) = \frac{1}{2} \begin{pmatrix} 21 \\ 18 \end{pmatrix} \quad w = \begin{pmatrix} 3 \\ -6 \end{pmatrix}$$

$$\therefore w^T(x - x_0) = 0$$

$$(3 - 6) \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} - \frac{(3 - 6)}{2} \begin{pmatrix} 21 \\ 18 \end{pmatrix} = 0$$

$$3x_1 - 6x_2 - \frac{1}{2} (63 - 108) = 0$$

$$3x_1 - 6x_2 + \frac{45}{2} = 0$$

$$6x_1 - 12x_2 + 45 = 0$$

$$\boxed{3x_1 - 4x_2 + 15 = 0}$$

Decision Boundary

For case II i. $\Sigma_i = \Sigma \neq \sigma^2 I$.

→ The component of feature vector may be statistically independent but not equal or same.

$$g_i(x) = -\frac{1}{2}(x - \mu_i)^T \Sigma^{-1} (x - \mu_i) - \frac{1}{2} \ln |\Sigma| + \ln P(w_i)$$

class independent

$$\therefore g_i(x) = \frac{1}{2}(x^T \Sigma^{-1} x - 2\mu_i^T \Sigma^{-1} x + \mu_i^T \Sigma^{-1} \mu_i) + \ln P(w_i)$$

$$\therefore g_i(x) = \mu_i^T \Sigma^{-1} x - \frac{\mu_i^T \Sigma^{-1} \mu_i}{2} + \ln P(w_i)$$

$$\boxed{g_i(x) \approx w_i^T x + w_{i0}}$$

$$w_i = \mu_i^T \Sigma^{-1} \quad w_{i0} = -\frac{1}{2} \mu_i^T \Sigma^{-1} \mu_i + \ln P(w_i)$$

For decision boundary,

$$g(x) = g_1(x) - g_2(x)$$

$$g_1(x) = \mu_1^T \Sigma^{-1} x - \frac{\mu_1^T \Sigma^{-1} \mu_1}{2} + \ln P(w_1)$$

$$g_2(x) = \mu_2^T \Sigma^{-1} x - \frac{\mu_2^T \Sigma^{-1} \mu_2}{2} + \ln P(w_2)$$

$$g(x) = \mu_1^T \Sigma^{-1} x - \mu_2^T \Sigma^{-1} x - \frac{\mu_1^T \Sigma^{-1} \mu_1}{2} + \frac{\mu_2^T \Sigma^{-1} \mu_2}{2} + \ln \frac{P(w_1)}{P(w_2)}$$

$$g(x) = (\mu_1 - \mu_2)^T \Sigma^{-1} x - \frac{1}{2} (\mu_1^T \Sigma^{-1} \mu_1 - \mu_2^T \Sigma^{-1} \mu_2) + \ln \frac{P(w_1)}{P(w_2)}$$

$$\boxed{g(x) \approx w^T (x - x_0) = 0} \rightarrow \text{Linear}$$

$$w = \Sigma^{-1}(\mu_1 - \mu_2)$$

$$x_0 = \frac{1}{2}(\mu_1 + \mu_2) - \frac{1}{(\mu_1 - \mu_2)^T \Sigma^{-1} (\mu_1 - \mu_2)} \left(\ln \frac{P(w_1)}{P(w_2)} \right) (\mu_1 - \mu_2)$$

(-) Example :-

$$w_1 \Rightarrow \begin{pmatrix} 6 \\ 2 \end{pmatrix} \begin{pmatrix} 9 \\ 3 \end{pmatrix} \begin{pmatrix} 7 \\ 5 \end{pmatrix} \begin{pmatrix} 10 \\ 6 \end{pmatrix}$$

$$w_2 \Rightarrow \begin{pmatrix} 6 \\ 11 \end{pmatrix} \begin{pmatrix} 9 \\ 12 \end{pmatrix} \begin{pmatrix} 7 \\ 14 \end{pmatrix} \begin{pmatrix} 10 \\ 15 \end{pmatrix}$$

$$u_1 \Rightarrow \begin{pmatrix} 32 \\ 16 \end{pmatrix} / 4 = \begin{pmatrix} 8 \\ 4 \end{pmatrix}$$

$$\boxed{u_1 = \begin{pmatrix} 8 \\ 4 \end{pmatrix}}$$

$$u_2 = \begin{pmatrix} 32 \\ 52 \end{pmatrix} / 4 = \begin{pmatrix} 8 \\ 13 \end{pmatrix}$$

$$\boxed{u_2 = \begin{pmatrix} 8 \\ 13 \end{pmatrix}}$$

$$(x_1 - u_1)(x_1 - u_1)^T = \begin{pmatrix} -2 \\ -2 \end{pmatrix} (-2 - 2) = \begin{pmatrix} 4 & 4 \\ 4 & 4 \end{pmatrix} = m_1$$

$$(x_2 - u_1)(x_2 - u_1)^T = \begin{pmatrix} 1 \\ -1 \end{pmatrix} (1 - 1) = \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix} = m_2$$

$$(x_3 - u_1)(x_3 - u_1)^T = \begin{pmatrix} 1 \\ 1 \end{pmatrix} (1 - 1) = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} = m_3$$

$$(x_4 - u_1)(x_4 - u_1)^T = \begin{pmatrix} 2 \\ 2 \end{pmatrix} (2 - 2) = \begin{pmatrix} 4 & 4 \\ 4 & 4 \end{pmatrix} = m_4$$

$$\Sigma_1 = \frac{1}{4} [m_1 + m_2 + m_3 + m_4] = \frac{1}{4} \begin{bmatrix} 10 & 6 \\ 6 & 10 \end{bmatrix} = \begin{bmatrix} 5/2 & 3/2 \\ 3/2 & 5/2 \end{bmatrix}$$

$$\Sigma_1^{-1} = \frac{1}{8} \begin{bmatrix} 5 & -3 \\ -3 & 5 \end{bmatrix}$$

$$(x_1 - u_2)(x_1 - u_2)^T = \begin{pmatrix} -2 \\ -2 \end{pmatrix} (-2 - 2) = \begin{pmatrix} 4 & 4 \\ 4 & 4 \end{pmatrix} = m_1'$$

$$(x_2 - u_2)(x_2 - u_2)^T = \begin{pmatrix} 1 \\ -1 \end{pmatrix} (1 - 1) = \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix} = m_2'$$

$$(x_3 - u_2)(x_3 - u_2)^T = \begin{pmatrix} 1 \\ 1 \end{pmatrix} (1 - 1) = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} = m_3'$$

$$(x_4 - u_2)(x_4 - u_2)^T = \begin{pmatrix} 2 \\ 2 \end{pmatrix} (2 - 2) = \begin{pmatrix} 4 & 4 \\ 4 & 4 \end{pmatrix} = m_4'$$

$$\Sigma_2 = \begin{bmatrix} 5/2 & 3/2 \\ 3/2 & 5/2 \end{bmatrix}$$

$$\Sigma_2^{-1} = \begin{bmatrix} 5/8 & -3/8 \\ -3/8 & 5/8 \end{bmatrix}$$

$$\therefore \boxed{\Sigma_1 = \Sigma_2 = \Sigma \neq \sigma^2 I}$$

$$\therefore x_0 = \left(\frac{u_1 + u_2}{2} \right) \quad \left[\begin{array}{l} \text{since } P(u_1) = P(u_2) \\ \therefore \ln \frac{P(u_1)}{P(u_2)} = 0 \end{array} \right]$$

$$\therefore w = \Sigma^T (u_1 - u_2)$$

$$w = \begin{bmatrix} 5/8 & -3/8 \\ -3/8 & 5/8 \end{bmatrix} \begin{bmatrix} 0 \\ 9 \end{bmatrix} = \begin{bmatrix} -27/8 \\ 45/8 \end{bmatrix}$$

$$\therefore g(x) = w^T x - w^T x_0$$

$$g(x) = \begin{bmatrix} -\frac{27}{8} & \frac{45}{8} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - \begin{bmatrix} -27/8 & 45/8 \end{bmatrix} \begin{bmatrix} 16/2 \\ 17/2 \end{bmatrix}$$

$$g(x) = -\frac{27x_1}{8} + \frac{45x_2}{8} - \left[-\frac{27 \times 16}{16} + \frac{45 \times 17}{8 \times 2} \right]$$

$$g(x) = -\frac{27x_1 + 45x_2}{8} - \frac{333}{16} = 0$$

$$\boxed{g(x) = 54x_1 - 90x_2 + 333 = 0} \quad \# \checkmark$$

Decision Boundary

For case III :- $\Sigma^i = \underline{\text{arbitrary}}$

$$g_i(x) = -\frac{1}{2}(x - \mu_i)^T \Sigma^i (x - \mu_i) - \frac{1}{2} \ln |\Sigma^i| + \ln P(w_i)$$

we can't ignore

Reason:- Different for different classes.

$$\therefore g_i(x) = -\frac{1}{2} (x^T \Sigma^i x - 2 \mu_i^T \Sigma^i x + \mu_i^T \Sigma^i \mu_i) - \frac{1}{2} \ln |\Sigma^i| + \ln P(w_i)$$

$$\therefore g_i(x) = -\frac{1}{2} x^T \Sigma^i x + \mu_i^T \Sigma^i x - \frac{1}{2} \mu_i^T \Sigma^i \mu_i - \frac{1}{2} \ln |\Sigma^i| + \ln P(w_i)$$

$$\boxed{g_i(x) \approx x^T A_i x + B_i^T x + C_i} \rightarrow \text{Quadratic}$$

$$A_i = -\frac{1}{2} \Sigma^i \quad B_i = \mu_i^T \Sigma^i$$

$$C_i = -\frac{1}{2} \mu_i^T \Sigma^i \mu_i - \frac{1}{2} \ln |\Sigma^i| + \ln P(w_i)$$

For Decision boundary,

$$g(x) = g_1(x) - g_2(x) = 0$$

$$g_1(x) = -\frac{1}{2} x^T \Sigma_1 x + \mu_1^T \Sigma_1 x - \frac{1}{2} \mu_1^T \Sigma_1 \mu_1 - \frac{1}{2} \ln |\Sigma_1| + \ln P(w_1)$$

$$g_2(x) = -\frac{1}{2} x^T \Sigma_2 x + \mu_2^T \Sigma_2 x - \frac{1}{2} \mu_2^T \Sigma_2 \mu_2 - \frac{1}{2} \ln |\Sigma_2| + \ln P(w_2)$$

$$\therefore g(x) = -\frac{1}{2} x^T (\Sigma_1 - \Sigma_2) x + (\mu_1^T \Sigma_1 - \mu_2^T \Sigma_2) x - \frac{1}{2} \ln \frac{|\Sigma_1|}{|\Sigma_2|} \\ - \frac{1}{2} \mu_1^T \Sigma_1 \mu_1 + \frac{\mu_2^T \Sigma_2 \mu_2}{2} + \ln \frac{P(w_1)}{P(w_2)}$$

$$\boxed{g(x) \approx x^T A_i x + B_i^T x + C_i}$$

$$A_i = -\frac{1}{2} (\Sigma_1 - \Sigma_2)$$

$$B_i = (\mu_1^T \Sigma_1 - \mu_2^T \Sigma_2)$$

$$C_i = -\frac{1}{2} \mu_1^T \Sigma_1 \mu_1 + \frac{1}{2} \mu_2^T \Sigma_2 \mu_2$$

$$-\frac{1}{2} \ln \frac{|\Sigma_1|}{|\Sigma_2|} + \ln \frac{P(w_1)}{P(w_2)}$$

(-) Example

$$w_1 \Rightarrow \begin{pmatrix} 5 \\ 4 \end{pmatrix} \begin{pmatrix} 9 \\ 2 \end{pmatrix} \begin{pmatrix} 9 \\ 6 \end{pmatrix} \begin{pmatrix} 13 \\ 6 \end{pmatrix}$$

$$u_1 = \begin{pmatrix} 36 \\ 18 \end{pmatrix} / 4$$

$$\boxed{u_1 = \begin{pmatrix} 9 \\ 9/2 \end{pmatrix}}$$

$$(x_1 - u_1)(x_1 - u_1)^T$$

$$= \begin{pmatrix} 4 & 1/2 \\ -1/2 \end{pmatrix} \begin{pmatrix} 4 & 1/2 \\ -1/2 \end{pmatrix}$$

$$= \begin{pmatrix} 16 & -2 \\ -2 & 4/4 \end{pmatrix} = m_1$$

$$(x_2 - u_1)(x_2 - u_1)^T$$

$$= \begin{pmatrix} 0 \\ -5/2 \end{pmatrix} \begin{pmatrix} 0 & -5/2 \end{pmatrix}$$

$$= \begin{pmatrix} 0 & 0 \\ 0 & 25/4 \end{pmatrix} = m_2$$

$$(x_3 - u_1)(x_3 - u_1)^T$$

$$= \begin{pmatrix} 0 \\ 3/2 \end{pmatrix} \begin{pmatrix} 0 & 3/2 \end{pmatrix}$$

$$= \begin{pmatrix} 0 & 0 \\ 0 & 9/4 \end{pmatrix} = m_3$$

$$(x_4 - u_1)(x_4 - u_1)^T$$

$$= \begin{pmatrix} 4 \\ 3/2 \end{pmatrix} \begin{pmatrix} 4 & 3/2 \end{pmatrix}$$

$$= \begin{pmatrix} 16 & 6 \\ 6 & 9/4 \end{pmatrix} = m_4$$

$$\Sigma_1 = \frac{1}{4} [m_1 + m_2 + m_3 + m_4]$$

$$= \frac{1}{4} \begin{bmatrix} 82 & 4 \\ 4 & 11 \end{bmatrix} = \begin{bmatrix} 8 & 1 \\ 1 & 11/4 \end{bmatrix}$$

$$\Sigma_1^T = \frac{1}{21} \begin{bmatrix} 11/4 & 1 \\ -1 & 8 \end{bmatrix} //$$

$$w_2 \Rightarrow \begin{pmatrix} 9 \\ 10 \end{pmatrix} \begin{pmatrix} 9 \\ 14 \end{pmatrix} \begin{pmatrix} 7 \\ 12 \end{pmatrix} \begin{pmatrix} 11 \\ 12 \end{pmatrix}$$

$$u_2 = \begin{pmatrix} 36 \\ 48 \end{pmatrix} / 4$$

$$\boxed{u_2 = \begin{pmatrix} 9 \\ 12 \end{pmatrix}}$$

$$(x_1 - u_2)(x_1 - u_2)^T$$

$$= \begin{pmatrix} 0 \\ -2 \end{pmatrix} \begin{pmatrix} 0 & 2 \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ 0 & 4 \end{pmatrix} = m_1'$$

$$(x_2 - u_2)(x_2 - u_2)^T$$

$$= \begin{pmatrix} 0 \\ 2 \end{pmatrix} \begin{pmatrix} 0 & 2 \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ 0 & 4 \end{pmatrix} = m_2'$$

$$(x_3 - u_2)(x_3 - u_2)^T$$

$$= \begin{pmatrix} -2 \\ 0 \end{pmatrix} \begin{pmatrix} -2 & 0 \end{pmatrix} = \begin{pmatrix} 4 & 0 \\ 0 & 0 \end{pmatrix} = m_3'$$

$$(x_4 - u_2)(x_4 - u_2)^T$$

$$= \begin{pmatrix} 2 \\ 0 \end{pmatrix} \begin{pmatrix} 2 & 0 \end{pmatrix} = \begin{pmatrix} 4 & 0 \\ 0 & 0 \end{pmatrix} = m_4'$$

$$\Sigma_2 = \frac{1}{4} [m_1' + m_2' + m_3' + m_4']$$

$$= \frac{1}{4} \begin{bmatrix} 8 & 0 \\ 0 & 8 \end{bmatrix} = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

$$\Sigma_2^T = \begin{bmatrix} 4/2 & 0 \\ 0 & 4/2 \end{bmatrix}$$

$$|\Sigma_1| = 21$$

$$|\Sigma_2| = 4$$

$$\therefore g(x) \neq$$

$$\Sigma_1^T - \Sigma_2^T = \begin{bmatrix} 11/84 & -1/21 \\ -1/21 & 8/42 \end{bmatrix} - \begin{bmatrix} 4e & 0 \\ 0 & 42 \end{bmatrix}$$

$$\Sigma_1^T - \Sigma_2^T = \begin{bmatrix} -31/84 & -1/21 \\ -1/21 & -5/42 \end{bmatrix}$$

$$\therefore g(x) = \frac{-1}{2} x^T (\Sigma_1^T - \Sigma_2^T) x + (u_1^T \Sigma_1^T - u_2^T \Sigma_2^T) x$$

$$- \frac{1}{2} \ln\left(\frac{21}{4}\right) + \ln\left(\frac{42}{42}\right) - \frac{1}{2} u_1^T \Sigma_1^T u_1 + \frac{1}{2} u_2^T \Sigma_2^T u_2$$

$$\therefore g(x) = \frac{-1}{2} [u_1 \ u_2] \begin{bmatrix} -31/84 & -1/21 \\ -1/21 & -5/42 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$$

$$+ \left(\begin{bmatrix} 9 & 9/2 \end{bmatrix} \begin{bmatrix} 11/84 & -1/21 \\ -1/21 & 8 \end{bmatrix} - \begin{bmatrix} 9 & 12 \end{bmatrix} \begin{bmatrix} 42 & 0 \\ 0 & 42 \end{bmatrix} \right) \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$$

$$- \frac{1}{2} \left(\begin{bmatrix} 9 & 9/2 \end{bmatrix} \begin{bmatrix} 11/84 & -1/21 \\ -1/21 & 8 \end{bmatrix} \begin{bmatrix} 9 \\ 9/2 \end{bmatrix} \right) + \frac{1}{2} \left(\begin{bmatrix} 9 & 12 \end{bmatrix} \begin{bmatrix} 42 & 0 \\ 0 & 42 \end{bmatrix} \begin{bmatrix} 9 \\ 12 \end{bmatrix} \right)$$

$$- \frac{1}{2} \ln\left(\frac{21}{4}\right) + 0$$

$$\therefore g(x) =$$

* Metrics in Classification *

NOTE: when we have balanced dataset at that time we usually prefer accuracy as our metric.
for imbalance dataset we can either choose precision or recall. or F1-score.

(1) confusion Metric :- → used in binary classification.
Actual values.

		1	0
Predicted values	1	TP	FP
	0	FN	TN

Type 1 error (FPR)

$$FPR = \frac{FP}{FP + TN}$$

→ our aim should be reducing type 1 and type 2 error.

$$\rightarrow \text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN}$$

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

(2) Recall :- → It also mentioned by TPR, Sensitivity.

→ used for imbalance dataset.

→ It signifies that out of all actual positive values how much we predicted correctly.

$$\rightarrow \text{Recall} = \frac{TP}{TP + FP}$$

- (3) Precision :- → used when there is imbalance
- It is also known as positive prediction value.
 - Out of ~~the~~ positively predicted values what % of values are actually positive.
 - Precision =
$$\frac{TP}{TP + FN}$$

*NOTE :- When to use what?

- 1) whenever (FP) is more important at that time we will use Precision.
- 2) whenever (FN) is more important at that time we will use Recall.

Example:- In spam detection we will use Precision as a measure and in cancer detection we will use recall as a measure.

(4) F-Beta Score :-

→ Our main aim is to select (β) value.

$$F_{Beta} = \frac{(1+\beta^2) \cdot \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$F_1\text{-score} = \frac{2 \cdot \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

$\uparrow \beta=1$

Harmonic mean.

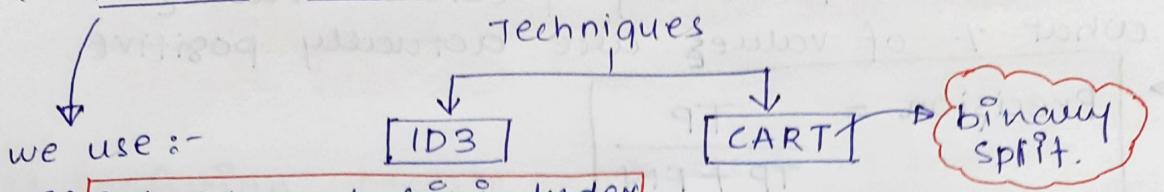
NOTE :- When to choose $\beta=1$?

- When $\beta \neq (FP)$ and (FN) both are equally important.
- When (FP) more imp. we reduce β value.
- When (FN) more imp we ↑ β value.

* Decision Tree *

NOTE:- Decision Tree classifier can be used for both classification and regression.

(*) decision Tree classifier :-



a) Entropy and Gini Index.

b) Information gain

purity of split.

helps us to decide that using which feature we need to split.

Example:- Basic Decision Tree

age = 14

if age ≤ 15

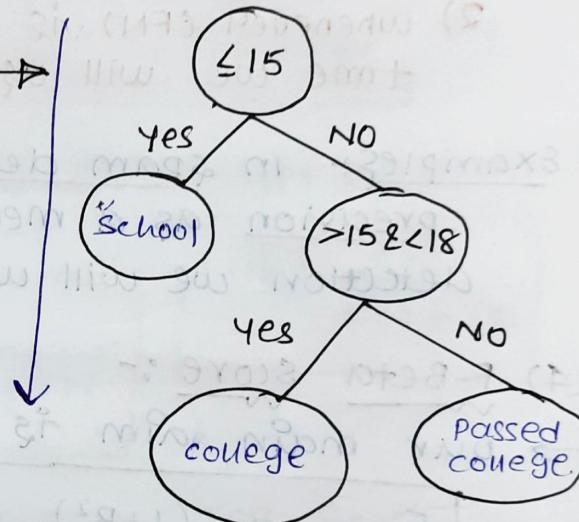
pf("School");

elif age ≥ 15 and ≤ 18

pf("college");

else

pf("Passed college");



- Leaf Node :-

→ They are also called as "terminal node".

→ "Nodes that don't split into more nodes."

→ classes are assigned by majority vote.

→ classes are purely split.

→ for impure split we will further split the node with other features.

Q. How to check the split is pure or not?
→ using entropy and gini index / gini impurity.

Q. How the features are selected?
→ information gain.

(*) Entropy :-

Assuming binary classification,

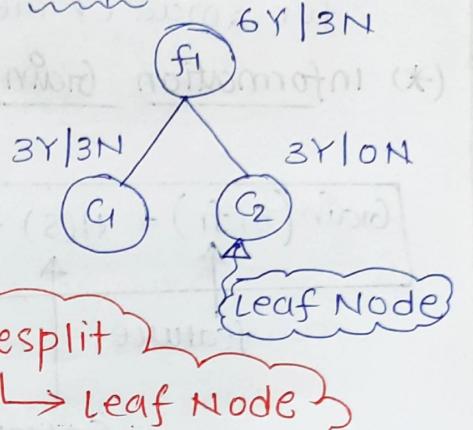
$$H(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

$$H(C_2) = -\frac{3}{3} \log_2 \frac{3}{3} = 0$$

$$\underline{\underline{H(C_2) = 0}}$$

{if $H(S) = 0$ → pure split}

Example :-



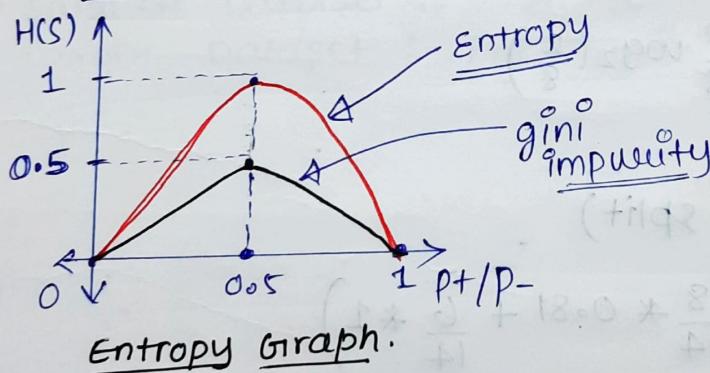
$$H(C_1) = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6}$$

$$= -\log_2 \frac{1}{2} = -(\log_2 1 - \log_2 2) = 1.$$

$$\underline{\underline{H(C_1) = 1}}$$

{if $H(S) = 1$ → highly impure split}

(*) Graphically {Entropy / gini} :-



*. NOTE:-

Entropy ranges from 0 to 1

(*) Gini Impurity :-

$$\text{Gini Impurity} = 1 - \sum_{i=1}^n (p_i)^2$$

Assume binary classification,

$$\text{Gini Impurity} = 1 - ((p_+)^2 + (p_-)^2) =$$

NOTE: Gini Impurity ranges from 0 to 0.5.

pure split highly impure split.

Q. what to use Entropy or Gini Impurity?

→ for large dataset :- use gini impurity
 Reason :- Log operation in entropy will take more time than the algebraic operation in gini impurity.

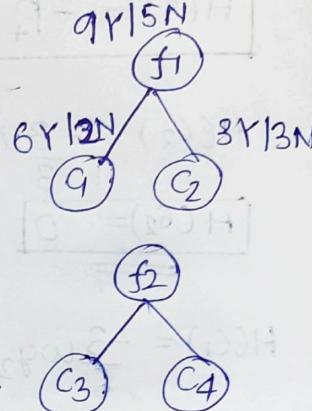
for small or medium dataset :- use entropy

(*) information Gain :- (It tells us about feature selection).

GAIN

$$\text{Gain}(S, f_i) = H(S) - \sum_{\text{Value}} \frac{|S_v|}{|S|} \cdot H(S_v)$$

↑ feature ↑ Entropy of root node ↑ Total sum. ↑ Entropy of category (v).



$$H(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-$$

$$H(S) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14}$$

$$H(S) \approx 0.94$$

$$H(C_1) = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \left(\frac{2}{8}\right)$$

$$H(C_1) \approx 0.81$$

$$H(C_2) = 1 \quad (\because \text{impure split})$$

$$\text{Gain}(S, f_1) = 0.94 - \left(\frac{8}{14} * 0.81 + \frac{6}{14} * 1 \right)$$

$$\text{Gain}(S, f_1) \approx 0.049$$

Similarly, calculate $\text{Gain}(S, f_2)$.

if $\text{Gain}(S, f_2) > \text{Gain}(S, f_1)$

go with f_2 feature.

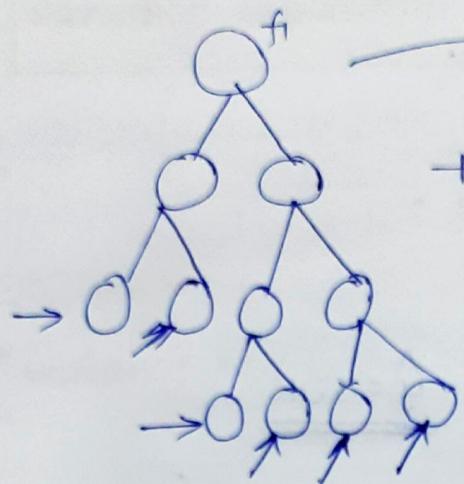
else

go with f_1 feature.

How to select feature?

* Post and Pre-pruning in Decision Tree *

Assume Decision Tree,



- This kind of D.T. may lead to overfitting.
- So, to overcome it.
we use i) Post Pruning.
ii) Pre Pruning.

(I) Post Pruning :- → also known as backward pruning.

→ Doing hyperparameter tuning after construction of tree.

→ Example :- max-depth., min-samples-split

(II) Pre-Pruning :-

→ Doing hyperparameter tuning simultaneously while constructing decision tree.

Q. when to use what?

- Smaller Dataset :- Post Pruning.
- Larger Dataset :- Pre Pruning.

PrePruning is more efficient among both.

Because it does not construct whole D.T.

*. Decision Tree Regressor *.

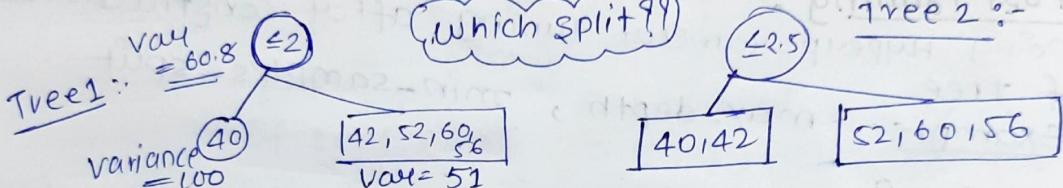
sample dataset :-

Experience Gap salary.

2	Yes	40K
2.5	Yes	42K
3	No	52K
4	No	60K
4.5	Yes	56K

Independent features

$$\text{output.} \quad \hat{y} = 50K$$



→ we will use variance reduction

$$\text{variance} = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \quad \bar{y} = \text{avg. output.}$$

$$\text{variance (Root)} = \frac{1}{5} (100 + 64 + 4 + 100 + 36) \\ = \frac{304}{5} = \underline{\underline{60.8}}$$

For Tree 1:-

Variance (child 1)

$$= \frac{1}{2} + (40 - 50)^2 \\ = \underline{\underline{100}}$$

For Tree 2:-

var(C₁)

$$= \frac{1}{2} [100 + 64] \\ = \frac{164}{2} = \underline{\underline{82}} \checkmark$$

Variance (child 2)

$$= \frac{1}{4} [64 + 4 + 100 + 36] \\ = \frac{204}{4} = \underline{\underline{51}}$$

var(C₂)

$$= \frac{1}{3} [4 + 100 + 36] \\ = \frac{140}{3} = \underline{\underline{46.67}}$$

* variance reduction :-

$$\text{variance reduction} = \text{Var}(\text{root}) - \sum w_i \text{Var}(c_i)$$

For Tree 1:

$$\begin{aligned}\text{Variance reduction} &= 60.8 - \left(\frac{1}{5} * 100 + \frac{4}{5} * 51 \right) \\ &= 60.8 - \frac{304}{5} = \underline{\underline{0}}\end{aligned}$$

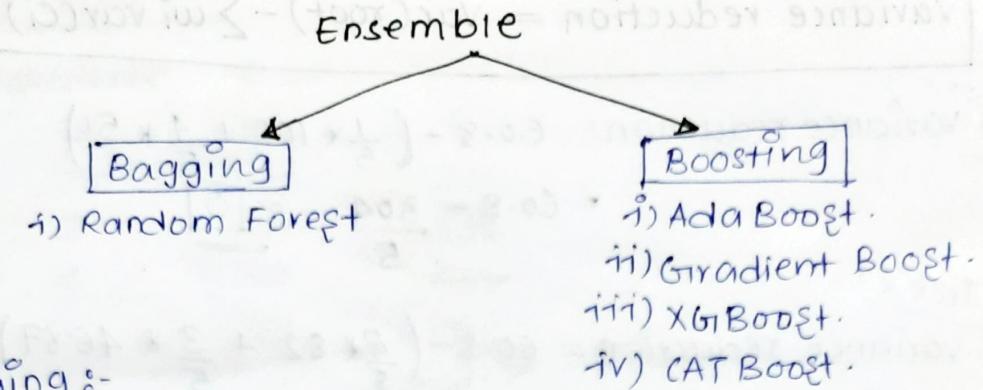
for Tree 2:

$$\begin{aligned}\text{Variance reduction} &= 60.8 - \left(\frac{2}{5} * 82 + \frac{3}{5} * 46.67 \right) \\ &= 60.8 - \left(\frac{164 + 139.98}{5} \right) = \underline{\underline{0.004}}\end{aligned}$$

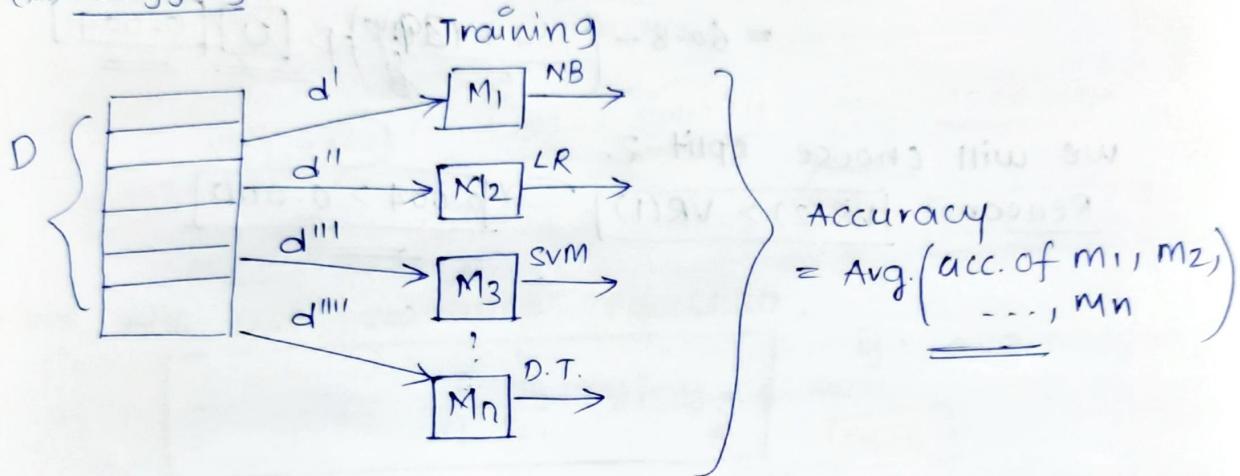
We will choose split-2.

Reason:- $\boxed{VR(2) > VR(1)}$ $\rightarrow \boxed{0.004 > 0.000}$

* Ensemble Techniques *



(*) Bagging :-

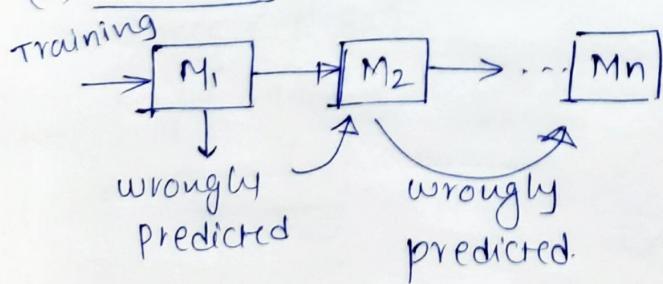


In case of Test Data:

(I) In case of classification, max-voting is used. maximum-voted label will be the output label.

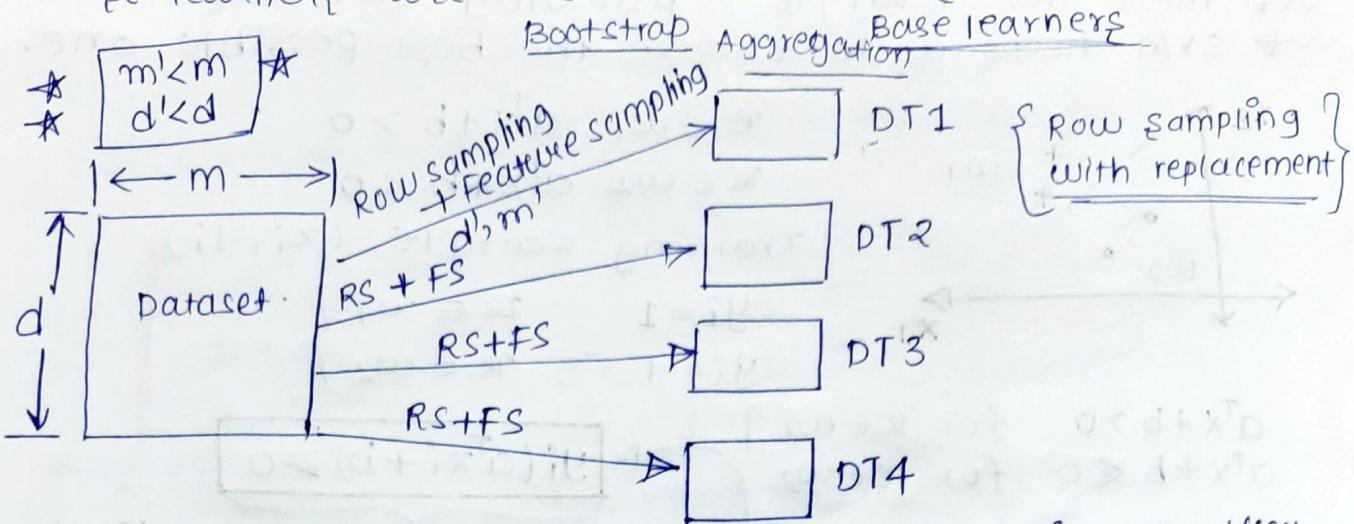
(II) In case of regression, avg of all output is the output of test data.

(*) Boosting :- {weak learners}.



* Random Forest Regression and classification *

→ Base learners are decision tree models.



- In case of classification, we use majority voting classifier.
- In case of regression, we compute avg of all outputs.

Q. why we are using Random Forest over Decision Trees?

→ In Decision Trees, we get

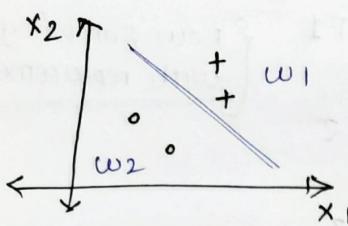
{ Low Bias }
{ High Variance } overfitting.

We need to convert this High variance to low variance in which ~~decisive~~ Random forest helps.

Because, in Random Forest we are using multiple decision trees with Row sampling and Feature sampling.

* Support Vector Machines *

→ There are multiple separating surfaces.
→ SVM helps us to choose the best possible one.



$$x \in w_1 \quad a^T x + b > 0$$

$$x \in w_2 \quad a^T x + b < 0$$

Training sample (x_i, y_i) .

$$y_i = 1$$

$$y_i = -1$$

$$\begin{cases} x \in w_1 \\ x \in w_2 \end{cases}$$

$$\begin{aligned} a^T x + b &> 0 \quad \text{for } x \in w_1 \\ a^T x + b &\leq 0 \quad \text{for } x \in w_2 \end{aligned}$$

$$y_i(a^T x_i + b) \geq 0$$

if correctly classified.

NOTE → Distance of feature vector belongs to class w_1 and class w_2 from the separating surface should be maximize.

- In other words, margin should be maximized.
→ Feature vector near to the separating surface will have less confidence of being correctly classified than the feature vector further from that.

• Separating Plane :-

$a^T x + b = 0$ gives the location of sep. plane.
orthogonal to separating plane.

- change in vector (a) leads to change in orientation of separating plane.
→ change in (b) leads to change in position / loc. of sep. plane.

(*) Benefits of SVM :-

- 1) Robust to outliers.
- 2) Can work on non-linear data.
- 3) Simple.
- 4) sparse solution
- 5) Global optimum solution.

$$a^T x_i + b > d$$

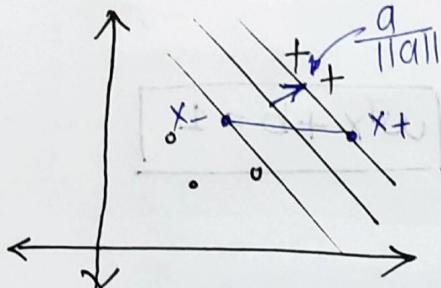
↑
threshold.

$$a^T x_i + b < -d.$$

$$y_i(a^T x_i + b) > d$$

correctly classifying.

$$y_i(a^T x_i + b) \geq 1.$$



$$\text{Margin} \Rightarrow \frac{a^T (x_+ - x_-)}{\|a\|}$$

$$a^T x_+ + b = 1 \quad \text{--- (1)}$$

$$a^T x_- + b = -1. \quad \text{--- (2)}$$

subtract (2) from (1).

$$a^T (x_+ - x_-) = 2 \quad \text{--- (3)}$$

substitute eqn-(3) in (*).

$$\boxed{\text{Margin} \Rightarrow \frac{2}{\|a\|}}$$

from here we observe
that to increase the
margin we have to
 $\downarrow a$.

Now, this becomes a
constraint optimization prblm. It follows $y_i(a^T x_i + b) \geq 1$.
It can be solved using
Langrangian.

$$L = \frac{\|a\|^2}{2} - \sum_i \alpha_i [y_i(a^T x_i + b) - 1]$$

$$\frac{\partial L}{\partial a} = a - \sum_i \alpha_i y_i x_i = 0$$

$$\frac{\partial L}{\partial b} = - \sum_i \alpha_i y_i = 0$$

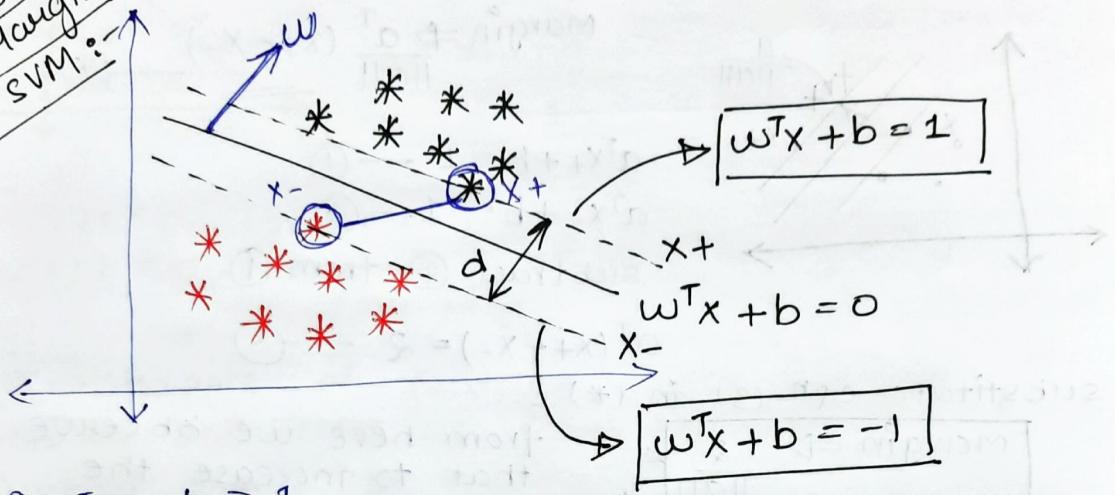
$$\Rightarrow a = \sum_i \alpha_i y_i x_i$$

$$\sum_i \alpha_i y_i = 0$$

for any x^i

$$y_i = \begin{cases} +1 & \text{if } \omega^T x_i + b \geq 0 \\ -1 & \text{if } \omega^T x_i + b < 0 \end{cases}$$

* Hard Margin SVM:



$$\begin{cases} \omega^T x_+ + b \geq 1 \\ \omega^T x_- + b \leq -1 \end{cases}$$

$$y_i(\omega^T x_i + b) = 1 \quad (*)$$

for support vector.

$$\begin{cases} y_i(\omega^T x_i + b) \geq 1 \end{cases}$$

for correctly classifying

from (*)

$$(1)(\omega^T x_+ + b) = 1 \quad (1)$$

$$(-1)(\omega^T x_- + b) = 1 \quad (2)$$

add (1) and (2)

$$\omega^T(x_+ - x_-) = 2 \quad (3)$$

$$\frac{\omega^T(x_+ - x_-)}{\|\omega\|} = d \quad (4)$$

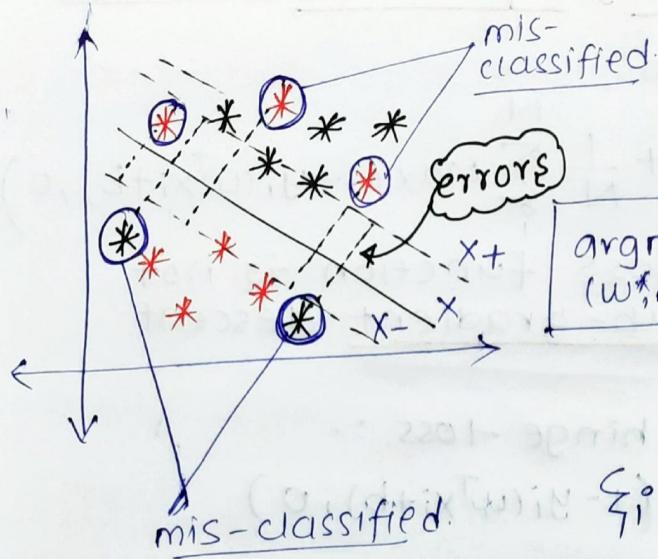
substitute (3) in (4)

$$d = \frac{2}{\|\omega\|}$$

width of margin

$$\boxed{\arg \max_{(\omega^*, b^*)} \frac{2}{\|\omega\|} \text{ such that } y_i(\omega^T x_i + b) \geq 1}$$

(*) Soft Margin SVM :-



$$\underset{(w^*, b^*)}{\operatorname{argmin}} \frac{\|w\|}{2}$$

Hard Margin.

$$\underset{(w^*, b^*)}{\operatorname{argmin}} \frac{\|w\|^2}{2} + C \sum_{i=1}^l \xi_i$$

ξ_i slack variable

$\xi_i = 0 \} \text{ for correctly classified points.}$

SVM error = Margin Error + classification Error

$$J = \underset{(w^*, b^*)}{\operatorname{argmin}} \lambda \frac{\|w\|^2}{2} + C \sum_{i=1}^l \xi_i$$

Hinge loss.
regularized...
subject to $y_i(w^T x_i + b) \geq 1 - \xi_i$ Hyperparameter.
 $\xi_i = \max(1 - y_i(w^T x_i + b), 0)$

~~$\nabla J(w, b) = \lambda w +$~~

NOTE:

- 1) $y_i(w^T x_i + b) > 1$.
 - point outside margin.
 - No contribution to loss.
- 2) $y_i(w^T x_i + b) < 1$.
 - point violates margin constraints.
 - contributes to loss.
- 3) $y_i(w^T x_i + b) = 1$
 - point is on margin.

* NOTE:

- 1) small C allow constraint to be easily ignored (large margin)
- 2) large $C \rightarrow$ narrow margin...
- 3) $C = \infty \rightarrow$ Hard Margin.

(*) Gradient Descent Algo for SVM :-

optimization problem

$$\min_{\omega} J(\omega) = \frac{\lambda}{2} \|\omega\|^2 + \frac{1}{N} \sum_{i=1}^N \max(1 - y_i(\omega^T x_i + b), 0)$$

Because, the hinge loss function is not differentiable, a sub-gradient descent is used.

\Rightarrow sub-gradient for hinge loss :-

$$L(x_i, y_i, \omega) = \max(1 - y_i(\omega^T x_i + b), 0)$$

$$\frac{\partial L}{\partial \omega} = -y_i x_i + \frac{\partial b}{\partial \omega}$$

\Rightarrow do iterative

$$\begin{aligned} \omega_{t+1} &\leftarrow \omega_t - \eta \nabla_{\omega} J(\omega_t) \\ &\leftarrow \omega_t - \eta \frac{1}{N} (\lambda \omega_t + \nabla_{\omega} L(x_i, y_i, \omega)) \end{aligned}$$

$$\omega_{t+1} \leftarrow \omega_t - \eta (\lambda \omega_t - y_i x_i) \text{ if } y_i(\omega^T x_i + b) \leq 1$$

otherwise

$$\leftarrow \omega_t - \eta (\lambda \omega_t)$$

(*) PCA :-

- assume that the data is centred.
which means every variable has mean (0).
- for 1-D projection.
let's take a line with unit vector along it \vec{w} . and projection data vector \vec{x}_i onto the line $\vec{x}_i \cdot \vec{w}$ which is scalar.
- The actual coordinate in p-dimensional space is $(\vec{x}_i \cdot \vec{w}) \vec{w}$.
- mean of the projection will be zero.
since the mean of data vector (\vec{x}_i) is zero
- $$\frac{1}{n} \sum_i^n (\vec{x}_i \cdot \vec{w}) \vec{w} = \left(\left(\frac{1}{n} \sum_i^n \vec{x}_i \right) \vec{w} \right) \vec{w} \quad (*)$$
- error of the projection for data vector \vec{x}_i given by,

$$\begin{aligned}
 \| \vec{x}_i - (\vec{w} \cdot \vec{x}_i) \vec{w} \|^2 &= (\vec{x}_i - (\vec{w} \cdot \vec{x}_i) \vec{w}) (\vec{x}_i - (\vec{w} \cdot \vec{x}_i) \vec{w}) \\
 &= \vec{x}_i \cdot \vec{x}_i - \vec{x}_i \cdot (\vec{w} \cdot \vec{x}_i) \vec{w} \\
 &\quad - (\vec{w} \cdot \vec{x}_i) \vec{w} \cdot \vec{x}_i \\
 &\quad + (\vec{w} \cdot \vec{x}_i) \vec{w} \cdot (\vec{w} \cdot \vec{x}_i) \vec{w} \\
 &= \|\vec{x}_i\|^2 - 2(\vec{w} \cdot \vec{x}_i)^2 + (\vec{w} \cdot \vec{x}_i)^2 \vec{w} \cdot \vec{w} \\
 &= \|\vec{x}_i\|^2 - 2(\vec{w} \cdot \vec{x}_i)^2 + (\vec{w} \cdot \vec{x}_i)^2 \quad [\vec{w} \cdot \vec{w} = 1] \\
 &= \vec{x}_i \cdot \vec{x}_i - (\vec{w} \cdot \vec{x}_i)^2
 \end{aligned}$$

NOW, $MSE(\vec{w}) = \frac{1}{n} \left(\underbrace{\sum_i^n \|\vec{x}_i\|^2}_{\text{independent of } (\vec{w})} - \sum_i^n (\vec{w} \cdot \vec{x}_i)^2 \right)$

will not matter in minimizing MSE.

Thus to minimize the MSE, we need to

maximize second term. i.e.,

$$\text{Max } \left\{ \frac{1}{n} \sum_{i=1}^n (\vec{w} \cdot \vec{x}_i)^2 \right\} \quad \text{Mean of } (\vec{w} \cdot \vec{x}_i)^2$$

(*) Mean of square = square of mean + variance.

$$\text{Max } \frac{1}{n} \sum_{i=1}^n (\vec{w} \cdot \vec{x}_i)^2 = \left(\frac{1}{n} \sum_{i=1}^n \vec{x}_i \cdot \vec{w} \right)^2 + \text{Var}[\vec{w} \cdot \vec{x}_i]$$

we know that mean of projection is zero [From (*)]
Thus, minimizing MSE would result in maximizing
the variance of projections.

⇒ Maximizing the variance:-

- n data vectors
- $X_{n \times p}$ matrix
- projection is given by $XW_{n \times 1}$ matrix.

$$\text{variance is } \sigma_w^2 = \frac{1}{n} \sum_{i=1}^n (\vec{x}_i \cdot \vec{w})^2$$

$$= \frac{1}{n} (XW)^T (XW)$$

$$= \frac{1}{n} W^T X^T X W$$

$$= W^T \underbrace{\frac{1}{n} X^T X}_{\Sigma} W$$

$$\boxed{\sigma_w^2 = W^T \Sigma W}$$

$$\Sigma = \text{covariance matrix}$$
$$\frac{X^T X}{n} = \Sigma \quad \left(\because X \text{ is centered} \right)$$

→ we need to find vector \vec{w} such that it maximizes σ_w^2 . for that we use lagragian multiplier (λ).

$$L(\vec{w}, \lambda) = \sigma_w^2 - \lambda(\vec{w}^T \vec{w} - 1) = W^T \Sigma W - \lambda(W^T W - 1)$$

$$\frac{\partial L}{\partial \lambda} = W^T W - 1 = 0 \quad \frac{\partial L}{\partial W} = 2\Sigma W - 2\lambda W = 0$$

$$\boxed{W^T W = 1}$$

$$\boxed{\Sigma W = \lambda W}$$

$$\boxed{\Sigma w = \lambda w}$$

w is the eigen vector of covariance matrix (Σ).
maximizing the vector is associated with
max. eigen value.

Thus,

The eigen vectors of Σ are the principal components of the data.

→ The 1st principle component with will be the largest eigen vector which goes with largest eigen value. (λ).

* Softmax Regression *

(Multinomial LR)

Softmax function =

$$\sigma(z)_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

$K = \# \text{ class.}$

$$0 < \sigma < 1$$

$$\sigma(z)_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$\sigma(z)_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$\sigma(z)_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$\sigma(z)_1 + \sigma(z)_2 + \sigma(z)_3 = 1$$

In softmax regression, loss is the sum of distance between labels and output probability distribution which is known as cross entropy.

$$\text{cross entropy}_{x_i} = - \sum_{j=1}^K (I(y_i^{\circ}=k) \cdot \log \left(\frac{e^{w^T x_i^{\circ}}}{\sum_{j=1}^K e^{w^T x_j^{\circ}}} \right))$$

$K = \text{classes}$

Total cross-entropy or loss,

$$\text{loss} = - \sum_{i=1}^m \sum_{j=1}^K (I(y_i^{\circ}=k) \cdot \log \frac{e^{w^T x_i^{\circ}}}{\sum_{j=1}^K e^{w^T x_j^{\circ}}})$$

$$\nabla_w \text{loss} = \frac{1}{m} \sum_{i=1}^m x_i^{\circ} \left(I(y_i^{\circ}=k) - \frac{e^{w^T x_i^{\circ}}}{\sum_{j=1}^K e^{w^T x_j^{\circ}}} \right)$$

New parameter for class k ,

$$W_{k+1} = W_k + \alpha \cdot \frac{1}{m} \sum_{i=1}^m x_i^{\circ} \left(I(y_i^{\circ}=k) - \frac{e^{w^T x_i^{\circ}}}{\sum_{j=1}^K e^{w^T x_j^{\circ}}} \right)$$

*. K-means clustering :-

Advantages :-

- 1) simple to implement.
- 2) easily adapt new examples.
- 3) guarantees convergence by trying to minimize the total SSE.
- 4) Fast and efficient. complexity $\rightarrow O(n * k * d)$

Disadvantage :-

- 1) clustering outliers.
- 2) Dependent on initial values.
- 3) choosing K-manually.
- 4) Scaling with number of dimensions.

objective Function is

$$\text{Min } J = \sum_{i=1}^m \sum_{j=1}^K w_{ik} \|x_i - \mu_k\|^2$$

↑
centroid of class K.

Two steps of minimization :-

- 1) minimize w_{ik} given that μ_k fixed. (cluster assignmt)
- 2) find μ_k given that w_{ik} fixed. (centroid compute)

$$\frac{\partial L}{\partial w_{ik}} = \sum_{i=1}^m \sum_{j=1}^K \|x_i - \mu_j\|^2$$

$$\Rightarrow w_{ik} = \begin{cases} 1 & \text{if } k = \operatorname{argmin}_j \|x_i - \mu_j\|^2 \\ 0 & \text{otherwise} \end{cases}$$

$$\frac{\partial L}{\partial \mu_k} = 2 \sum_{i=1}^m w_{ik} (x_i - \mu_k) = 0$$

$$\Rightarrow \mu_k = \frac{\sum_{i=1}^m w_{ik} x_i}{\sum_{i=1}^m w_{ik}}$$