

28th - March →

Classification

Bayesian Method →

2 type of methods →

- ①. Parametric → To estimate parameters of given distribution
- ②. Non-

- ①. Parametric Method → Data follows some known distribution.

Duda-Hart → chapter (1) & (2).

$x_1 \rightarrow$ length $\not\rightarrow$ 2 features.
 $x_2 \rightarrow$ lightness.

$y \rightarrow$ classes $\rightarrow \{w_1, w_2\}$.

$\rightarrow P(x/w_i)$ → } Probability-Density function for R.I. x^i
 class-conditional of lightness for a particular value x^i given class w_i .
conditional
Densities

$\rightarrow P(w_i/x) \rightarrow$ Given the feature x^i , what is the probability that it belongs to w_i class.

$P(w_2/x) \rightarrow$ "

$P(w_1) \rightarrow$ Probability of fish being salmon (w_1)

$P(w_2) \rightarrow$ " " " " " sea-bass (w_2).

Prior probabilities

So, Given a feature 'n' →

$$\text{Bayesian formula} \rightarrow P(w_j/n) = \frac{P(n/w_j) \cdot P(w_j)}{P(n)}$$

$$\underline{P \text{ for 2 categories} \rightarrow P_{(M)} = \sum_{j=1}^2 P(M/w_j) \cdot P(w_j)}$$

Bayesian formula can be informally expressed in English as →

$$\text{Posterior} \rightarrow (\text{Likelihood} \times \text{Prior})$$

Decision Criteria → whichever is greater.

$$P(w_1/n) \geq P(w_2/n) \rightarrow \text{class } w_1$$

$$P(\text{error}/n) = \begin{cases} P(w_1/n) & ; \text{ if we decide } w_2 \\ P(w_2/n) & ; w_3 \end{cases}$$

for multiclass classification →

if $P(\omega_{j/n}) \geq P(\omega_i/n)$: $\forall i = 1, 2, \dots, k$
 $i \neq j$

\Rightarrow Average probability of error \rightarrow

$$P(\text{error}) = \int_{-\infty}^{+\infty} P(\text{error}, m) dm = \int_{-\infty}^{+\infty} P(\text{error}/m) \cdot P(m) dm$$

$$\text{Thus } \rightarrow P(\text{error}/n) = \min [P(w_1/n), P(w_2/n)]$$

8 $P(\omega_1/x) + P(\omega_2/x) = 1$: {for binary classification}

③ Decide ω_1 , if $\rightarrow P(x/\omega_1) \cdot P(\omega_1) > P(x/\omega_2) \cdot P(\omega_2)$:

Decision based on conditional & prior probability. otherwise decide ω_2
(evidence is same).

→ Basically Bayesian Decision rule combines conditional & prior probability to achieve minimum probability of error.

⇒ We can extend it for 'c' classes.

~~Q4 H.R.P.~~ \Rightarrow The results that we have seen in case of bayesian approach can be biased. Therefore, we need to attach another factor

e.g. $L_1 \rightarrow w_1$ } of (e.g. $w_1 = \text{gold}$, $w_2 = \text{coal}$) in that
 $L_2 \rightarrow w_2$
 $L_j \rightarrow w_j$ } $\rightarrow (x_i)$ is the action to classify observation (x) to class (w_i).
Loss w.r.t. actions taken

①. $\lambda(L_1/w_1) = 0 \rightarrow$ if accurately classified.

②. $\lambda(L_1/w_2) = 1 \rightarrow$ ~~if wrongly classified~~. loss regarding vice-versa of ①, very high loss. { If stone was coal

③. $\lambda(L_2/w_1) = K \rightarrow$ vice-versa of ②, very high loss. { It is thrown in the bucket of diamond

Conditional ④. $\lambda(L_2/w_2) = 0 \rightarrow$ accurate.

Risk

$$R(L_2/x) = \sum \lambda(L_2/w_j) P(w_j/x)$$

basically \rightarrow all shifts. λ : loss values. P : probability of assignment.

\Rightarrow And we'll make the decision with min. risk.

$$\text{overall Risk} \Rightarrow R(\alpha_i/n) = \sum_{j=1}^c \lambda(\alpha_i/\omega_j) \cdot P(\omega_j/x)$$

$\frac{f}{\text{overall}} \rightarrow \text{then } R = \int R(\alpha_i/n) \cdot P_{\text{out}} \cdot dx$: {where date is our notation
for a d-space volume element-
where the integral extends
over the entire
feature space}

Risk. → Suppose we are given 2 classes →

2-category classification

$$R(\alpha_1/n) = \lambda_{11} \cdot P(\omega_1/n) + \lambda_{12} \cdot P(\omega_2/n)$$

$$\lambda(\alpha_1/\omega_1)$$

$$\lambda(\alpha_1/\omega_2)$$

$$\text{Similarly, } R(\alpha_2/n) = \lambda_{21} \cdot P(\omega_1/n) + \lambda_{22} \cdot P(\omega_2/n)$$

$$\lambda(\alpha_2/\omega_1)$$

$$\lambda(\alpha_2/\omega_2)$$

$$\text{So, if } R(\alpha_1/n) \leq R(\alpha_2/n)$$

this will be selected

$$\Rightarrow \lambda_{11} \cdot P(\omega_1/n) + \lambda_{12} \cdot P(\omega_2/n) \leq \lambda_{21} \cdot P(\omega_1/n) +$$

$$\lambda_{22} \cdot P(\omega_2/n)$$

this will form the decision-criteria.

04th April

Decision Criteria → Decide w_1 if $R(\alpha_1/x) < R(\alpha_2/x)$,

otherwise w_2 , so if $R(\alpha_1/x) < R(\alpha_2/x) \Rightarrow w_1$

To decide w_1 → $(\lambda_{\alpha_1} - \lambda_{11}) \cdot P(w_1/x) > (\lambda_{12} - \lambda_{\alpha_2}) \cdot P(w_2/x)$

$$[(\lambda_{\alpha_1} - \lambda_{11}) \cdot P(x/w_1) \cdot P(w_1)] > [(\lambda_{12} - \lambda_{\alpha_2}) \cdot P(x/w_2) \cdot P(w_2)]$$

positive factor

Also, under assumption that $\lambda_{\alpha_1} > \lambda_{11}$ → we can say that:

* Likelihood Ratio: $\frac{P(x/w_1)}{P(x/w_2)} > \frac{(\lambda_{12} - \lambda_{\alpha_2}) \cdot P(w_2)}{(\lambda_{\alpha_1} - \lambda_{11}) \cdot P(w_1)}$;

& we can consider $P(x/w_j)$ as a f_j of w_j , which is likelihood f_j .

Minimum error-rate classification →

if action is α_i & true state of nature is w_j ,
then provided, the cond.^m. that we are using
symmetrical (or) 0/1 loss fn.

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

So → Conditional Risk →

$$R(\alpha_i/x) = \sum_{j=1}^c \lambda(\alpha_i/w_j) \cdot P(w_j/x)$$

$$\Rightarrow \sum_{j \neq i} P(w_j/x) ; \quad \text{in case of 0/1 loss}$$

$$\Rightarrow \{1 - P(w_i/x)\}$$

using other way → we can maximise the posterior prob. for decision criteria.

Select (w_i) if $P(w_i/x) > P(w_j/x) \therefore$ for all $j \neq i$

$$\frac{P(x/w_i)}{P(x/w_j)}$$

θ_2

θ_1

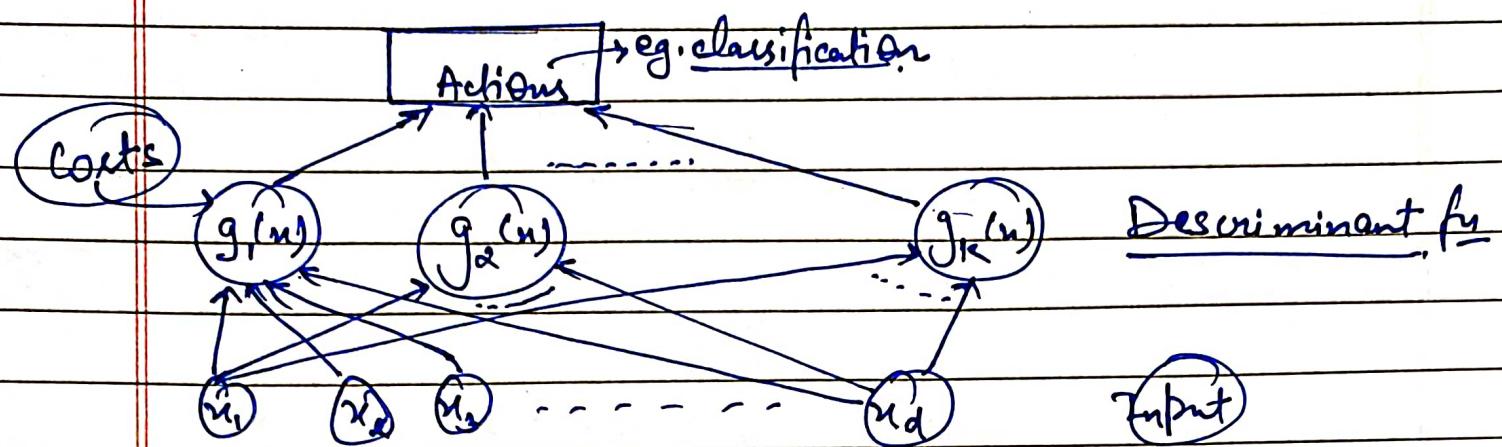
→ θ_0 is threshold, if loss fn penalizes

misclassifying w_2 as w_1 , patterns more than the
universe (i.e. $\lambda_{12} > \lambda_{21}$),

$\leftarrow R_2 \rightarrow \leftarrow R_1 \rightarrow \leftarrow R_2 \rightarrow R_1 \rightarrow x$ we get larger
threshold θ_0 and

hence Range of x values for which we classify
a pattern as ' w_i ' gets smaller.

Classifiers & Discriminant functions →



classifier is said to assign a feature vector ' x ' to class ' w_i '. if →

$$\boxed{g_i(x) > g_j(x) \text{ for all } j \neq i}$$

→ Minimum error rate classification

$$R(\alpha_i/n) = 1 - P(w_i/n)$$

→ all the model displayed in slide comes under umbrella of generating models.

⇒ $g_1(u), g_2(u), \dots, \dots, g_K(u)$

$$\boxed{\arg \max_k (g_k(u))}$$

most probably of 2 class

Date _____

Page No. _____

$$\text{So } \rightarrow g_j(x) = -R(\alpha_j/x)$$

$\Rightarrow -(1 - P(w_j/x))$: { check on what cond. this is true }.

So if $g_i(x) > g_j(x)$: if $i \neq j$
 $\& i, j = 1, 2, \dots, k$.

Then we'll select class 'i'.

$$g_i(x) = P(w_i/x) = \frac{P(x/w_i) \cdot P(w_i)}{P(x)}$$

$$g_j(x) = \frac{P(x/w_j) \cdot P(w_j)}{P(x)}$$

ignore this.

→ Take log

$g_i(x)$

$$\log g_i(x) = \log P(x/w_i) + \log(P(w_i)) \quad \text{--- (1)}$$

$g_j(x)$

$$\log g_j(x) = \log P(x/w_j) + \log(P(w_j))$$

Now, how to calculate $P(x/w_i)$ & $P(x/w_j)$

Decision Boundaries are formed.

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

so, if $g(x) > 0 \Rightarrow w_1$
otherwise $\Rightarrow w_2$

Decision Boundary $\rightarrow g(u) = g_1(u) - g_2(u)$

So $\rightarrow g(u) = P(w_1/u) - P(w_2/u)$

$$g(u) = \ln \left[\frac{P(x/w_1)}{P(u/w_2)} \right] + \ln \left[\frac{P(w_1)}{P(w_2)} \right]$$

→ Discriminant fct for Normal Density \rightarrow

So \rightarrow if $x \in \mathbb{R}^d$

& $P(x/w_i) \sim N(\mu_i, \Sigma_i)$

So \rightarrow

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

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

if we have 'c' classes \rightarrow

such that:-

$$w_1 \rightarrow \mu_1, \Sigma_1$$

$$w_2 \rightarrow \mu_2, \Sigma_2$$

⋮

$$w_c \rightarrow \mu_c, \Sigma_c$$

There are multiple cases, that can occur \rightarrow

Case(1) \rightarrow $\sum_i = \sigma^2 \cdot I$ \rightarrow when features are statistically independent. each feature has same variance σ^2 & co-variance of α features are 0.

in 2-D: $x \in \mathbb{R}^2$. so \rightarrow

$$\sum_i = \begin{bmatrix} \sigma_i^2 & 0 \\ 0 & \sigma_i^2 \end{bmatrix}$$

Note \rightarrow ①. if variance in X & Y directions are same & co-variance is 0, then contours will be co-centric circles.

②. if co-variance is 0 but variance in both directions are different, then contours will be ellipse

③. if co-variance is also (non), the contours will be tilted ellipse.

See book for all 3-cases \rightarrow

Duda-Hart page

chapter 2 - 2.6

Q \rightarrow

Convert the discriminant f_2 into a linearly Discriminant f_1 .

Date _____
Page No. _____

So, Discriminant f_1 for case(i) is \rightarrow

$$g_i(u) = -\frac{1}{2} \frac{(u - \mu_i)^T}{\sigma^2 I} - \frac{1}{2} \ln |\sigma^2 I| + \ln P(\omega_i)$$

↓

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

(a) Similarly. \rightarrow

$$g_2(m) = \frac{-1}{\alpha} (n - \mu_2) - \frac{1}{2} \cdot \ln |\sigma^2| + \ln P(\omega_2)$$

~~$\text{Assignment} \rightarrow \Sigma_i = \Sigma_\alpha = \begin{bmatrix} 0.6 & 0 \\ 0 & 0.6 \end{bmatrix}$~~

$$g_1(n) = \log P(\omega_1 | n) = \frac{P(n | \omega_1) * P(\omega_1)}{P(n)}$$

$$\text{so } \log \bar{g}_1(n) = \log P(n | \omega_1) + \log P(\omega_1)$$

$$\text{Similarly } \log \bar{g}_2(n) = \log P(n | \omega_2) + \log P(\omega_2)$$

$$| g(n) = g_1(n) - g_2(n) |$$

if $g(n) \geq 0 \Rightarrow g_1(n)$
else $\Rightarrow g_2(n)$

$$-(x - \mu_1)^T \cdot \Sigma^{-1} (n - \mu_1)$$

$$g_1(n) = \log \left(\frac{1}{(2\pi)^{\frac{D}{2}} \cdot \sum^2} \right) \exp$$

$$+ \log P(\omega_1)$$

\Rightarrow expand

Similarly for $g_2(n) \Rightarrow$

$$\Sigma = \begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix}$$

$$\Sigma^{-1} = \frac{1}{\sigma^2}$$

$$g_1(n) = -\frac{1}{2\sigma^2} \cdot (x^T x - 2\mu_1^T x + \mu_1^T \mu_1) + \log P(\omega_1)$$

$$g_2(n) = -\frac{1}{2\sigma^2} (n^T n - 2\mu_2^T n + \mu_2^T \mu_2) + \log P(\omega_2)$$

$$\text{So } \rightarrow g(n) = g_1(n) - g_2(n) = 0$$

$$g(n) = \left(\frac{1}{\sigma^2} \mu_1 \right)^T x - \frac{1}{2\sigma^2} \mu_1^T \mu_1 + \log P(\omega_1)$$

$$\text{So } \rightarrow$$

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

then, we'll get sol. \rightarrow

$$g(n) \Rightarrow \boxed{\omega^T (n - n_0) = 0} \quad \text{and } \omega = (\mu_i - \mu_j)$$

$$\& n_0 = \frac{1}{2}(\mu_i + \mu_j) - \frac{\sigma^2}{\|(\mu_i - \mu_j)\|^2} \ln \frac{P(\omega_2)}{P(\omega_1)} (\mu_i - \mu_j)$$

example from assign. 05 \rightarrow

$$g_1(n) = -\frac{\|n - \mu_1\|_2^2 + \log P(\omega_1)}{2\sigma^2}$$

$$g_2(n) = -\frac{\|n - \mu_2\|_2^2 + \log P(\omega_2)}{2\sigma^2}$$

Decision separating plane

Mahalanobis Distance

$$x_0 = \frac{1}{2}(\mu_1 + \mu_2) - \frac{\sigma^2}{\|(\mu_1 - \mu_2)\|^2} \ln \frac{P(\omega_1)}{P(\omega_2)} (\mu_1 - \mu_2)$$

$$\sum_i = \sum$$

variance is same in both the dir. but co-variance exists.

Date	_____
Page No.	_____

Case ② $\rightarrow g_L(u) = w^T x + w_0$

$$w_0 = \sum^{-1} u_1$$

$$w_0 = -\frac{1}{2} u_1^T \sum^{-1} u_1 + \ln P(w_1)$$

$$g(u) = g_L(u) - g_R(u) = 0$$

$$\Rightarrow w^T(x - x_0) : \{ \text{Decision criteria}\}$$

$$w = \sum^{-1}(u_1 - u_2)$$

$$x_0 = \frac{1}{2}(u_1 + u_2) - \frac{\ln [P(w_1)/P(w_2)]}{(u_1 - u_2)^T \sum^{-1} (u_1 - u_2)}$$

(missing)

This min.

Complete

Case ③ \rightarrow arbitrary \rightarrow

$$g_L(u) = w_1^T u_1 + w_0$$

$$w_1 = -\frac{1}{2} \cdot \sum^{-1}$$

$$w_0 = \sum^{-1} u_1$$

$$\text{and} \rightarrow w_{10} = -\frac{1}{2} u_1^T \sum^{-1} u_1 - \frac{1}{2} \ln |\sum| + \ln P(w_1)$$

Similarly $\rightarrow g_R(u) = -$

$$\rightarrow g(u) = g_L(u) - g_R(u) = 0$$

Ques. → What is generative model?

Ques. → " " " Discriminant fn./model?

→ Regression Methods for binary classification →

we'll try to contain the regression values b/w (0 & 1). So, that the o/p. values denote the probability values.

$$x \xrightarrow{\theta} y$$

$$\cancel{P(Y=1/x)} P(Y=1/x)$$

④

$$\begin{cases} \beta_1 x + \beta_0 \\ w^T x + b \end{cases}$$

it is not going to work

$$\log \frac{P(Y=1/x)}{1-P(Y=1/x)}$$

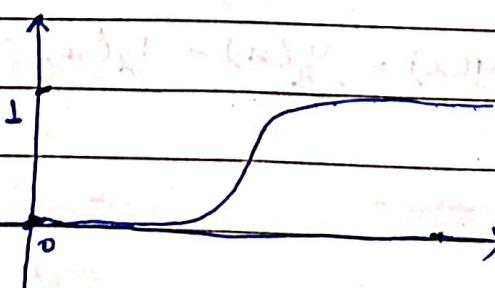
$$f \rightarrow \log \frac{P(Y=1/x)}{1-P(Y=1/x)}$$

our estimate is that, this is linear

$$\log \left(\frac{P(x)}{1-P(x)} \right) = \beta^T x, \quad x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, \quad \beta = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_n \end{bmatrix}$$

$$\Rightarrow \text{it is a } \beta^T x + \beta_0 = \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \beta_0$$

$$\begin{aligned} \text{Linear fn. of } x \\ \frac{P(x)}{1-P(x)} &= e^{\beta^T x} \\ P(x) &= \frac{e^{\beta^T x}}{1+e^{\beta^T x}} = \frac{1}{1+e^{-\beta^T x}} \end{aligned}$$



$$\sum_{i=1}^L \left(f(x_i) - \left(\frac{e^{\beta^T x_i}}{1 + e^{\beta^T x_i}} \right) \right)^2$$

\downarrow \downarrow \downarrow
 x_i \sim $f(x_i)$
 actual estimated

If (Y/X) follows Normal distribution, then least square loss function is optimal.

$$(x_i, y_i) \in \mathbb{R}^n \times \{0, 1\}$$

$$\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

$$L(\beta) = \prod_{i=1}^n p(x_i) \prod_{i=1}^n 1 - p(x_i)$$

$\begin{cases} \prod \rightarrow \text{product} \\ \sum \rightarrow \text{summation} \end{cases}$

for a given value of X , Y can have 2 values, either '0' or '1', so it follows Bernoulli distribution. If B is most of the data pts. are around '0' & '1'. y.

$$d(\beta) = \prod_{i=1}^m p(x_i)^{y_i} (1-p(x_i))^{1-y_i}$$

$$\text{Max } \log L(\beta) = \sum_{i=1}^n y_i \log P(X_i) + (1 - y_i) \log (1 - P(X_i))$$

$$= \sum_{i=1}^n y_i \log\left(\frac{P(X_i)}{1 - P(X_i)}\right) + \log(1 - P(X_i))$$

$$= \sum_{i=1}^n y_i (\beta^T x_i) - \log (1 + e^{\beta^T x_i})$$

$$\underset{\beta}{\operatorname{Min}} \quad - \sum_{i=1}^n y_i (\beta^T x_i) - \log (1 + e^{\beta^T x_i})$$

$$\nabla_{\beta} f = - \sum_{p=1}^n \left(y_p x_p - \left(\frac{e^{\beta^T x_p}}{1 + e^{\beta^T x_p}} \right) x_p \right) = 0$$

$$= - \sum_{p=1}^n \left(y_p - \frac{e^{\beta^T x_p}}{1 + e^{\beta^T x_p}} \right)$$

Choose β^0

Do until

$$\|\nabla f(\beta)\| \leq \epsilon$$

$$\beta^{t+1} = \beta^t + \eta_p \sum_{p=1}^n \left(x_p - \frac{e^{x_p \beta^T}}{1 + e^{x_p \beta^T}} \right) x_p$$

end

$$\eta_p = 0.01, 0.001$$

$$\eta_p = \frac{1}{1+p}$$

$$\eta_p = \frac{1}{p^2}$$

$$\hat{P}(x) = \frac{e^{\beta^T x}}{1 + e^{\beta^T x}}$$

$$f(x) = \begin{cases} 0 & \hat{P}(x) \leq 0.5 \\ 1 & \text{otherwise} \end{cases}$$

$$\text{logit } \hat{P}(x) = \beta^T x = \log \left(\frac{\hat{P}(x)}{1 - \hat{P}(x)} \right)$$

$$f(x) = \begin{cases} 0 & \text{if } \beta^T x \leq 0 \\ 1 & \text{otherwise} \end{cases}$$

$$\beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \beta_0 \leq 0$$

$$\beta^T x = 0$$

$$P(Y=1/X) = \frac{e^{\beta_0^T x}}{1 + \sum_{j=1}^{k-1} e^{\beta_j^T x}}$$

$$P(Y=2/X) = \frac{e^{\beta_2^T x}}{1 + \sum_{j=1}^{k-1} e^{\beta_j^T x}}$$

$$P(Y=k/X) = \frac{e^{\beta_k^T x}}{1 + \sum_{j=1}^{k-1} e^{\beta_j^T x}}$$

Weighted
regression technique

$$\text{Min } f_0(x)$$

$$\text{subject to } f_i(x) \leq 0, i=1, 2, \dots, m$$

$$h_i(x) = \alpha_i^T x - b_i = 0, i=1, 2, \dots, p$$

$$\text{Min } 3x_1^2 + 4x_2^2$$

such that

$$4x_1 + 5x_2 \leq 5$$

$$6x_1 + 6x_2 = 4$$

$$x_1 \geq 0, x_2 \geq 0$$

$$p^* = \inf \{f_0(x) \mid x \in D\}$$

where, $D = \{x \mid f_i(x) < 0$

$$i=1, 2, \dots, m$$

$$h_i(x) = 0$$

$, i = 1, 2, \dots$

$$f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{j=1}^p \mu_j$$

For $x \in D$, $\lambda \geq 0$

$$\underline{L}(x, \lambda, \mu) \leq f_0(x)$$

$$\inf_{\lambda \in \mathbb{R}^m} \underline{L}(x, \lambda, \mu) \leq \inf_{\lambda \in D} \underline{L}(x, \lambda, \mu)$$

$$\leq \inf_{x \in D} f_0(x) = p^*$$

$$g(\lambda, \mu) = \inf_{x \in D} f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{j=1}^p \mu_j h_i(x)$$

subject

$$\lambda \geq 0$$

$$g(\lambda, \mu) \leq p^*$$

→ Binary classification through S.V.M. ⇒

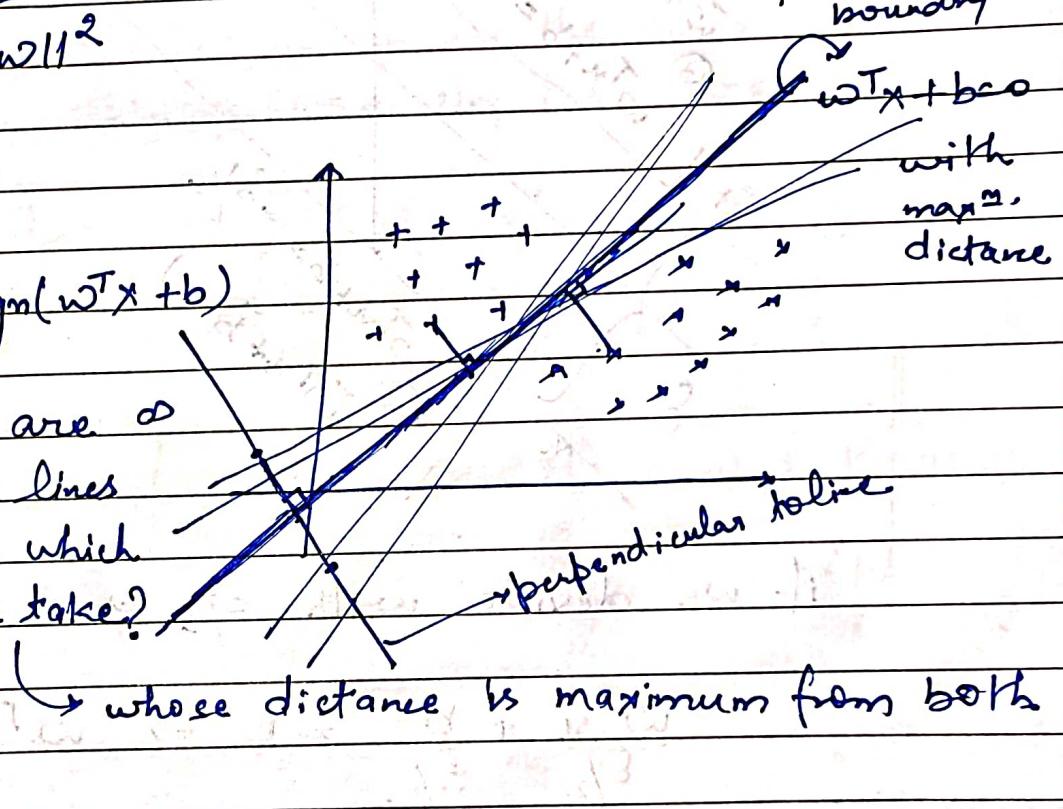
$$T = \{(x_1, y_1), (x_2, y_2), \dots, (x_L, y_L)\}; x_i \in \mathbb{R}^n \\ y_i \in \{-1, 1\} \\ i=1, 2, \dots, L.$$

$$\text{Min } \frac{\alpha}{2} \|w\|_p + C \sum_{i=1}^L L(y_i, x_i, w, b)$$

we are mostly minimizing $\frac{1}{2} \|w\|^2$

$$f(w) = \text{Sign}(w^T x + b)$$

But there are ∞ possibility of lines in this region. which one should be take?



⇒ Take projections of pts. on the perpendicular line

$$w^T x + b$$

$$x_+ = \text{Min}_{x_+ \in I} (w^T x_+ + b),$$

$$x_- \in I$$

So, we need to find a dim. in which \rightarrow margin.

$$\frac{\text{Max.}}{w, b} (w^T x_+ + b) - (w^T x_- + b) \geq \frac{1}{\|w\|_2}$$

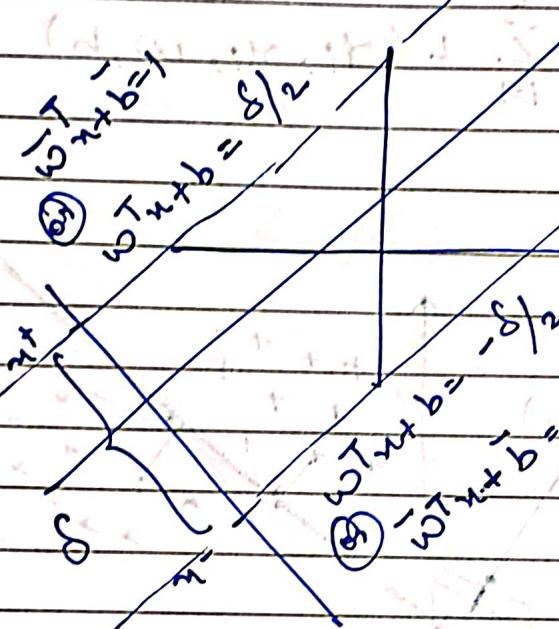
$$\Rightarrow \text{Max} \left\{ \frac{\delta}{\|w\|_2} \right\} = \text{Max} \left\{ \frac{1}{\|w\|_2} \right\}$$

margin.

$$\Rightarrow \text{Min} \left\{ \frac{1}{2} \|w\|_2^2 \right\}$$

$$\Rightarrow \text{Min} \left\{ \frac{1}{2} \|w\|_2^2 \right\}$$

$\hat{w}^T x + b = 0$

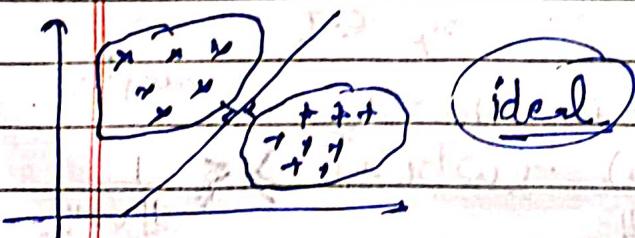


If we divide with $\delta/2 \Rightarrow$

$$\frac{w^T x + b}{\delta/2} = 1 \Rightarrow \bar{w}^T x + \bar{b} = 1$$

$$\text{Sim}'' \Rightarrow \bar{w}^T x + \bar{b} = -1$$

Ques → What happens when data is not linearly separable?



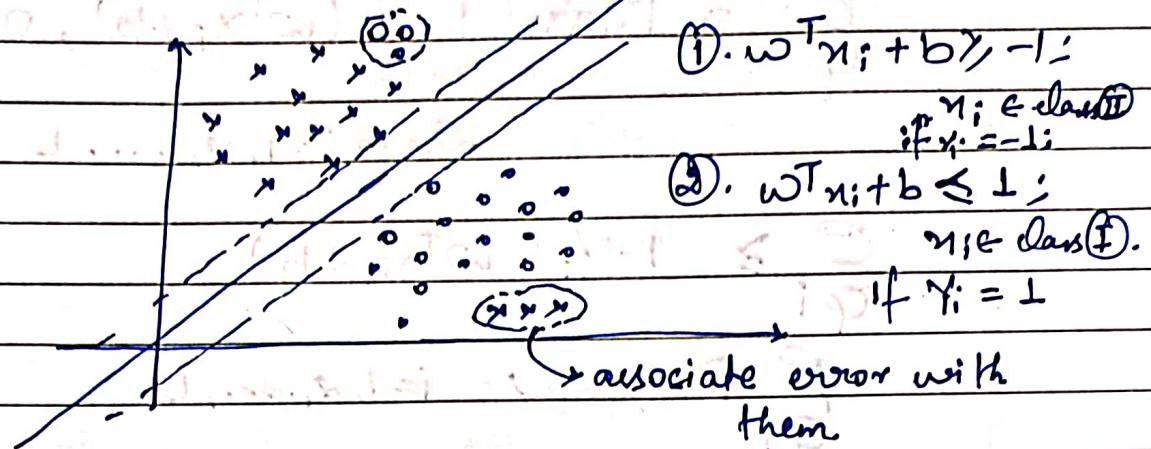
In general $\rightarrow Y_i(\omega^T x_i + b) \leq 1$

for ever to occur

Date _____
Page No. _____

error case \rightarrow

Real word data \rightarrow



final optimization problem

① Try minimizing the error.

② Try maximizing the margin

$$\underset{\omega, b}{\text{Min}} \frac{1}{2} \omega^T \omega + C \cdot \sum_{i=1}^n \text{Max}(1 - Y_i(\omega^T x_i + b), 0)$$

so, for error \rightarrow

$$1 - Y_i(\omega^T x_i + b) \geq 0 \rightarrow \text{error}$$

$< 0 \rightarrow \text{No error}$

$$\underset{\omega, b}{\text{Min}} \frac{1}{2} \|\omega\|^2 + C \cdot \sum_{i=1}^n (1 - Y_i(\omega^T x_i + b))$$

$\uparrow L(u)$

$\rightarrow \text{loss fn. (Hinge)}$

$$L(u) = \text{Max}(u, 0)$$



$$\mathcal{E}_i = \max(1 - \gamma_i(\omega^T x_i + b), 0)$$

where $i = 1, 2, \dots, L$

$$\mathcal{E}_i \geq 1 - \gamma_i(\omega^T x_i + b)$$

$$\mathcal{E}_i \geq 0 ; i = 1, 2, \dots, L$$

$$\Rightarrow \underset{\omega, b, \mathcal{E}_i}{\text{Min}} \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^L \mathcal{E}_i$$

(convex programming problem)

Final

Optimization

$$\text{problem is reduced to } \underset{\omega, b, \mathcal{E}_i}{\text{Min}} \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^L \mathcal{E}_i$$

$$\text{s.t. } \gamma_i(\omega^T x_i + b) \geq 1 - \mathcal{E}_i,$$

$$\mathcal{E}_i \geq 0 ; i = 1, 2, \dots, L.$$

Lagrangian fn. \rightarrow

$$L(\omega, b, \mathcal{E}_i, \alpha, \alpha^*) = \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^L \mathcal{E}_i$$

$$\begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_L \end{bmatrix} \begin{bmatrix} \alpha_1^* \\ \alpha_2^* \\ \vdots \\ \alpha_L^* \end{bmatrix} \alpha^*$$

$$+ \sum_{i=1}^L -\alpha_i (\gamma_i(\omega^T x_i + b) - 1 + \mathcal{E}_i)$$

$$+ \sum_{i=1}^L -\alpha_i^* \cdot \mathcal{E}_i$$

So \rightarrow

$$L(w, b, \epsilon_i, \alpha, \alpha^*) = \frac{1}{2} w^T w + C \sum_{i=1}^L \epsilon_i$$

$$= \sum_{i=1}^L \alpha_i \cdot (y_i \cdot (w^T x_i + b) - 1 + \epsilon_i)$$

$$- \sum_{i=1}^L \epsilon_i \cdot \alpha_i^*$$

$$\frac{\partial L}{\partial w} = 0 \Rightarrow w = \sum_{i=1}^L \alpha_i \cdot y_i \cdot x_i$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{i=1}^L \alpha_i \cdot y_i = 0$$

$$\frac{\partial L}{\partial \epsilon_i} = 0 \Rightarrow C - \alpha_i - \alpha_i^* = 0 \quad ; \text{ for } i=1, 2, \dots, L$$

$$\alpha_i \cdot (y_i \cdot (w^T x_i + b) - 1 + \epsilon_i) = 0$$

$$\alpha_i^* \epsilon_i = 0 ; i=1, 2, 3, \dots, L$$

$$y_i \cdot (w^T x_i + b) \geq 1 - \epsilon_i ; i=1, 2, \dots, L$$

$$\epsilon_i \geq 0 ; i=1, 2, \dots, L$$

$$\Rightarrow \text{Max} \sum_{i=1}^L \sum_{j=1}^L \alpha_i \alpha_j y_i y_j x_i^T x_j - \sum_{i=1}^L \alpha_i y_i \left(\sum_{j=1}^L \alpha_j y_j x_i^T x_j \right) + \sum_{i=1}^L \alpha_i$$

$$- \sum_{i=1}^L \sum_{j=1}^L \dots \dots \dots$$

$$\text{Subject to} \rightarrow \sum_{i=1}^L \alpha_i y_i = 0$$

$$C - \alpha_i - \alpha_i^* ; i=1, 2, \dots, L \quad \{ \alpha_i \leq C \}$$

$$\alpha_i = C - \alpha_i^*$$

Min

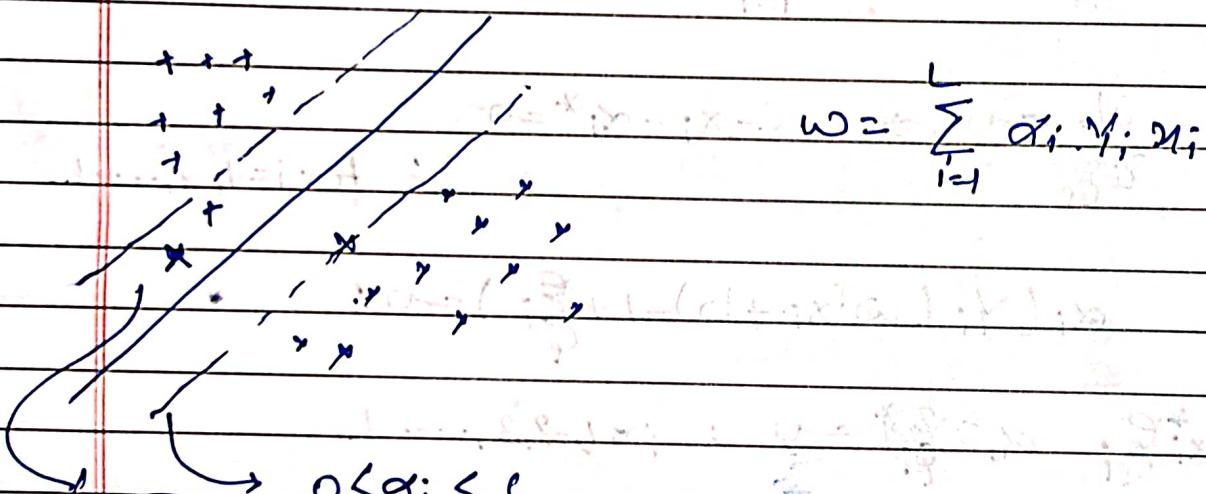
$$\rightarrow \underset{(\alpha)}{\text{Max}} + \frac{1}{2} \sum_{i=1}^L \sum_{j=1}^L \alpha_i \alpha_j y_i y_j (\omega^T \gamma_j) + \sum_{i=1}^L \alpha_i$$

Subject to →

$$\sum_{i=1}^L \alpha_i y_i = 0$$

$$0 \leq \alpha_i \leq C ; i=1, 2, \dots, L.$$

Geometrical Representation →



$$\alpha_i = C$$

→ they are contributing in w.

$$\text{if } \alpha_i \in (0, C) \rightarrow \alpha_i * > 0$$

$$\text{then } \alpha_i * = 0$$

$$\text{then } y_i * (\omega^T \gamma_i + b) - 1 = 0$$

$$y_i b = 1 - y_i \cdot (\omega^T \gamma_i)$$

$$\left[b = y_i - \omega^T \gamma_i \right] \left[b = \frac{1}{y_i} - \omega^T \gamma_i \right] \quad \forall 0 < \alpha_i < C$$

(Q1).

$$\left\{ b = \gamma_i - \sum_{j=1}^L \alpha_j \gamma_j \cdot \gamma_j^T \cdot \gamma_i \right\}.$$

final value of b is mean of all values of b .

for non-linearity \rightarrow

changes.

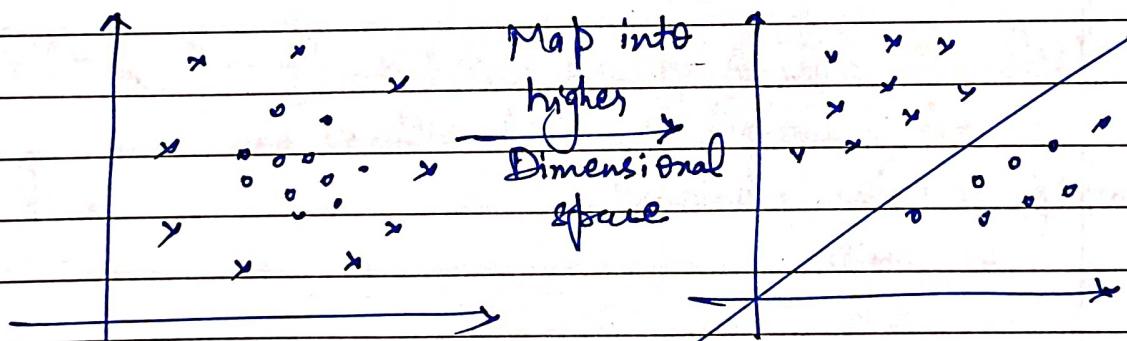
$$\gamma_i \rightarrow \phi(\gamma_i)$$

$$\gamma_i^T \gamma_j \rightarrow \phi(\gamma_i)^T \cdot \phi(\gamma_j)$$

$$w \rightarrow \sum_{i=1}^L \alpha_i \gamma_i \phi(\gamma_i)$$

$$w^T \gamma = -$$

$$\gamma_i^T \gamma_j = \text{Kernel} \Rightarrow K(\gamma_i, \gamma_j)$$



To make it linearly
separable