

④ Bayes rule minimizes avg prob error / Bayes rule optimum prob

→ probability of error for x , $P(\text{error}_1|x) = \begin{cases} P(w_1|x) & \text{if } w_2 \\ P(w_2|x) & \text{if } w_1 \end{cases}$

Bayes rule ensures $P(\text{error}_1|x)$ is as small as possible so.

that integral will be small,

$$\begin{aligned} P(\text{error}) &= \int_{-\infty}^{\infty} P(\text{error}, x) dx \\ &= \int_{-\infty}^{\infty} P(\text{error}|x) P(x) dx \end{aligned}$$

$$\text{Thus, } P(\text{error}|x) = \min [P(w_1|x), P(w_2|x)]$$

⑤ Select optimal decision (math).

$$\text{Given, } P(w_1) = P(w_2) = 0.5$$

$$P(x|w_1) = N(0, 0.5)$$

$$P(x|w_2) = N(1, 0.5)$$

Here, loss, λ has diagonals 0 but $\lambda_{12} \neq \lambda_{21}$

$$\text{so, we will choose } w_1 \text{ if } \frac{P(x|w_1)}{P(x|w_2)} > \frac{(\lambda_{12} - \lambda_{21}) P(w_2)}{(\lambda_{21} - \lambda_{12}) P(w_1)}$$

Calculate likelihood:

~~Given~~ we know, $N = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$ and $N(\mu, \sigma^2)$

$$\text{so for } P(x|w_1) = \frac{1}{\sqrt{2\pi(0.5)}} e^{-\frac{x^2}{2}} = \frac{e^{-\frac{(x-0)^2}{2}}}{\sqrt{\pi}}$$

$$\text{for, } P(x|w_2) = \frac{\sqrt{2}}{\sqrt{2\pi}} \cdot e^{-\frac{(x-1)^2}{2}} = \frac{e^{-\frac{(x-1)^2}{2}}}{\sqrt{\pi}}$$

So, likelihood ratio:

$$\frac{P(x|w_1)}{P(x|w_2)} = \frac{\frac{1}{\sqrt{2\pi}} e^{-x^2 + (x-1)^2}}{\frac{1}{\sqrt{2\pi}} e^{-x^2 + x^2 - 2x + 1}} = e^{-2x+1} = \ell$$

Now,

$$e^{-2x+1} > \frac{(0.5 - 0) \times 0.5}{(1 - 0) \times 0.5}$$

$$\Rightarrow e^{-2x+1} > \frac{1}{2}$$

$$\Rightarrow -2x+1 > \ln(\frac{1}{2})$$

$$\Rightarrow -2x+1 > -0.69$$

$$\Rightarrow -2x+1.69 > 0$$

$$\Rightarrow 1.69 > 2x$$

$$\Rightarrow x < 0.845$$

Thus, if $x < 0.845$ then w_1 , else w_2 .

$$\text{Likelihood Ratio} = \frac{P(x|w_1)}{P(x|w_2)} = \frac{e^{-2x+1}}{e^{-x^2 + x^2 - 2x + 1}} = e^{2x-1}$$

* Proof 2D Gaussian/Normal Df are two 1D Gaussians if uncorrelated

$$\text{Let, } \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad \boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$$

$$p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2} \left[\begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix}^T \begin{bmatrix} \frac{1}{\sigma_1^2} & 0 \\ 0 & \frac{1}{\sigma_2^2} \end{bmatrix} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix} \right]}$$

$$\text{Denominator: } (2\pi)^{1/2} \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}^{1/2} = 2\pi \sqrt{(\sigma_1^2 \sigma_2^2)} \\ = \sqrt{4\pi^2} \sqrt{\sigma_1^2 \sigma_2^2} \\ = \sqrt{2\pi \sigma_1^2} \cdot \sqrt{2\pi \sigma_2^2}$$

$$\text{Exponential index: } \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix}^T \begin{bmatrix} \frac{1}{\sigma_1^2} & 0 \\ 0 & \frac{1}{\sigma_2^2} \end{bmatrix} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix}$$

$$= \left[\frac{x_1 - \mu_1}{\sigma_1^2} - \frac{x_2 - \mu_2}{\sigma_2^2} \right] \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix}$$

$$= \frac{(x_1 - \mu_1)^2}{\sigma_1^2} + \frac{(x_2 - \mu_2)^2}{\sigma_2^2}$$

$$\text{Exponent: } e^{-\frac{1}{2} \left[\frac{(x_1 - \mu_1)^2}{\sigma_1^2} + \frac{(x_2 - \mu_2)^2}{\sigma_2^2} \right]} \\ = e^{-\frac{(x_1 - \mu_1)^2}{2\sigma_1^2}} \cdot e^{-\frac{(x_2 - \mu_2)^2}{2\sigma_2^2}}$$

$$\text{So, } p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{2\pi \sigma_1^2}} e^{-\frac{(x_1 - \mu_1)^2}{2\sigma_1^2}} \cdot \frac{1}{\sqrt{2\pi \sigma_2^2}} e^{-\frac{(x_2 - \mu_2)^2}{2\sigma_2^2}}$$

* Bayes decision rule interprets $\rightarrow w_1$ if likelihood ratio exceeds a threshold value independent of observation of x .

Here

actions = a_1, a_2

states / class = w_1, w_2

loss, $\lambda_{ij} \Rightarrow \lambda(a_i | w_j)$ so $\lambda(a_1 | w_2) \rightarrow$ loss of w_1 when we decide w_2
 $\lambda(a_2 | w_1) \rightarrow$ " " w_2 " " w_1

Conditional risk :

$$R(a_1 | x) = \sum_{j=1}^c \lambda(a_1 | w_j) P(w_j | x)$$

$$\text{So, } R(a_1 | x) = \lambda_{11} P(w_1 | x) + \lambda_{12} P(w_2 | x)$$

$$R(a_2 | x) = \lambda_{21} P(w_1 | x) + \lambda_{22} P(w_2 | x)$$

Minimum risk decision :

Decide w_1 if $R(a_1 | x) < R(a_2 | x)$:

$$\text{or, } \lambda_{11} P(w_1 | x) + \lambda_{12} P(w_2 | x) < \lambda_{21} P(w_1 | x) + \lambda_{22} P(w_2 | x)$$

$$\text{or, } (\lambda_{21} - \lambda_{11}) P(w_1 | x) > (\lambda_{22} - \lambda_{12}) P(w_2 | x)$$

$$\text{or, } \frac{P(w_1 | x)}{P(w_2 | x)} > \frac{(\lambda_{12} - \lambda_{22})}{(\lambda_{21} - \lambda_{11})} \frac{P(w_2)}{P(w_1)}$$

Decision rule

$$\frac{(w_1 - w_2)^T \theta}{(w_1 - w_2)^T \theta + (w_2 - w_1)^T \theta}$$

Multivariate Dis fn :- (general)

$$q_i(x) = -\frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln(\Sigma^{-1}) + \ln P(w_i)$$

Case 1 $\Sigma = \sigma^2 I = \begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix}$ $\Sigma^{-1} = \begin{bmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{1}{\sigma^2} \end{bmatrix} = \frac{1}{\sigma^2} I$

$$(x - \mu_i)^T \Sigma^{-1} (x - \mu_i) = (x_1 - \mu_{1i})^2 + (x_2 - \mu_{2i})^2 = \|x - \mu_i\|^2$$

$$q_i(x) = \frac{\|x - \mu_i\|^2}{2} + \ln P(w_i)$$

$$= \frac{(x_1 - \mu_{1i})^2 + (x_2 - \mu_{2i})^2}{2} + \ln P(w_i)$$

Decision boundary: $w^T(x - x_0) = 0$

$$w = (\mu_1 - \mu_2)$$

$$x_0 = \frac{(\mu_1 - \mu_2)}{2} - \frac{\sigma^2}{\|\mu_1 - \mu_2\|^2} \cdot \ln \left(\frac{P(w_1)}{P(w_2)} \right) (\mu_1 - \mu_2)$$

Case 2 $\Sigma_i = \Sigma \rightarrow \begin{bmatrix} a & \text{cov} \\ \text{cov} & b \end{bmatrix}$ (use mahalanobis)

Decision boundary: $D_m = (x - \mu)^T \Sigma^{-1} (x - \mu) = 0$

Decision boundary: $w^T(x - x_0) = 0$

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

case 3] Σ = Arbitrary

যদি Σ অবশ্যই

general equation for case 3*

We know,

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

$$\begin{aligned} \text{Here, } (x - \mu_i)^T \Sigma^{-1} (x - \mu_i) &= x^T \Sigma^{-1} x + \mu_i^T \Sigma^{-1} \mu_i - \mu_i^T \Sigma^{-1} x - x^T \Sigma^{-1} \mu_i \\ &= x^T \Sigma^{-1} x - 2 \mu_i^T \Sigma^{-1} x + \mu_i^T \Sigma^{-1} \mu_i \end{aligned}$$

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

~~Let,~~ $W_i = -\frac{\Sigma^{-1}}{2}$

$w_i = \mu_i^T \Sigma^{-1}$

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

$$\text{Thus, } g_i(x) = x^T W_i x + w_i x + w_{i0}$$

$$\text{Decision boundary} \rightarrow g_1(x) = g_2(x)$$

CS CamScanner

* Unknown mean μ : soft margins bound probabilistic approach
 assume $p(x|w_j) \rightarrow N(\mu_j, \Sigma_j)$.

$$\text{let } \Theta_j = [\mu_j, \Sigma_j]^T$$

Now, in Dataset D,

$$\begin{aligned} p(D|\Theta) &= p(x_1|\Theta) * p(x_2|\Theta) * \dots * p(x_n|\Theta) \\ &= \prod_{k=1}^n p(x_k|\Theta). \end{aligned}$$

loglikelihood, $L(\Theta) = \ln p(D|\Theta) = \sum_{k=1}^n \ln p(x_k|\Theta)$

$$\text{or, } \nabla_\Theta L = \prod_{k=1}^n \nabla_\Theta \ln p(x_k|\Theta) = -\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu).$$

we know, $p(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$

for a sample point,

$$p(x_k|\mu) = \underbrace{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}_{\text{unknown } \mu} - \underbrace{\frac{1}{2} \ln [(2\pi)^d |\Sigma|]}_{\Sigma \text{ known}}$$

$$\begin{aligned} \text{Here, } \nabla_\mu \ln p(x_k|\mu) &= \frac{\partial}{\partial \mu} \left[-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu) \right] \\ &= \frac{\partial}{\partial \mu} \left[-\frac{1}{2} (x^T \Sigma^{-1} x - 2\mu^T \Sigma^{-1} x + \mu^T \Sigma^{-1} \mu) \right] \\ &= -\frac{1}{2} (-2 \Sigma^{-1} x + \Sigma^{-1} \mu) \\ &= \Sigma^{-1} x - \Sigma^{-1} \mu \\ &= \Sigma^{-1} (x-\mu) \end{aligned}$$

$$\text{Thus, } \sum_{k=1}^n \Sigma^{-1} (\mathbf{x}_k - \boldsymbol{\mu}) = 0 .$$

$$\Rightarrow \sum_{k=1}^n \Sigma^{-1} \mathbf{x}_k - \sum_{k=1}^n \Sigma^{-1} \boldsymbol{\mu} = 0 .$$

$$\Rightarrow \sum_{k=1}^n \mathbf{x}_k = \sum_{k=1}^n \boldsymbol{\mu} .$$

$$\Rightarrow \sum_{k=1}^n \mathbf{x}_k = \boldsymbol{\mu} \cdot n .$$

$$\Rightarrow \boldsymbol{\mu} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k .$$

* Unknown mean $\boldsymbol{\mu}$ and covariance Σ :-

$$\text{Let } \Theta_1 = \boldsymbol{\mu} \text{ and } \Theta_2 = \Sigma$$

$$\text{We know, } p(\mathbf{x}) = \frac{1}{(2\pi)^d |\Sigma|^{1/2}} e^{-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})} .$$

$$p(\mathbf{x}) = \frac{e^{-\frac{(\mathbf{x} - \boldsymbol{\mu})^2}{2\sigma^2}}}{\sqrt{2\pi|\Sigma|}}$$

$$p(x_k | \Theta) = \frac{1}{2} (\mathbf{x} - \boldsymbol{\theta}_1)^T .$$

$$\text{For univariate, } \ln p(x_k | \Theta) = -\frac{1}{2\sigma^2} (x_k - \boldsymbol{\mu})^2 - \frac{1}{2} \ln (2\pi|\Sigma|) .$$

$$\text{my work} = -\frac{1}{2\theta_2} (x_k - \theta_1)^2 - \frac{1}{2} \ln (2\pi \theta_2) .$$

$$\text{So, } \nabla_{\Theta} \ln p(x_k | \Theta) = \left[\begin{array}{l} \frac{1}{\theta_2} (x_k - \theta_1) \\ \frac{1}{2\theta_2} (x_k - \theta_1)^2 - \frac{1}{2\theta_2} \end{array} \right]$$

KKT wise,

$$\sum_{k=1}^n \frac{1}{\theta_2} (x_k - \theta_1) = 0 \quad \text{--- ①}$$

$$\sum_{k=1}^n \left[\frac{1}{2\theta_2} (x_k - \theta_1)^2 - \frac{1}{2\theta_2} \right] = 0 .$$

For ①,

$$\sum_{k=1}^n (x_k - \theta_1) = 0$$

$$\Rightarrow \sum_{k=1}^n x_k = \sum_{k=1}^n \theta_1$$

$$\Rightarrow \sum_{k=1}^n x_k = n\theta_1$$

$$\Rightarrow \theta_1 = \frac{1}{n} \sum_{k=1}^n x_k$$

$$\Rightarrow \hat{\mu} = \frac{1}{n} \sum_{k=1}^n x_k$$

For ②,

$$\sum_{k=1}^n \left[\frac{(x_k - \theta_1)^2}{2\theta_2^2} - \frac{1}{2\theta_2^2} \right] = 0$$

$$\Rightarrow \sum_{k=1}^n \left[\frac{(x_k - \theta_1)^2}{2\theta_2^2} - \frac{1}{\theta_2^2} \right] = 0$$

$$\Rightarrow \sum_{k=1}^n \frac{(x_k - \theta_1)^2}{\theta_2^2} = \sum_{k=1}^n \frac{1}{\theta_2^2}$$

$$\Rightarrow \sum_{k=1}^n (x_k - \theta_1)^2 = \sum_{k=1}^n \theta_2^2$$

$$\Rightarrow \sum_{k=1}^n (x_k - \hat{\mu})^2 = n\theta_2^2$$

$$\Rightarrow \sum_{k=1}^n (x_k - \hat{\mu})^2 = n\hat{\theta}_2^2$$

$$\Rightarrow \hat{\theta}_2^2 = \frac{1}{n} \sum_{k=1}^n (x_k - \hat{\mu})^2$$

* Perceptron convergence theorem:

Assumptions:

1. Dataset is linearly separable with γ margin
2. Dataset D where for all x belongs to D , $\forall x \in D$ and $\|x_i\| \leq R$. when $R > 0$.
3. If w is the solⁿ vector that separates data with γ margin then each data can be rescaled $\rightarrow w^*$ and $\|w^*\| = 1$.
 $\|x_i\| \leq 1, x_i \in D$.

4. $x=1$

5. $x \rightarrow$ features

$y \rightarrow$ labels

6. $w_0 \leftarrow 0$.

Proof: w^L is weight vector wrt (x, y) , for which mistake again if $y(w^L \cdot x) \leq 0$. occur.
 update: $w^{L+1} = w^L + xy$.

Now, upperlimit:

$$\begin{aligned} \|w^{L+1}\|^2 &= \|w^L + xy\|^2 \\ &= (w^L + xy)^T (w^L + xy) \\ &= \|w^L\|^2 + w^{L T} xy + x^T w^L y + x^T x y^2 \\ &= \|w^L\|^2 + 2((w^L)^T x)y + \|x\|^2 y^2 \end{aligned}$$

$$\begin{aligned} \text{Now, } \|w^{L+1}\|^2 &\leq \|w^L\|^2 + R^2 \\ &\leq (\|w^{L-1}\|^2 + R) + R^2. \end{aligned}$$

$$\therefore \|w^{L+1}\|^2 \leq \boxed{LR^2} \quad \text{--- 1}$$

Again lowerlimit:

$$\begin{aligned} (w^{L+1})^T \cdot w^* &= (w^L + xy)^T \cdot w^* \\ &= w^L w^* + ((w^*)^T x)y \end{aligned}$$

$$\begin{aligned} \text{So, } (w^{L+1})^T w^* &= (w^L)^T w^* + y \\ &= (w^{L-1})^T w^* + y + y \end{aligned}$$

$$\therefore (w^{L+1})^T w^* = (w^0)^T w^* + 2y \quad \text{--- 2}$$

In ②,

$$L^2 \gamma^2 \leq \| \omega^{L+1} \|^2 \| \omega^* \|^2$$

$$\text{or, } L^2 \gamma^2 \leq \| \omega^{L+1} \|^2 \cdot [(\because \| \omega^* \|^2 = 1)] \quad \text{--- ③}$$

In ① and ③,

$$L^2 \gamma^2 \leq \| \omega^{L+1} \|^2 \leq L^2 R^2$$

$$\text{or, } L^2 \gamma^2 \leq L^2 R^2$$

$$\text{or, } L^2 \leq R^2 \leq \frac{1}{\gamma^2}$$

$$\text{or } L \leq \frac{1}{\gamma^2}$$

(Proved)

Perceptron: (Math)

Batch : $I_{tr} \rightarrow a^T u = [\quad] \rightarrow \text{sign}(a^T u) \rightarrow \text{classify}$

if misclassified \rightarrow update

$$a = a + \eta \sum y_i u_i$$

Fixed Incon:

$I_{tr} \rightarrow a^T u = [\quad] \rightarrow \text{sign}(a^T u) \rightarrow \text{mis.}$

$$a = a + \eta y_i u_i$$

$$= [\quad] + () \times () \times [\quad]$$

(Update করে আবারু $x_i - \text{তত্ত্ব. Start}$)

ROC - AUC Math

$\text{Score} < \text{threshold} \rightarrow \text{not spam}$
 $\text{Score} \geq " " \rightarrow \text{spam.}$

Forc Thres = 1.0,

+	0	5
-	1	4

$$FPR = \frac{FP}{TN+FP} = \frac{1}{4+1} = 0.2$$

$$TPR = \frac{TP}{TP+FN} = \frac{0}{0+1} = 0$$

$$(0.2, 0)$$

Forc Thres = 0.5

+	1	4
-	1	4

$$FPR = \frac{1}{1+4} = 0.2$$

$$TPR = \frac{1}{1+4} = 0.2$$

$$(0.2, 0.2)$$

Forc Thres = 0.8

+	2	3
-	1	4

$$FPR = \frac{1}{1+4} = 0.2$$

$$TPR = \frac{2}{2+3} = \frac{2}{5} = 0.4$$

$$(0.2, 0.4)$$

Forc Thres = 0.7

+	2	3
-	2	3

$$FPR = \frac{2}{2+5} = 0.4$$

$$TPR = 0.4$$

$$(0.4, 0.4)$$

Forc Thres = 0.6

+	3	2
-	2	3

$$FPR = \frac{2}{5} = 0.4$$

$$TPR = \frac{3}{5} = 0.6$$

$$(0.4, 0.6)$$

Truth	Prediction
N	0.1
N	0.2
P	0.3
N	0.4
P	0.5
P	0.6
N	0.7
P	0.8
P	0.9
N	1.0

Thres = 0.5	FPR = $\frac{2}{2+3} = 0.4$	(0.4, 0.8)
-	TPR = $\frac{4}{5} = 0.8$	

Thres = 0.4	FPR = $\frac{3}{4+1} = 0.6$	(0.6, 0.8)
-	TPR = $\frac{4}{5} = 0.8$	

Thres = 0.3

+	-	FPR = $\frac{3}{5} = 0.6$	(0.6, 1)
-	3 2	TPR = $\frac{5}{5} = 1$	

Thres = 0.2

+	-	FPR = $\frac{4}{5} = 0.8$	(0.8, 1)
-	4 5	TPR = $\frac{5}{5} = 1$	

Thres = 0.1

+	-	FPR = 1	(1, 1)
-	5 0	TPR = 1	

AUC (area under curve) $\frac{(F_2 - F_1)(T_1 + T_2)}{2}$

$$\text{Forc } (0.2, 0) (0.2, 0.2) \rightarrow \frac{(0.2 - 0.2)(0 + 0.2)}{2} = \frac{0}{2} = 0$$

$$\text{Forc } (0.2, 0.2) (0.2, 0.4) \rightarrow \frac{(0.2 - 0.2)(0.2 + 0.4)}{2} = 0$$

$$\text{Forc } (0.2, 0.4) (0.4, 0.4) \rightarrow \frac{(0.4 - 0.2)(0.4 + 0.4)}{2} = 0.08$$

$$\text{Forc } (0.4, 0.4) (0.4, 0.6) = 0$$

$$\text{Forc } (0.4, 0.6) (0.4, 0.8) = 0$$

$$\text{Forc } (0.4, 0.8) (0.6, 0.8) = \frac{0.2 \times 1.6}{2} = 0.16$$

$$\text{Forc } (0.6, 0.8) (0.6, 1) = 0$$

$$\text{Forc } (0.6, 1) (0.8, 1) = \frac{0.2 \times 2}{2} = 0.2$$

$$\text{Forc } (0.8, 1) (1, 1) = \frac{0.2 \times 2}{2} = 0.2$$

$$\sum \text{Area} = 0 + 0 + 0.08 + 0 + 0 + 0.16 +$$

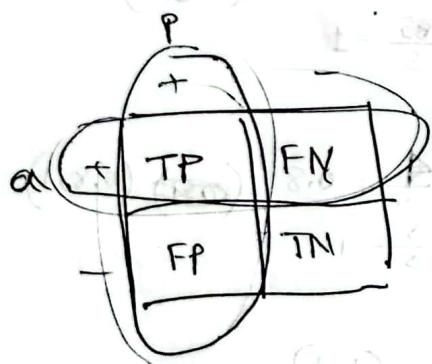
$$0.2 = 0.6$$

$$= 0.69$$

good but not perfect

as $AUC > 0.5$.

but ~~AUC~~ < 1 .



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

$$\text{recall} = \frac{TP}{TP + FN}$$

$$\text{F1 Score} = \frac{2TP}{2TP + FN + FP}$$

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

$$\text{pre} = \frac{TP}{TP + FP}$$

TP	FP
FN	TN

(P.O, N.O)

TP	FP
FN	TN

(P.O, P.O)

- * KMeans Clusters
- ① choose k
 - ② initialize C .
 - ③ C - ~~average~~ distance of all points. + assign to nearest
 - ④ update $C \rightarrow \frac{\text{all points}}{N} \left(\frac{\Sigma x}{N}, \frac{\Sigma y}{N} \right)$
 - ⑤ repeat 3 and 4,
(convergence \rightarrow no data points change cluster)

Given, $P_1(2, 10)$
 $P_6(6, 4)$

$P_2(2, 5)$ $P_3(8, 1)$ $P_4(5, 8)$ $P_5(7, 5)$
 $P_7(1, 2)$ $P_8(4, 9)$ $P_9(3, 6)$ $P_{10}(6, 3)$.

Clusters :-

① $K=3$.

② Clusters

③ Distance:

$C_1(2, 10)$

$C_2(5, 8)$

$C_3(1, 2)$

C_1

$$P_1 = \sqrt{(2-2)^2 + (10-10)^2} = 0.$$

$$P(2, 5) = 5.$$

$$3.6055$$

$$8.06$$

$$4.24$$

$$3.17$$

C_3 .

$$④ \quad S_1 = \{P_1\}$$

$$S_2 = \{P_3, P_4, P_5, P_6, P_8, P_9, P_{10}\}$$

$$S_3 = \{P_2, P_7\}$$

update centroids:

$$C_1 = (2, 10)$$

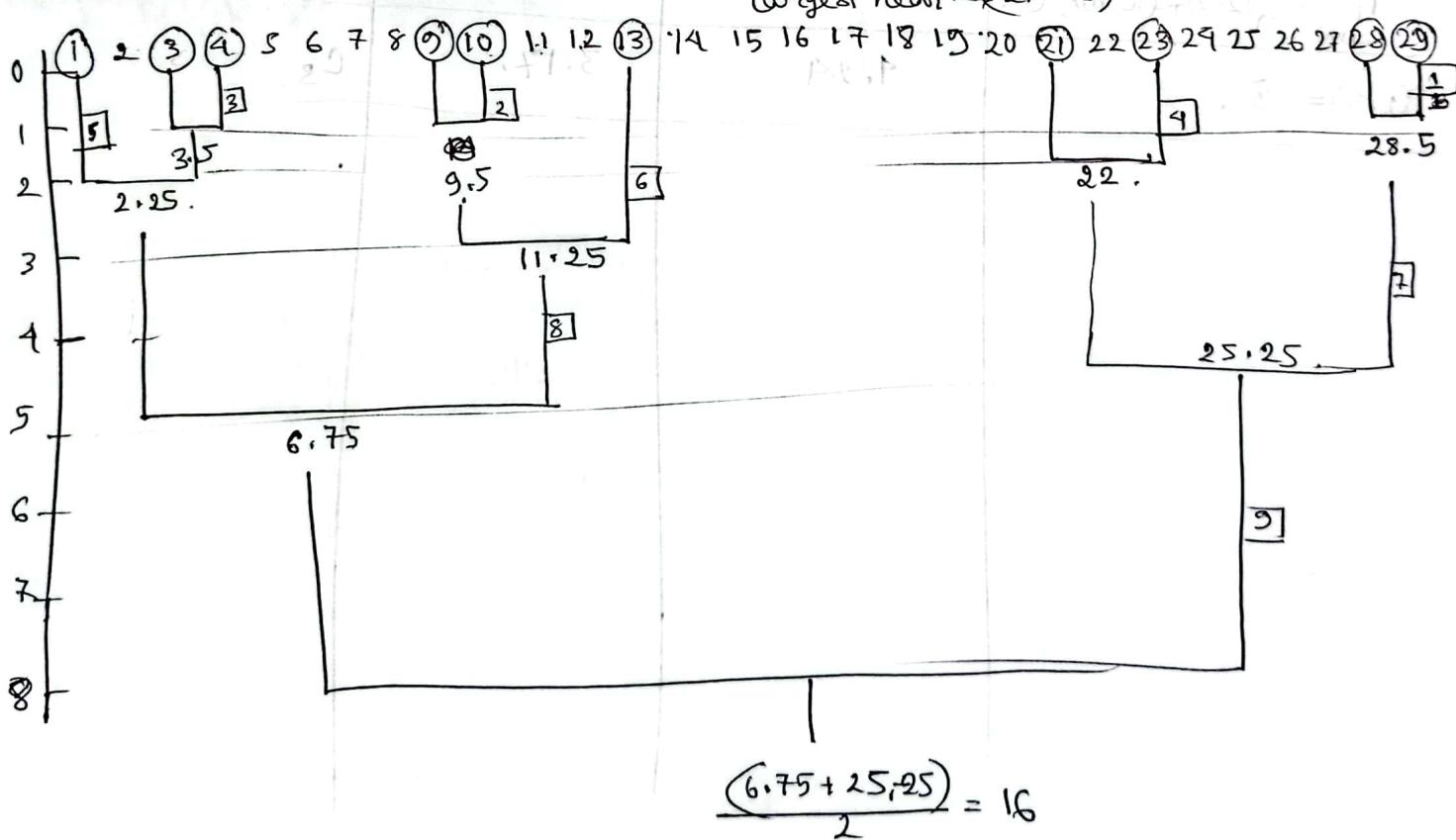
$$C_2 = \left(\frac{8+5+7+6+4+3+9}{7}, \frac{4+8+5+4+9+6+3}{7} \right)$$

$$= (6, 5.57)$$

$$C_3 = \left(\frac{2+1}{2}, \frac{5+2}{2} \right) = (1.5, 3.5)$$

⑤ repeat with new C .

Agglomerative Hierarchical Cluster ? Given, (1, 3, 4, 9, 10, 13, 21, 23, 28, 29).
largest near $(21-13) = 8 \Rightarrow \text{Dist} = 8$



PCA \rightarrow max CA that maximizes variance
 LDA \rightarrow maximize CA for class separation.

LDA \rightarrow 7 steps :-

- ① Class μ_0, μ_1
- ② Overall μ
- ③ Within Class scatter matrix (S_w) $S_i = \sum (x - \mu_i)^T (x - \mu_i)$; $S_w = \sum S_i$
- ④ Between " " (S_B) $S_{B,i} = \sum (\mu_i - \mu)^T (\mu_i - \mu)$; $S_B = \sum S_{B,i}$
- ⑤ $A = S_w^{-1} S_B$.
- ⑥ $|A - \lambda I| = 0$; $\lambda = ?$. Eigenvectors, $w = ?$
- ⑦ project $\rightarrow xw$.

Example:

x_1	x_2	Class
1	2	0
2	3	0
3	1	0
6	5	1
7	7	1
8	6	1

① Class μ_0, μ_1

$$\mu_0 = \left(\frac{1+2+3}{3}, \frac{2+3+1}{3} \right) = (2, 2)$$

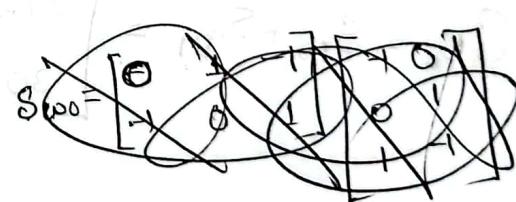
$$\mu_1 = \left(\frac{6+7+8}{3}, \frac{5+7+6}{3} \right) = (7, 6)$$

② Overall μ :

$$\mu = (4.5, 4)$$

③ S_w :

$$S_w = \begin{bmatrix} 1 & 2 \\ 2 & 3 \\ 3 & 1 \end{bmatrix} - [2, 2] = \begin{bmatrix} -1 & 0 \\ 0 & 1 \\ 1 & -1 \end{bmatrix}$$



$$S_{w0} = \begin{bmatrix} -1 & 0 & 1 \\ 0 & 1 & -1 \end{bmatrix} \begin{bmatrix} -1 & 0 \\ 0 & 1 \\ 1 & -1 \end{bmatrix} = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$$

$$\mu_1 - \mu_1 = \begin{bmatrix} 6 & 5 \\ 7 & 7 \\ 8 & 6 \end{bmatrix} - \begin{bmatrix} 7 & 6 \end{bmatrix} = \begin{bmatrix} -1 & -1 \\ 0 & 1 \\ 1 & 0 \end{bmatrix}; S_{W1} = \begin{bmatrix} -1 & 0 & 1 \\ 1 & 1 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} -1 & -1 \\ 0 & 1 \\ 1 & 0 \end{bmatrix} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$$

$$S_W = S_{W0} + S_{W1} = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} + \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} = \begin{bmatrix} 4 & 0 \\ 0 & 4 \end{bmatrix}$$

④ S_B :

$$\mu_0 - \mu = \begin{bmatrix} 2 & 2 \end{bmatrix} - \begin{bmatrix} 4.5 & 1 \end{bmatrix} = \begin{bmatrix} -2.5 & -2 \end{bmatrix}$$

$n_0 = 3$

$$S_{B0} = \begin{bmatrix} -2.5 & -2 \end{bmatrix} \times \frac{1}{2} = \begin{bmatrix} -2.5 \\ 2 \end{bmatrix} = \begin{bmatrix} 2.5 \\ -2 \end{bmatrix}$$

$$S_{B0} = 3 \times \begin{bmatrix} -2.5 \\ -2 \end{bmatrix} \begin{bmatrix} -2.5 & -2 \end{bmatrix} = \begin{bmatrix} 6.25 \\ 5 \end{bmatrix}$$

$$\begin{bmatrix} 5 \\ 4 \end{bmatrix} \times 3 = \begin{bmatrix} 18.75 \\ 15 \end{bmatrix}$$

$$\mu_1 - \mu = \begin{bmatrix} 7 \\ 6 \end{bmatrix} - \begin{bmatrix} 4.5 \\ 4 \end{bmatrix} = \begin{bmatrix} 2.5 \\ 2 \end{bmatrix}$$

$$S_{B1} = 3 \times \begin{bmatrix} 2.5 \\ 2 \end{bmatrix} \begin{bmatrix} 2.5 & 2 \end{bmatrix} = \begin{bmatrix} 6.25 \\ 5 \end{bmatrix}$$

$$\begin{bmatrix} 5 \\ 4 \end{bmatrix} \times 3 = \begin{bmatrix} 18.75 \\ 15 \end{bmatrix}$$

$$S_{B0} = S_{B0} + S_{B1} = \begin{bmatrix} 37.5 & 30 \\ 30 & 24 \end{bmatrix}$$

⑤ $S_W^{-1} S_B \approx A$

$$S_W^{-1} = \frac{\text{adj}(S_W)}{|S_W|} = \frac{1}{16} \begin{bmatrix} 4 & 0 \\ 0 & 4 \end{bmatrix} = \begin{bmatrix} 0.25 & 0 \\ 0 & 0.25 \end{bmatrix}$$

$$\therefore S_W^{-1} S_B = \begin{bmatrix} 0.25 & 0 \\ 0 & 0.25 \end{bmatrix} \begin{bmatrix} 37.5 & 30 \\ 30 & 24 \end{bmatrix} = \begin{bmatrix} 9.375 & 7.5 \\ 7.5 & 6 \end{bmatrix}$$

$$\textcircled{6} \quad |A - \lambda I| = \begin{bmatrix} 9.375 - \lambda & 7.5 \\ 7.5 & 6 - \lambda \end{bmatrix} = 0$$

$$\Rightarrow (9.375 - \lambda)(6 - \lambda) - 56.25 = 0$$

$$\Rightarrow 56.25 - 15.375\lambda + \lambda^2 - 56.25 = 0$$

$$\Rightarrow \lambda^2 = 15.375\lambda$$

$$\Rightarrow \lambda = 15.375 \rightarrow 0$$

Eigen vectors, $(A - \lambda I)v = 0$

$$\begin{bmatrix} 9.375 - 15.375 & 7.5 \\ 7.5 & 6 - 15.375 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = 0$$

$$\Rightarrow -6v_1 + 7.5v_2 = 0 \quad \Rightarrow 7.5v_1 - 9.375v_2 = 0$$

$$\Rightarrow v_2 = 0.8v_1 \quad \Rightarrow v_2 = 0.8v_1$$

$$\text{So, } w = \begin{bmatrix} 1 \\ 0.8 \end{bmatrix}$$

LDA scores.

$$\textcircled{7} \quad Xw = \begin{bmatrix} 1 & 2 \\ 2 & 3 \\ 3 & 1 \\ 6 & 5 \\ 7 & 7 \\ 8 & 6 \end{bmatrix} \begin{bmatrix} 1 \\ 0.8 \end{bmatrix} = \begin{bmatrix} 10.6 \\ 2.6 \\ 2 \\ 3.8 \\ 10 \\ 12.6 \\ 12.8 \end{bmatrix}$$

LDA classifier:

$$\textcircled{1} \quad \bar{x}_{new} = (5, 4).$$

$$z_{new} = (1 \times 5) + (0.8 \times 4) = (5 + 3.2) = 8.2.$$

$$\textcircled{2} \quad \cancel{\text{mean}}, \text{ class 0: } \cancel{\frac{[5 \ 7] + [2]}{2}} =$$

$$\text{mean, class 0: } \frac{2.6 + 2 + 3.8}{3} = 3.6$$

$$\text{class 1: } \frac{10 + 12.6 + 12.8}{3} = 11.8$$

\textcircled{8} Dist:

$$\textcircled{0} \quad |8.2 - 3.6| = 4.6$$

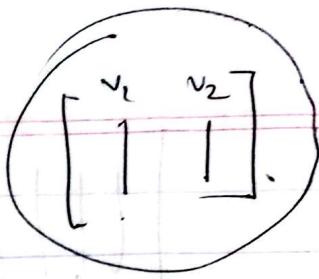
$$\textcircled{1} \quad |8.2 - 11.8| = 3.6$$

Class 1 closer

PCA :-

Steps →

- ① Standardize $\rightarrow S = \sqrt{\frac{\sum (x_i - \bar{x})^2}{N-1}}$
- ② Cov Matrix $\rightarrow \frac{1}{N-1} (X^T X)$
- ③ Eigenvalues, $\lambda \rightarrow |A - \lambda I| = 0$
Eigenvectors, $(A - \lambda_1 I) (v_1) = 0 \rightarrow v_1$
 $(A - \lambda_2 I) (v_2) = 0 \rightarrow v_2$



④ Sort Eigenvectors $\rightarrow \begin{bmatrix} v_1 & v_2 \\ \vdots & \vdots \end{bmatrix}$

⑤ Projection $\rightarrow X_{\text{normalized}} V = \begin{bmatrix} \cdot & \cdot & \cdot \end{bmatrix} \times \begin{bmatrix} \cdot & \cdot & \cdot \end{bmatrix} = \begin{bmatrix} PC_1 & PC_2 \\ \vdots & \vdots \end{bmatrix}$

$$\text{ans} \quad PC_1 = \frac{\sigma_{PC_1}^2}{(\sigma_{PC_1}^2 + \sigma_{PC_2}^2)} \times 100 \quad \checkmark$$

$$PC_2 = \text{same} - \quad \checkmark$$

① Standardized:

Given,	x_1	x_2
1	126	78
2	128	80
3	128	82
4	130	82
5	130	84
6	132	86
mean, $\bar{x} = \frac{\sum x}{N}$	129	82

$$N=1 \quad ; \quad \bar{x} = \frac{\sum x}{N}$$

Formula: $\sigma = \sqrt{\frac{\sum (x - \bar{x})^2}{N-1}}$

$x_1 \text{ center}$	$x_2 \text{ center}$
-3	-4
-1	-2
-1	0
1	0
1	2
3	4

$$\frac{\sum (x - \bar{x})^2}{(N-1)} =$$

$$2.1 \quad 2.8$$

② Correlation Matrix, A :-

$$\text{Correlation Matrix, } A = \frac{\sum X_{\text{centered}}^T X_{\text{centered}}}{(N-1)} = \frac{1}{N-1} \begin{bmatrix} -3 & -1 & -1 & 1 & 3 \\ -4 & -2 & 0 & 0 & 4 \end{bmatrix} \begin{bmatrix} -3 & -4 \\ -1 & -2 \\ -1 & 0 \\ 1 & 0 \\ 1 & 2 \\ 3 & 4 \end{bmatrix}$$

$$= \frac{1}{5} \begin{bmatrix} 22 & 28 \\ 28 & 40 \end{bmatrix} = \begin{bmatrix} 4.4 & 5.6 \\ 5.6 & 8 \end{bmatrix}$$

③ Eigen value, λ :- $|A - \lambda I| = 0$.

$$|A - \lambda I| = \begin{bmatrix} 4.4 & 5.6 \\ 5.6 & 8 \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} = \begin{bmatrix} 4.4 - \lambda & 5.6 \\ 5.6 & 8 - \lambda \end{bmatrix}$$

$$80, \begin{bmatrix} 4.4 - \lambda & 5.6 \\ 5.6 & 8 - \lambda \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\Rightarrow (4.4 - \lambda)(8 - \lambda) - (5.6)^2 = 0$$

$$\Rightarrow 35.2 - 12.4\lambda + 8\lambda^2 + \lambda^2 - 31.36 = 0$$

$$\Rightarrow \lambda^2 - 12.4\lambda + 3.84 = 0$$

$$\Rightarrow \lambda_1 = 12.08 \quad \lambda_2 = 0.32$$

Eigen vectors v_1, v_2 : $(A - \lambda I)v = 0$.

$$\text{For } \lambda_1, \begin{cases} A - \lambda_1 I = \begin{bmatrix} 4.4 - 12.08 & 5.6 \\ 5.6 & 8 - 12.08 \end{bmatrix} = \begin{bmatrix} -7.68 & 5.6 \\ 5.6 & -4.08 \end{bmatrix} \end{cases}$$

$$\text{So, } (A - \lambda_1 I) v_1 = 0$$

$$\text{or, } \begin{bmatrix} -7.68 & 5.6 \\ 5.6 & -4.08 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\text{Here, } -7.68v_1 + 5.6v_2 = 0 \Rightarrow v_1 = \frac{5.6}{7.68} v_2$$

$$5.6v_1 - 4.08v_2 = 0 \Rightarrow v_1 = \frac{4.08}{5.6} v_2$$

$$\Rightarrow v_2 = \frac{7.68}{5.6} v_1 \approx 1.37 v_1$$

$$\Rightarrow v_2 = \frac{4.08}{5.6} v_1 \approx 1.37 v_1$$

$$\text{For, } \lambda_2, (A - \lambda_2 I) v_2 = 0.$$

$$\begin{bmatrix} 4.08 & 5.6 \\ 5.6 & 7.68 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = 0$$

$$\text{Here, } 4.08v_1 + 5.6v_2 = 0$$

$$5.6v_1 + 7.68v_2 = 0$$

$$\text{So, } v_2 = -\frac{4.08}{5.6} v_1 \approx -0.72 v_1$$

$$v_2 = -\frac{5.6}{7.68} v_1 \approx -0.72 v_1$$

④ Sort eigenvectors by value.

$$\text{So, } v_1 = \begin{bmatrix} 0.59 \\ 0.807 \end{bmatrix}$$

$$v_2 = \begin{bmatrix} 0.81 \\ -0.59 \end{bmatrix}$$

$$\text{So, if } v_1 = 1; v_2 = 1.37. \text{ Thus } v_1 = \begin{bmatrix} 1 \\ 1.37 \end{bmatrix}$$

$$\text{normalized, } v_1 = \frac{1}{\sqrt{1^2 + 1.37^2}} \begin{bmatrix} 1 \\ 1.37 \end{bmatrix}$$

$$\text{or, } v_1 = \begin{bmatrix} 0.59 \\ 0.807 \end{bmatrix}$$

$$\text{So, if } v_1 = 1, v_2 = -0.72$$

$$\text{normalize, } v_2 = \begin{bmatrix} 0.81 \\ -0.59 \end{bmatrix}$$

together,

$$\begin{bmatrix} \lambda_1 & \lambda_2 \\ 0.59 & 0.81 \\ 0.807 & -0.59 \end{bmatrix}$$

$$\sigma^2 = \frac{\sum (x - \bar{x})^2}{N-1}$$

⑤ Projection: $X_{\text{centerd.}} V = \begin{bmatrix} -3 & -4 \\ -1 & -2 \\ -1 & 0 \\ 1 & 0 \\ 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 0.59 & 0.81 \\ 0.807 & -0.59 \end{bmatrix}$

$$= \begin{bmatrix} -5 & -0.07 \\ -2.2 & 0.4 \\ -0.6 & -0.81 \\ 0.6 & 0.81 \\ 2.2 & -0.37 \\ 5 & 0.07 \end{bmatrix}$$

$$\sigma^2_{PC1} = 12.4$$

$$\sigma^2_{PC2} = 0.26$$

PC1 contains $= \frac{12.4}{(0.26+12.4)} \times 100 = 97.87\% \text{ of variation}$

PCA $\mu = \frac{0.26}{(0.26+12.4)} \times 100 = 2.05\%$

Sof Hard Margin SVM

$$\text{minimize}_{w,b} : \frac{1}{2} \|w\|^2 = \frac{1}{2} w^T w$$

$$\text{Constraint} : y_i(w^T x_i + b) \geq 1$$

$$\text{Here, } g_i(w, b) = y_i - (y_i(w^T x_i + b) - 1) \leq 0$$

Lagrangian func:

$$L(w, b, \lambda) = \frac{1}{2} w^T w - \sum_{i=1}^N (y_i(w^T x_i + b) - 1) \lambda_i$$

So, according to KKT,

$$\frac{\partial L}{\partial w} = \frac{1}{2} w - \sum_{i=1}^N y_i x_i \lambda_i = w - \sum_{i=1}^N \lambda_i y_i x_i$$

$$\frac{\partial L}{\partial b} = - \sum_{i=1}^N \lambda_i y_i$$

$$\text{So, } w - \sum_{i=1}^N \lambda_i y_i x_i = 0$$

$$\Rightarrow w = \sum_{i=1}^N \lambda_i y_i x_i$$

$$\text{and } \sum_{i=1}^N \lambda_i y_i = 0$$

Substitute value:

$$L(w, b, \lambda) = \frac{1}{2} w^T w - \sum_{i=1}^N \lambda_i y_i w^T x_i - \sum_{i=1}^N \lambda_i y_i b + \sum_{i=1}^N \lambda_i$$

$$= \frac{1}{2} w^T w - \cancel{w^T w} + \sum_{i=1}^N \lambda_i$$

$$= -\frac{1}{2} w^T w + \sum_{i=1}^N \lambda_i$$

$$= \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_i x_j$$

Soft Margin

$$\text{Minimize}_{w, b} : \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

$$\text{Constraint: } y_i(w^T x_i + b) \geq 1 - \xi_i \quad \text{and} \quad \xi_i \geq 0.$$

$$\text{So, } L(w, b, \xi, \lambda, \mu) = \frac{1}{2} w^T w + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \mu_i \xi_i - \sum_{i=1}^n x_i (y_i(w^T x_i + b) - 1 + \xi_i)$$

$$= \frac{1}{2} w^T w + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \mu_i \xi_i - \sum_{i=1}^n x_i y_i w^T x_i - \sum_{i=1}^n \lambda_i y_i b + \sum_{i=1}^n x_i - \sum_{i=1}^n \xi_i x_i$$

Partial derivatives:

$$\frac{\partial L}{\partial w} = w - \sum_{i=1}^n x_i y_i x_i$$

$$\frac{\partial L}{\partial b} = - \sum_{i=1}^n x_i y_i$$

$$\frac{\partial L}{\partial \xi_i} = C - \sum_{i=1}^n \mu_i - \sum_{i=1}^n x_i = C - \mu_i - \lambda_i = 0.$$

$$\sum_{i=1}^n x_i [y_i(w^T x_i + b) - 1 + \xi_i] = 0.$$

$$\mu_i \xi_i = 0$$

$$\text{So, } \mu_i \geq 0 \text{ and } \lambda_i \geq 0.$$

$$\text{Here, } C - \mu_i - \lambda_i = 0.$$

$$\text{or, } \mu_i = C - \lambda_i \geq 0.$$

$$\text{If } \xi_i \geq 0, \mu_i = 0.$$

$$\text{or, } C - \lambda_i = 0.$$

$$\text{or, } C = \lambda_i$$

$$\text{If } \xi_i = 0; \mu_i \geq 0.$$

$$\text{or, } C - \lambda_i \geq 0.$$

$$\text{or, } C \geq \lambda_i$$

Eliminate ξ_i ,
we get Lagrangian fn
 $L(w, b, \lambda)$

তো এখন আরেও কৃত
কোটি কোটি পার্শ্ব পরিস্থিতি

* Kernel Method :



To avoid HD projecting data directly to HD space - compute dot prod

$$\phi: x \rightarrow \phi(x)$$

$$K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$$

Find w

We know,

$$\max_{\lambda} \left[\sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j K(x_i, x_j) \right], \quad w = \sum_{i=1}^n \lambda_i y_i x_i$$

Find w:

$$g(x_{\text{new}}) = \left(\sum_{i=1}^n \lambda_i y_i \phi(x_i) \right) \phi(x_{\text{new}}) + b$$

$$w = \sum_{i=1}^n \lambda_i y_i k(x_i, x_{\text{new}}) + b$$

For b,

$$b = y_i - w^T x_i$$

$$= y_i - \sum_{i=1}^n \lambda_i y_i k(x_i, x_{\text{new}})$$

④ Euclidean distance $\rightarrow A = (3, 4) \quad B = (0, 0)$. in 2D space
 $\sqrt{(3-0)^2 + (4-0)^2} = \sqrt{9+16} = 5$

⑤ Normalized Euclidean distance :-

ব্যবহার করা মূল্য ২। একে $d(x,y) = \frac{\sum (x-y)^2}{(\max-\min)^2} \cdot \sigma^2$

$$A = (2, 1000) \quad B = (1, 800)$$

$$\text{সু, } d(A, B) = \sqrt{\frac{(1-2)^2}{(1-2)^2} + \frac{(1000-800)^2}{200^2}} = \checkmark$$

⑥ Mahalanobis Distance : dist btw point \rightarrow distribution - taking correlation account

generalized Euclidean dist

$$d_M(x, \mu) = \sqrt{(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

Euclidean	28	Mahalanobis
① dist btw 2 points in 2D space		① generalized euclidean dist btw point and its distribution
② Doesn't consider correlation		② Considers correlation
③ Treats feature independent		③ accounts feature correlation
④ assumes spherical cluster.		④ adapts to elliptical cluster

Both criterion for to solve linear inequalities $a^T y > 0$

Batch Perceptron

(vs)

Fixed Increment Single Sample per

1. Computes cumulative correction
(sum of misclassified samples)

2. updates after cumu... correcn

3. Computationally expensive

4. All samples are checked in each iteration

5. less noisy

6. Converges earlier than fixed increment.

1. Computes misclassified samples iteratively and updates them if found.

2. updates when misclassified sample found

3. cheaper

4. Each sample is handled in each iteration and keeps iterating if misclassified

5. More noisy

6. Convergence needs more iteration as each step is small

1. Initialize $\leftarrow a(0), \eta \leftarrow 1$
2. while (True) :
3. mis $\leftarrow 0$,
4. for all sample y in dataset
5. if label($a^T y$) ≤ 0
6. mis \leftarrow mis + 1
7. store (mis) in np.array
8. endif
9. end for
10. $a(\text{new}) \leftarrow a(\text{old}) + \sum \text{label.sample}$
11. if mis == 0 .
12. break;
13. endif
14. end while

1. Initialize $\leftarrow a(0), \eta \leftarrow 1$.
2. while (True) :
3. mis $\leftarrow 0$
4. for all sample y in D.
5. if label($a^T y$) ≤ 0
6. mis \leftarrow mis + 1
7. $a(\text{new}) \leftarrow a(\text{old}) + \eta \cdot \text{label}$.
8. end if
9. end for
10. if mis == 0
11. break
12. endif
13. end while

Calculation of ω in kernel:-

$g(\text{new}) = \dots$

$$g(x_{\text{new}}) = \left(\sum_{i=1}^n \lambda_i y_i \Phi(x_i) \right) \Phi(x_{\text{new}}) + b.$$

$$= \sum_{i=1}^n \lambda_i y_i K(x_i, x_{\text{new}}) + b.$$

Dual form:-

$$\max_{\lambda} \left[\sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j K(x_i, x_j) \right]$$

$$\omega = \sum \lambda_i y_i \Phi(x_i)$$

$$b = y_i - \omega^T \Phi(x_i).$$

for b ,

$$\cancel{y_i \left(\sum_{i=1}^n \lambda_i y_i \Phi(x_i) + b \right) = 1}$$

$$\Rightarrow \cancel{\sum_{i=1}^n \lambda_i y_i K(x_i, x_i) + b} = y_i$$

$$\Rightarrow \cancel{b} = y_i - \sum_{i=1}^n \lambda_i y_i K(x_i, x_i).$$

Hard if $g(x_{\text{new}}) > 0$

- ① assumes linearly separable data
- ② No misclassification allowed
- ③ widest-possible margin gap searches.
- ④ Rigid, strict
- ⑤ fails if noise/overlapping found as can't tolerate mistakes.



Soft

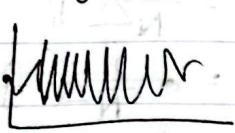
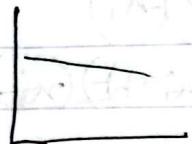
- ① messy data
- ② might face some.
- ③ adds penalty for mis points (C) keeping large margin regu.
- ④ flexible, balances.
- ⑤ gives penalties, tries to regularize using C , balances.



LR :

small

Larger (variance)



wild fluctuate



drastic increase

↓
noisy

Stochastic GD \rightarrow batch size = 1 per itro. (random chose example).

min batch SGD \rightarrow ~~1 <~~ 1 < batch size $< N$.

\hookrightarrow no. of sample

small batch \rightarrow SGD -> 驚

along " \rightarrow full GD

* Full batch \rightarrow 1000 sample, 20 epoch, 20 updates $\frac{\text{once}}{\text{per epoch}}$.

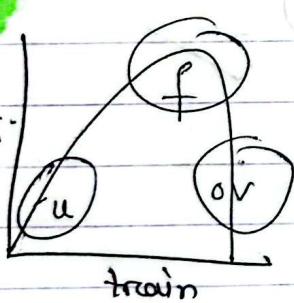
* SGD \rightarrow " " " 20,000 " w, b.

* minibatch SGD \rightarrow " " " ", 100 batch size,

$$\text{batches per epoch} = \frac{1000}{100} = 10.$$

$$\text{thus update} \cdot 10 \times 20 = 200.$$

* Regularization - 20, rule's avoid overfitting

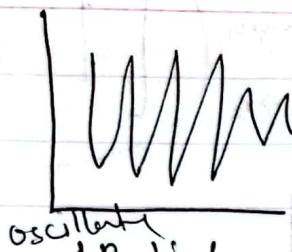


$$\min(\text{loss } f + \lambda(\text{regu}))$$

$$\text{regu (log loss} = \text{log loss} + \lambda \Omega(w))$$

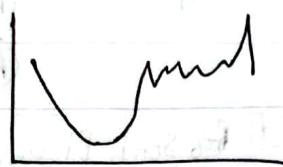
λ small \rightarrow penalty minimal
fit train data — large w.
overfits.
w is capturing
complex mod

large $\lambda \rightarrow$.
penalty — significantly
reduces w, closer to 0
don't fit noise
reduces complexity
oversimplifies
underfit



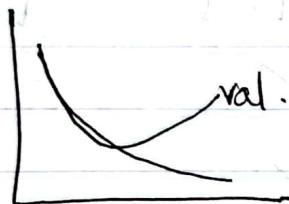
oscillate
LR high

→ clip grad, optimize



sharp rise.

LR too high, grad exploded
data noisy → model unstable.



sharp val loss rise

Overfit, data leak

concept drift



chaotic repeated loss

LR high.

overfit (no dropout / no regu).

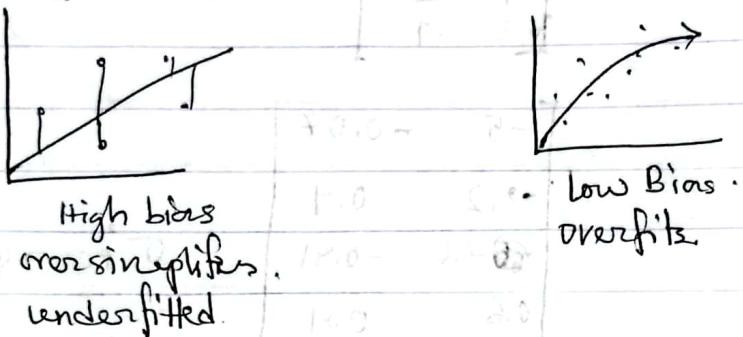
numeric instability.

unstable batch norm

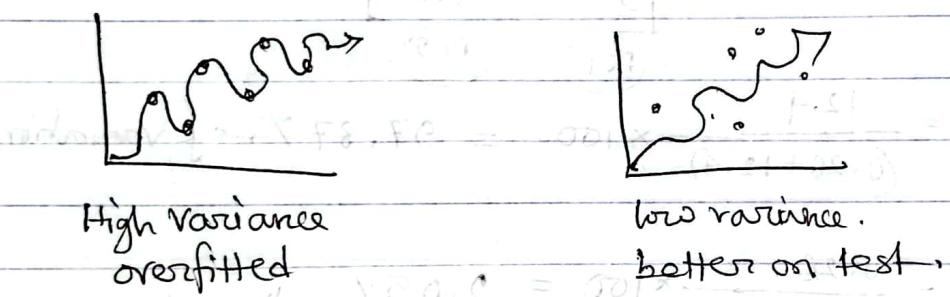
14

- * Algo do not outperform each other always.
- * compare - error rate on test data instead train.

Bias: diff btw predictions and corr. value



Variance: measures how much pred chgs if on diff dataset

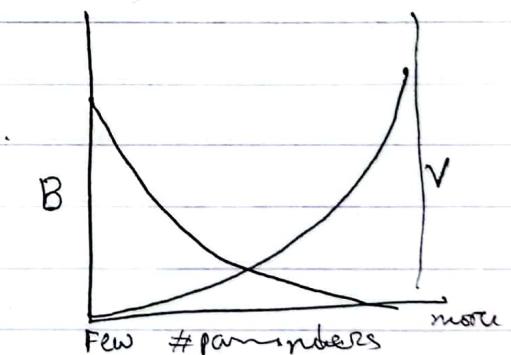


Bias-Variance tradeoff :-

HB - LV : model simple.
underfitted; both sets - loss high.

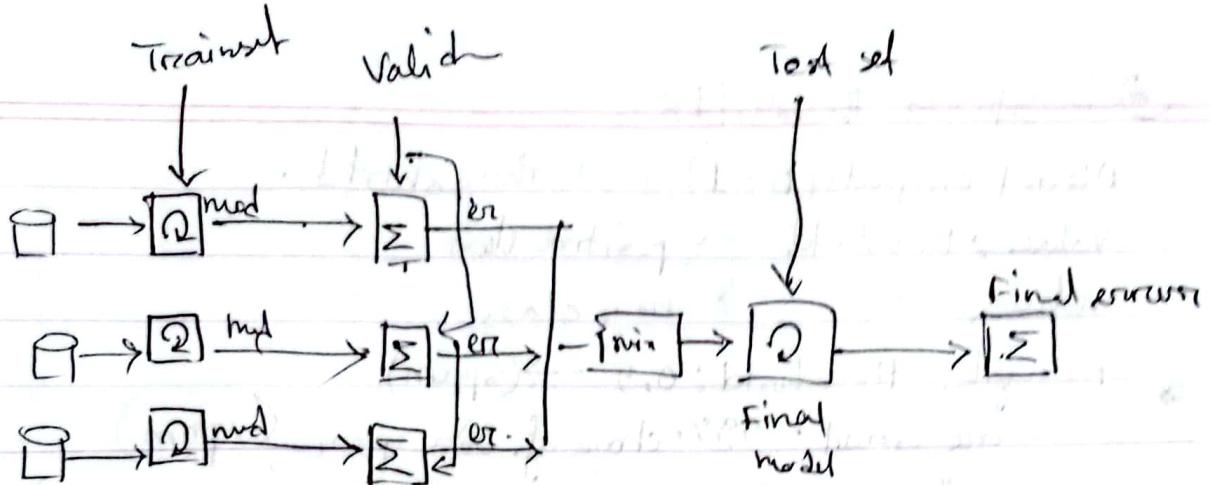
LB - HV : model complex.
overfitted

Balance B - V : moderate complexity.



Cross validation:

- train 1. Train set : optimizes parameters (w, b)
- evaluate 2. Validation set : \Rightarrow hyperpare (no. of NN, length of train)
- assess 3. Testing " : final assessment

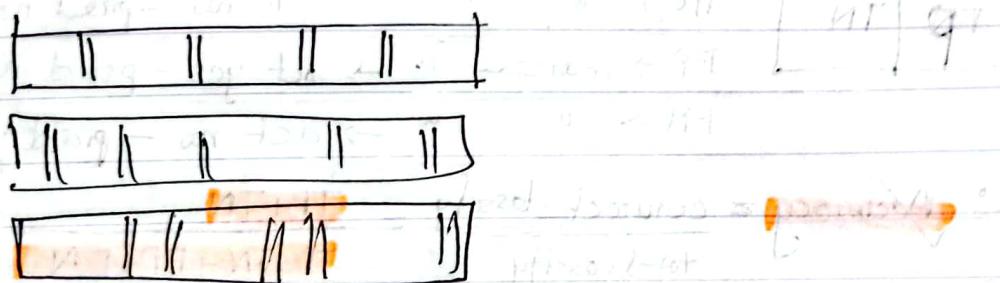


CV Methods:

① Holdout : simpler, test-train
learns \downarrow \rightarrow performs on test set.

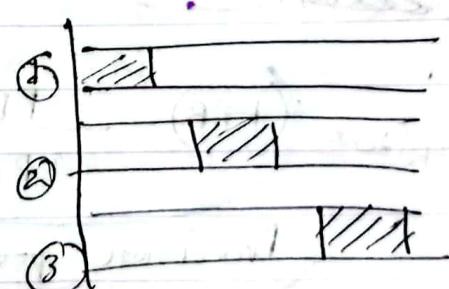
② Random Subsampling :

- * repeats holdout randomly selecting example.
- * error rate from test sample



③ k-Fold CV :

- * k equal sized subsets .
- * one subset - test
- * rest n - train
- * repeated k times on each subset .
- * Total errors = Avg of all errors .
- * like random subsampling but .
all examples are used for both train test



④ Classify & Thresholds

Model outputs (0 - 1), set threshold.

value > threshold → positive class.

value < " → neg class.

Example: threshold: 0.5 → (spam)

one email = 99% chance of being spam (spam)

another " = 0.44% " " " " (not spam)

third " = 0.51% " " " " (\$spam)

* Confusion Matrix: summarizes performance comparing actual vs predicted labels

		A	
		P	N
P	P	TP	FN
	N	FP	TN

TP → correct prediction → actually yes - pred yes

TN → " " " " → " no - pred no.

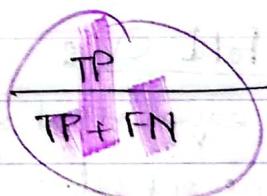
FP → incorrect " → act yes - pred no.

FN → " " " " → act no - pred yes

$$\text{Evaluation: Accuracy} = \frac{\text{correct classify}}{\text{total classify}} = \frac{TP+TN}{TP+TN+FP+FN}$$

perfect model - 90%, FP=FN=0.

$$\text{Recall (TPR)} = \frac{\text{correct - actual positive}}{\text{all actual pos}}$$



$$\text{FPR (FNR)} = \frac{FP}{FP+TN}$$

$$\text{Precision} = \frac{\text{positive correct positives}}{\text{total classified positive}} = \frac{TP}{TP+FP}$$

$$\text{F1 score} = \frac{\text{precision} + \text{recall}}{2}$$

$$= \frac{2TP}{2TP+FN+FP}$$