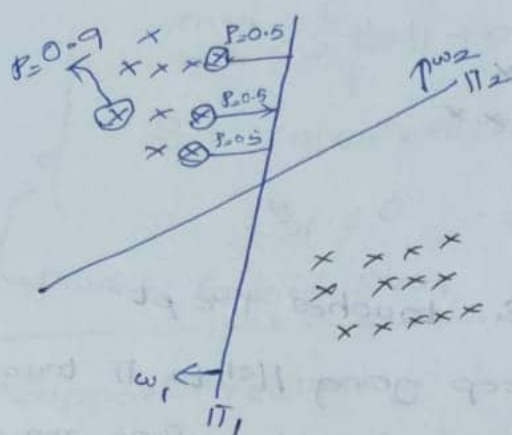


23/12/19

# Support Vector Machines.

Geometric Intuition :-



We can draw so many hyperplanes to separate these two (+ve, -ve) pts.

Actually many hyperplanes that separate my +ve & -ve pts.

# We need to choose the best hyperplane.

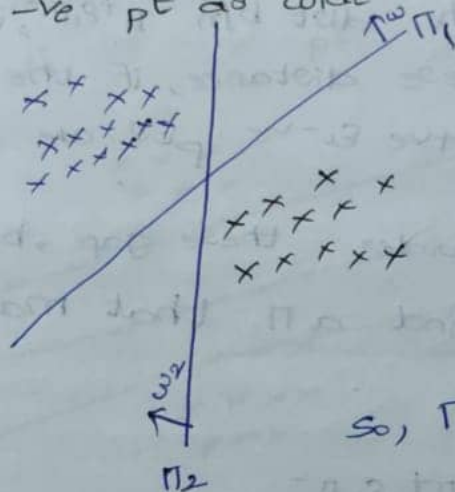
Suppose if i take  $\pi_1$ , so many pts are close to hyperplane

→ In Log. Reg we said  $P(y_i = 1) = \sigma(w^T x_i)$

So, pts close to hyperplane have low prob. (about 0.3)

& pts that are far away have high probability.

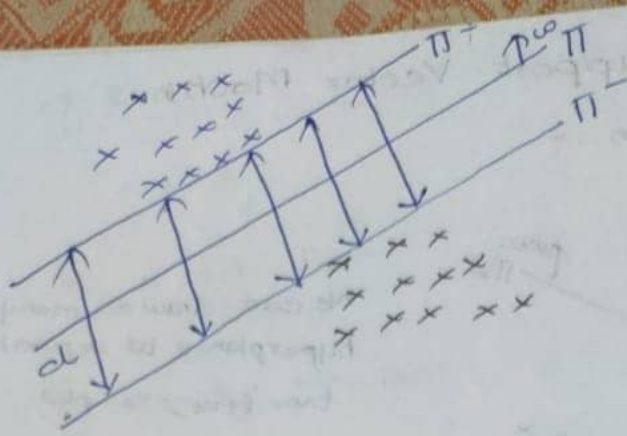
# Key Idea: - Find a hyperplane that separates my +ve & -ve pt as wide as possible.



Blue & Black pts are as far away as possible to  $\pi_1$  as compared to  $\pi_2$

So,  $\pi_1$  better than  $\pi_2$

→ So such a hyperplane is called Margin Maximising hyperplane



→  $\pi^+$  is  $\parallel$  to  $\pi$  & touches +ve pt.

Similarly, if i keep going  $\parallel$  to  $\pi$  towards my -ve pts i get  $\pi^-$  (touches first neg. pt)

$\pi^-$  is  $\parallel$  to  $\pi$ .

Since  $\pi^+$  &  $\pi^-$  are  $\parallel$  to each other, so dist. b/n  $\pi^+$  &  $\pi^-$  are having const dist.

Now, we want to find a hyperplane  $\pi$  s.t. if i go  $\parallel$  to  $\pi$  & touch +ve pt with  $\pi^+$  &  $\parallel$  to  $\pi$ , touch -ve pt with  $\pi^-$ , the margin is the dist b/n  $\pi^+$  &  $\pi^-$ , i need to maximise these distance, if the dist. is maximized +ve & -ve pts are quite faraway from  $\pi$ ,

→ The wider these gap, better for us.

→ SVM: Try to find a  $\pi$  that maximises the margin.

↓  
dist b/n  $\pi^+$  &  $\pi^-$

So, if margin is high chances of misclassification will decrease & my Gener. Error will improve.

As Margin Incr. my generaliz. acc improve  
↓  
Future/unseen



## → Dual Form of SVM:-

Soft-margin:-

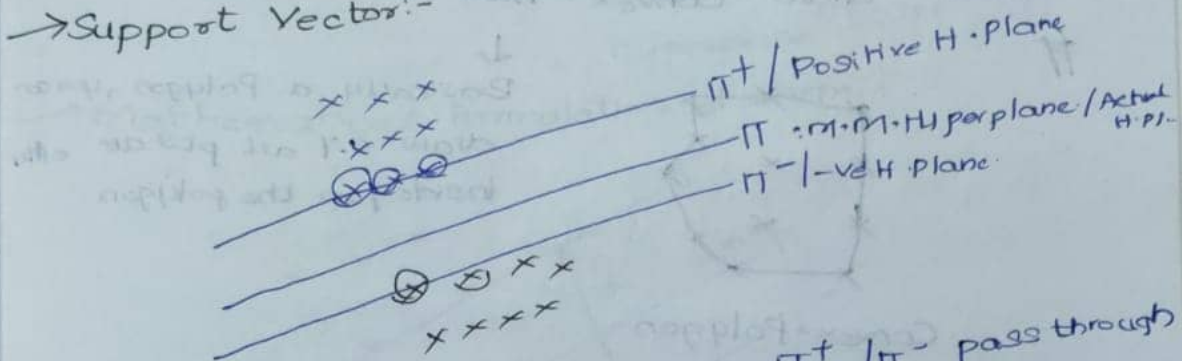
$$\min_{w, b} \frac{1}{2} \|w\|^2 + C \cdot \sum_{i=1}^n \xi_i$$

$$\text{s.t. } y_i (w^T x_i + b) \geq 1 - \xi_i, \forall i$$

$$\xi_i \geq 0$$

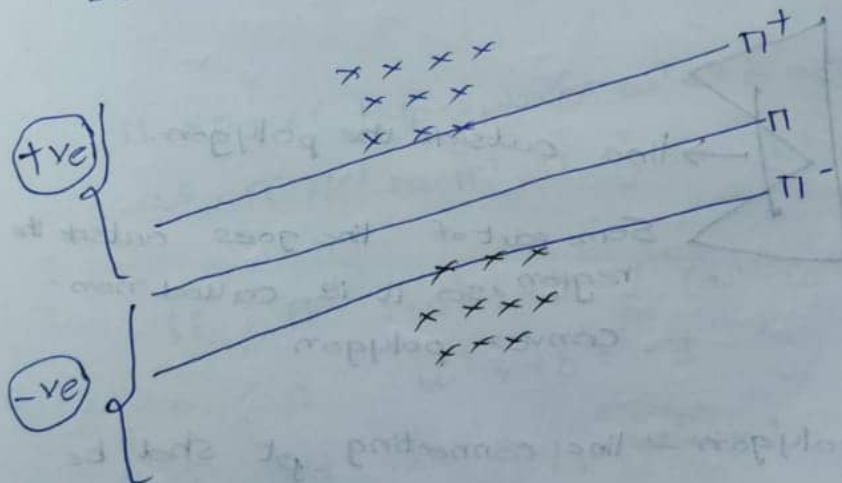
→ Primal Form of SVM

## → Support Vector:-



The pts through which  $\pi^+ / \pi^-$  pass through are called Support Vectors.

I not only use  $\pi^+ \& \pi^-$ , Any pt above  $\pi$  is +ve. Below  $\pi$  is -ve pt.



## Alternative Geometric Interpretation of SVM:-



First draw Convex-hull

#

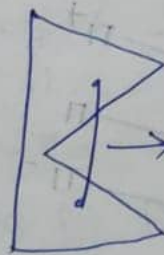


Convex-Polygon

↓  
Basically a Polygon, you can draw s.t all pts are either inside/on the polygon.



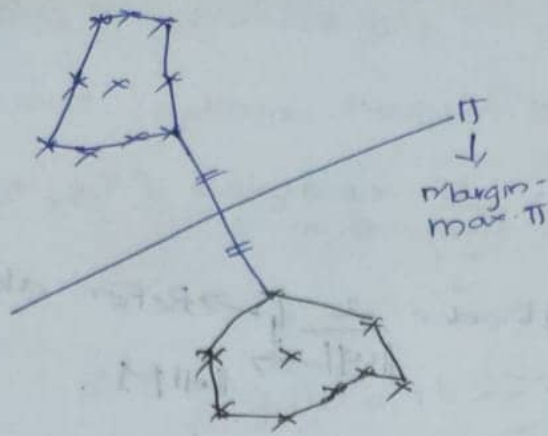
If I need to go from one pt to another pt the line is going inside the triangle only.  
So  $\Delta$  is Convex-Polygon.



→ line outside the polygon.

Some part of line goes outside the region, so it is called non-convex polygon.

→ So, convex polygon - line connecting pt shd be inside the polygon.



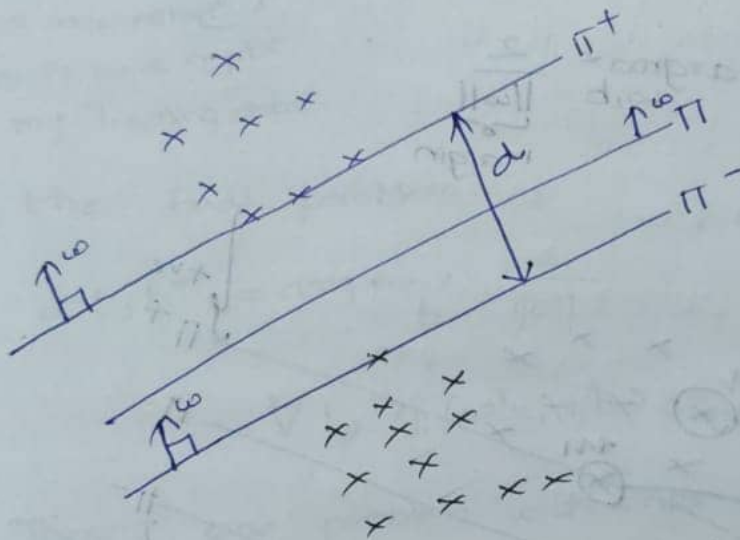
① Construct a convex hull for +ve pts & -ve pts.

② Find the <sup>shortest</sup> line connecting these hulls

③ Bisect the line

The plane that bisects it is called margin maximising hyperplane.

→ Mathematical Formulation of SVM:-



$\Pi \rightarrow$  Eq. of Hyperplane:  $w^T x + b = 0$  ( $w \perp \Pi$ )

as,  $\Pi \parallel$  to  $\Pi^+$ ,  $\Pi^-$ , so  $w \parallel$  to  $\Pi^+$ ,  $\Pi^-$

If  $\Pi^+$ :  $w^T x + b = 1$  (as  $w$  is  $\parallel$  to  $\Pi^+$ )

$\Pi^-$ :  $w^T x + b = -1$

I did not say,  $w^T w = 1$ .

#  $w^T w \neq 1$ ;  $w$  is not a unit vector.

Let,  $w$  be some vector, not necessarily unit vector



2 1 → Refer alexander p...

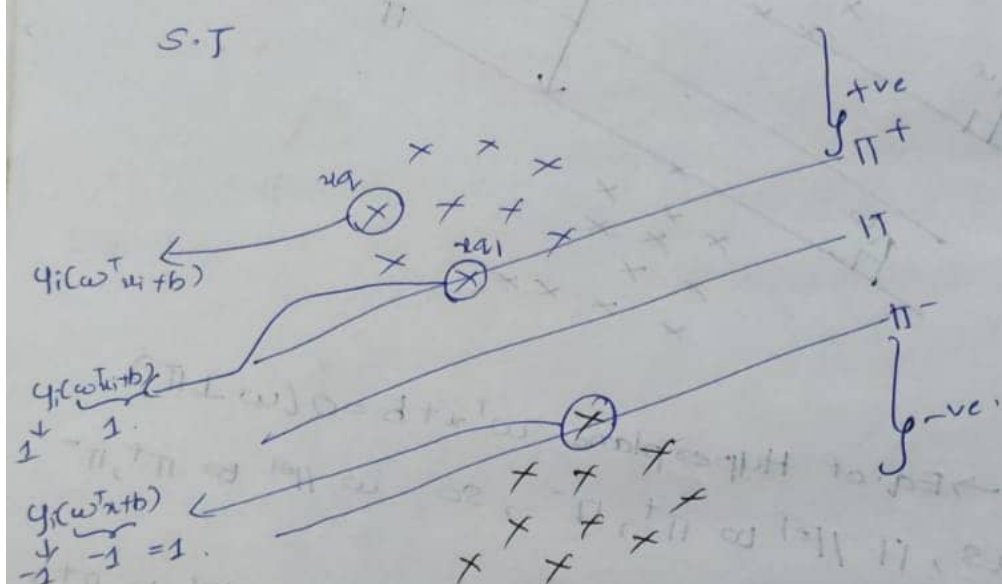
→ We want to find  $w^* \in \mathbb{R}^n$  s.t. which maximise

S.T (constraint)

$$= \arg \max_{\omega, b} \frac{2}{\|\omega\|}$$

margin

S.T



→ For  $2q$ :-

$$y_i(\underbrace{\omega^T x_i + b}_{\text{dist}}) > 1$$

For 24.1.7

$$\begin{aligned} y_i(a^T x_i + b) &= 1 \\ \downarrow & \quad \text{dist} = 1 \\ +1 & \\ \downarrow & \\ w^T x + b &= 1 \end{aligned}$$

So, For S.V it's 1.

→ So, our optimiz. Problem is

$$(w^*, b^*) = \arg \max_{w, b} \underbrace{\frac{2}{\|w\|}}_{\text{margin}}$$

S.T  $y_i(w^T x_i + b) \geq 1$  for all  $x_i$

↓

If all +ve pts are above  $\pi^+$   
& all -ve pts are above  $\pi^-$

It's n-constraint  
because we've n-pts  
in my Training data.

So, the Final problem is

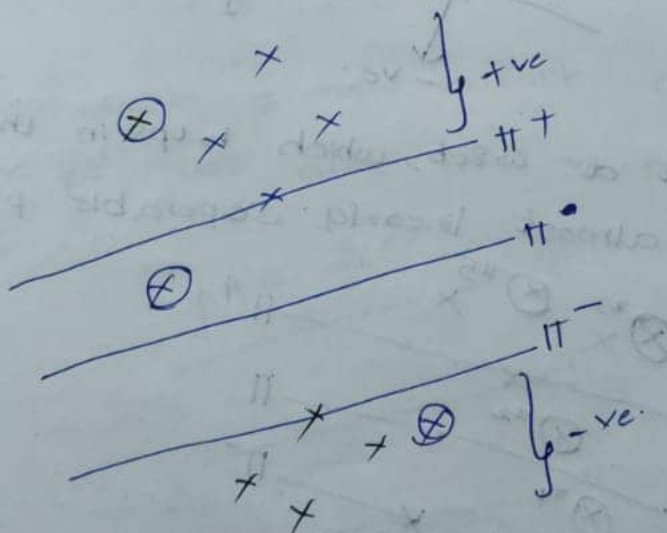
$$w^*, b^* = \arg \max_{w, b} \frac{2}{\|w\|}$$

s.t  $\forall i, y_i(w^T x_i + b) \geq 1$

Consts.  
Optimiz. Problem.  
of S.V.M.

There's one problem with this.

→ This only works  
if data is linearly  
separable.





So, here there are 3 pts, which is violating our rule, +pt below  $\pi^-$

-ve pt above  $\pi^+$

& another pt (ve) in no man's land (ie, b/w  $\pi^+$ ,  $\pi^-$ )

So, these 3 pts never satisfy our constraint

So, here it's almost Linearly Seperable, not linearly seperable, in such a case it's difficult to find  $w, b$ .

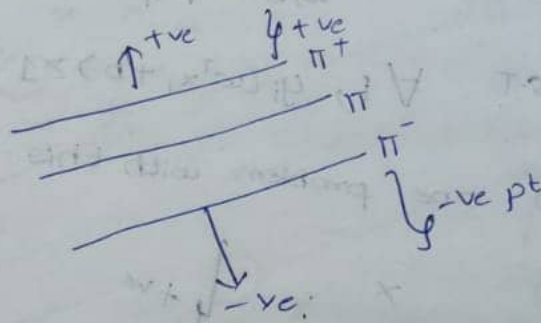
So,

$$(w^*, b^*) = \arg \max_{w, b} \frac{2}{\|w\|}$$

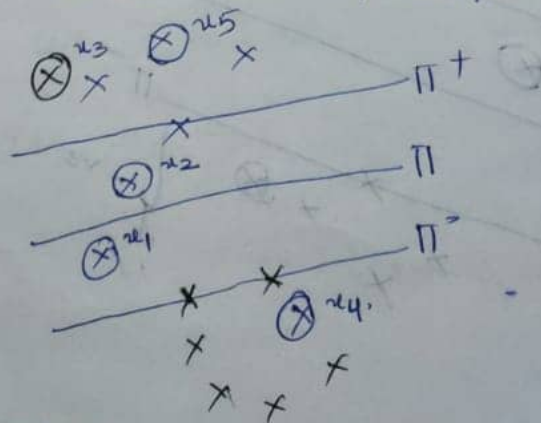
$$s.t. \quad y_i(w^T x_i + b) \geq 1 \rightarrow \text{is called}$$

Hard-Margin S.V.M

# Very strong saying that



→ Can we find a  $w$  &  $b$ , which helps in the problem  
& useful for almost linearly Seperable pts.





→ For  $x_1 \rightarrow y_i = 1$

$$w^T x_i + b$$

$$\Pi: w^T x_i + b = 0$$

$$\Pi^+: w^T x_i + b = +1$$

$$\Pi^-: w^T x_i + b = -1$$

$x_1$  is b/n  $\Pi$  &  $\Pi^-$

So, for  $x_1, y_i (w^T x_i + b) = -0.5$  (b/n  $\Pi$  &  $\Pi^-$ )

$\downarrow$   $\downarrow$   
 $+1$  half way b/n  $\Pi$  &  $\Pi^-$

So,  $y_i (w^T x_i + b) = 1 - (1.5)$

$\downarrow$   $\downarrow$   
 $3_1$   $3_1$

This 1 because, as it's the pt, & so,  $y_i (w^T x_i + b) = 1$  for  $\Pi^+$ , actually  $x_1$  need to be above  $\Pi^+$ , but it's not there so, we'll write it as follows

→ For,  $x_4 \rightarrow y_i = 1$

$$y_i (w^T x_i + b) = -1.5$$

$$y_i (w^T x_i + b) = 1 - (2.5)$$

$\downarrow$   
 $3_4$

→ For  $x_2 \rightarrow y_i = 1$

$$y_i (w^T x_i + b) = 0.5 \rightarrow \text{half way b/n } \Pi \text{ \& } \Pi^+ \text{ so } 0.5$$

$$y_i (w^T x_i + b) = 1 - (0.5)$$

$\downarrow$   
 $3_2$

→ For  $x_5 \rightarrow y_i = 1$

$$y_i (w^T x_i + b) = 1.5 \rightarrow \text{it's } > 1 ; \text{ so no worries}$$

(ie,  $y_i w^T x_i + b = 1$  for the pt)

→ We create a new variable called  $z_i$  s.t if

the pt lies above  $\Pi^+$  & the pt lies below  $\Pi^-$

$$z_i = 0$$

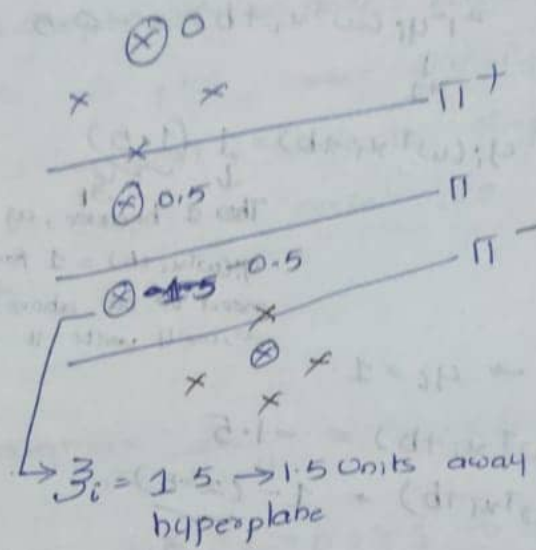
↑ +ve  $z_i = 0$

↓ -ve  $z_i = 0$

↑ +ve  $z_i = 0$

So, if a pt lies anywhere except in the regions showed above it

$$\xi_i \geq 0$$



as  $\xi_i \uparrow$ ; pt is farther away from the correct hyperplane in the incorrect reg direction.

So, for every pt  $x_i \rightarrow \xi_i$

$$\xi_i = 0, \text{ if } y_i(\omega^T x_i + b) > 1$$

$$\xi_i \geq 0 \text{ \& it is}$$

Correctly classified

Equal to the some units of dist. away from the correct hyper plane with  $\pi^+ / \pi^-$  in the incorrect direction.

So, now,

$$(\omega^*, b^*) = \arg \max_{\omega, b} \frac{2}{\|\omega\|} = \arg \min_{\omega, b} \frac{\|\omega\|}{2}$$

$$\# \arg \max f(u) = \arg \min -f(u) / \frac{1}{f(u)}$$



So, now

$$\omega^*, b^* = \arg \min_{\omega, b} \frac{\|\omega\|}{2} + c \cdot \frac{1}{n} \sum_{i=1}^n \xi_i$$

$\nearrow$  margin

s.t

$$y_i(\omega^T x_i + b) \geq 1 - \xi_i \quad \forall i$$

$$\xi_i \geq 0$$

# Previously we saw for every misclassified pt,  $i$   
can write  $y_i(\omega^T x_i + b) = 1 - \xi_i$  where  $\xi_i$  is the  
 $\xi_i > 0$ .

So, for all misclassified pt

$$y_i(\omega^T x_i + b) \geq 1 - \xi_i \quad \forall i \quad \left\{ \begin{array}{l} \text{Incorrect classified} \\ \text{pt} \end{array} \right.$$

where  $\xi_i > 0$

For correctly classified pt  $\xi_i = 0$ .

# We want to minimise errors / misclassifications  
It means,  $\xi_i > 0$  for misclassified pts.

||

as per  $\pi^+, \pi^-$   $\min \sum \xi_i \rightarrow$  we want to minimise sum of  $\xi_i$

$\rightarrow c \cdot \frac{1}{n} \sum_{i=1}^n \xi_i \rightarrow$  dist.   
  $\downarrow$    
 hyper parameter   
  $\nearrow$  avg. dist of misclassified pts from correct hyperplanes

$$\rightarrow (\omega^*, b^*) = \arg \min_{\omega, b} \frac{\|\omega\|}{2} + c \cdot \frac{1}{n} \sum_{i=1}^n \xi_i$$

$\nearrow$  margin

s.t  $y_i(\omega^T x_i + b) \geq 1 - \xi_i \quad \forall i$    
  $\xi_i \geq 0$

$\rightarrow$  avg. dist for misclassified pts.   
  $\rightarrow$  Constraint

We can think it as,

$$(\omega^*, b^*) = \underset{\omega, b}{\operatorname{argmin}} \underbrace{\frac{\|\omega\|^2}{2}}_{\text{regularization (like } L_2)} + \underbrace{c \cdot \frac{1}{n} \sum_{i=1}^n \xi_i}_{\text{Loss}}$$

In Log. reg. we had

$$\min_{\omega} (\log\text{-loss}) + \lambda (\text{reg}).$$

Here,

$$\min_{\omega, b} \quad c \cdot \text{hinge loss} + \text{Regularization}$$

↓  
hyper-parameter

as  $C \uparrow$ , loss  $\uparrow$ ; so - tendency to make mistakes  $\downarrow$   
 $\Rightarrow$  Overfit

$C \downarrow$ , tendency to make mistakes  $\uparrow$   
 $\Rightarrow$  Underfit

$\therefore C \uparrow$ , C-dominates; Overfit (high-var)

$C \downarrow$ ; reg-dominates; Underfit (high-bias)

$$\rightarrow \therefore (\omega^*, b^*) = \underset{\omega, b}{\operatorname{argmin}} \frac{\|\omega\|^2}{2} + c \cdot \frac{1}{n} \sum_{i=1}^n \xi_i$$

$$\text{S-T } y_i(\omega^T x_i + b) \geq 1 - \xi_i \quad \forall i$$
$$\xi_i \geq 0$$

It's Soft margin-SVM

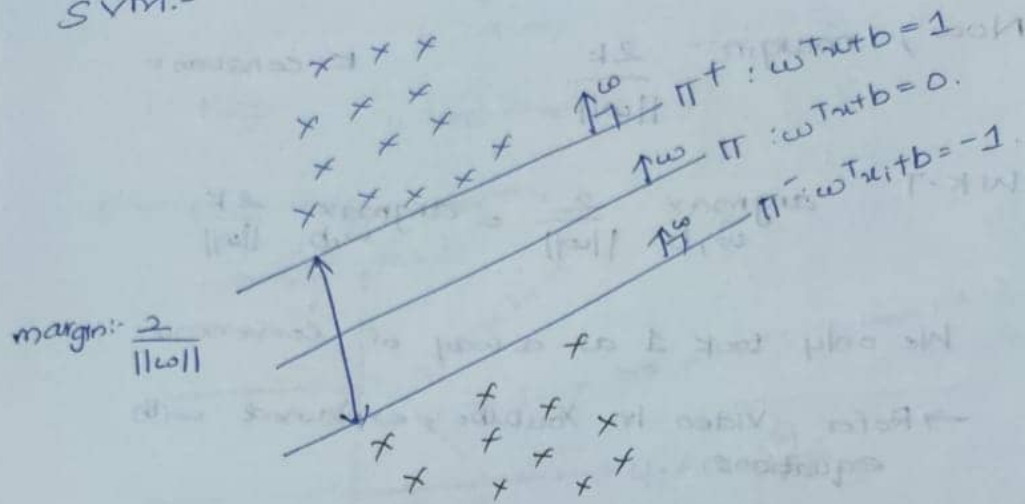
$\rightarrow$  Hard-Margin doesn't allow errors.

But Soft-Margin allows errors but it says minimize errors.



→ Why take +1 & -1:-

SVM:-



Q) Why +1 & -1 on R.H.S of  $\pi^+ \& \pi^-$

Here, we're not saying  $\|w\| \neq 1$ ;  $w$  could be any vector. It need not be a unit vector.

$$w \perp \pi, \pi^+, \pi^-$$

$$\rightarrow \text{margin} = \frac{2}{\|w\|}$$

Our whole task is to maximise the margin

$$w^*, b^* = \arg \max_{w, b} \frac{2}{\|w\|}$$

$$\text{s.t. } y_i(w^T x_i + b) \geq 1$$

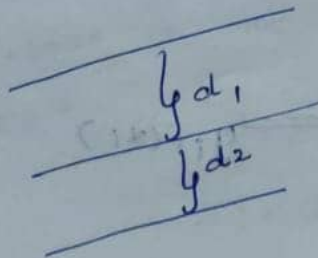
① Let's assume,

$$\pi^+ : w^T x + b = k$$

$$k > 0$$

$$\pi^- : w^T x + b = -k$$

we want dist shd be same



$d_1 = d_2$ , that's why we took  $k \& -k$ , but not  $k_1, -k_2$ .

So, +ve, -ve hyperplane are equivalently far away from  $\Pi$

Now, margin:-  $\frac{2k}{\|w\|}$   $k$ : constant

$$\text{W.K.T } \arg\max_{w,b} \frac{2}{\|w\|} = \arg\max_{w,b} \frac{2k}{\|w\|}$$

We only took 1 as a way of convenience.  
 → Refer Video In Youtube, explained with equations.

②  $\Pi$ :  $w^T x + b = k$

$$\left(\frac{w}{k}\right)^T x + \left(\frac{b}{k}\right) = 1$$

$$(w')^T x + b' = 1$$

$w \perp \Pi$   
 $\|w\|$  need not be 1.

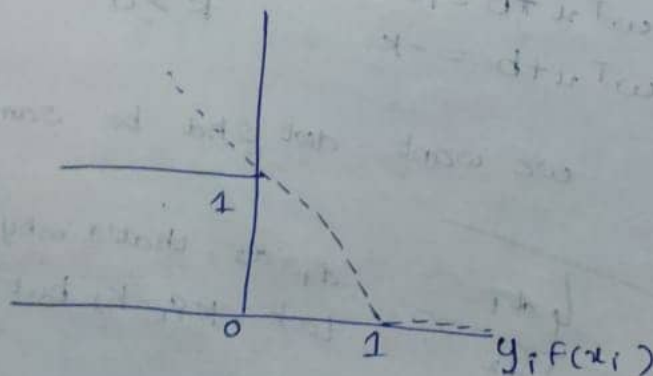
The reason to take +1 & -1 is to simplify math.

It's possible to take any value

All is because  $\|w\| \neq 1$ .

If  $\|w\| = 4$ ; margin = 2. It's not possible

→ Loss-minimization : Hinge-loss:-



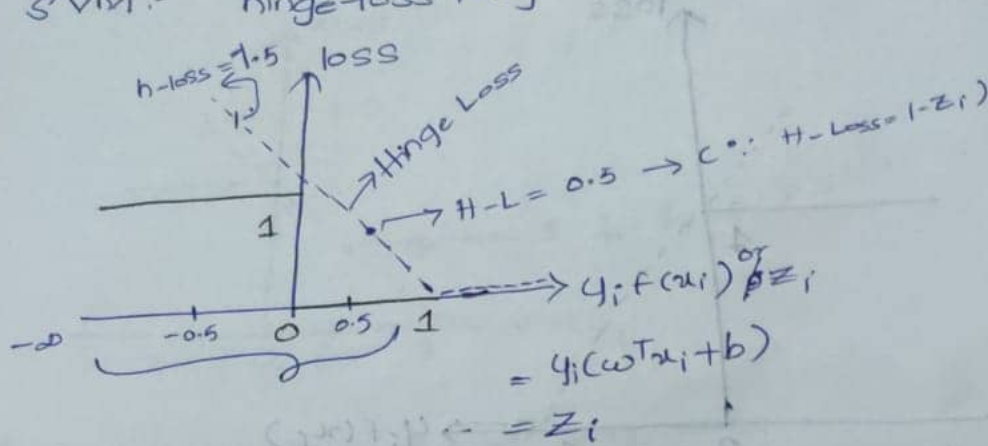


log-Reg:-  $\min \log\text{-loss} + \text{reg.}$

↓  
0-1-loss

lr-Reg:-  $\text{lr-loss} + \text{reg.}$

SVM:-  $\text{hinge-loss} + \text{reg.}$



when  $z_i$  is +ve,  $x_i$  is correctly classified

$z_i > 0$ ;  $x_i$  correctly classified

$$\rightarrow y_i (w^T x_i + b)$$

when  $z_i$  is -ve

$z_i < 0$ ;  $x_i$  is incorrectly classified.

Hence in 0-1, Loss we'll give 1 (Incorrect)

0 (Correct)

$\therefore$  0-1 loss is not continuous at 0, so 0-1 loss is not differentiable.

Hinge-Loss is a S-L from  $-\infty$  to 1,  $\Sigma$  from 1 it is 0.

Hinge-Loss is not differentiable at 1.

$\rightarrow$  Hinge-Loss:- if  $z_i \geq 1$ ; hinge-loss = 0

$z_i < 1$ ; hinge-loss =  $1 - z_i$

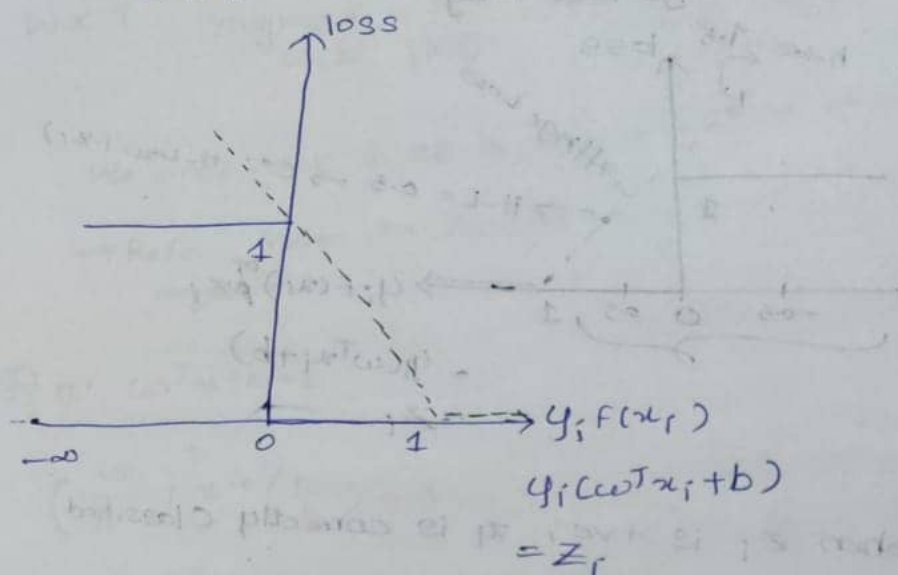
24/12 :-

→ Loss-Minimization : Hinge-loss

log-reg in loss point of view :- log-loss + reg

lr-reg :-  $l_2$ -loss + reg

SVM :- Hinge-loss + reg



when  $z_i$  is +ve,  $z_i > 0$ ,  $x_i$  is correctly classified

$z_i < 0$ ,  $x_i$  is incorrectly classified

# In 0-1 loss ; loss 1 if 1 incorrectly classified.

0-1 loss is not continuous, so not differentiable.

•

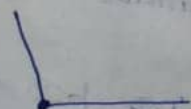
1

0

S.L from

→ Hinge-loss is  $-\infty$  to 1 & from 1 it's 0.

H-L is continuous, but it's not diff.



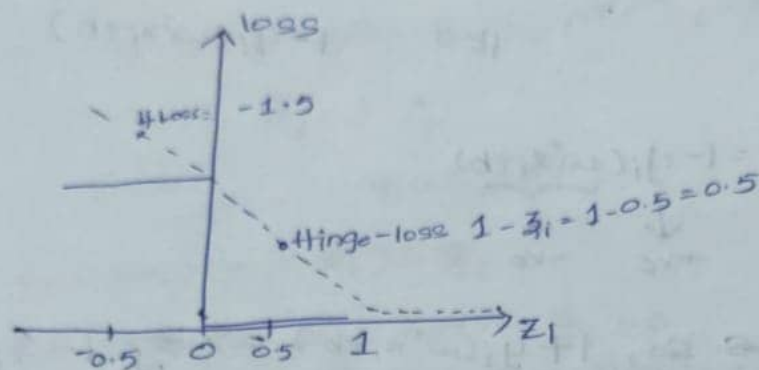


→ Let's look at how hinge loss behaves from  $z_i$  point of view

Hinge-Loss:-

if  $z_i \geq 1$ ; hinge-loss = 0

if  $z_i < 1$ ; hinge-loss =  $1 - z_i$  ①



The alternative way of hinge-loss is  $\max(0, 1 - z_i)$

Case 1:-  $z_i \geq 1$

$$1 - z_i = -ve$$

So,  $\max(0, -ve \text{ value}) = 0$ .

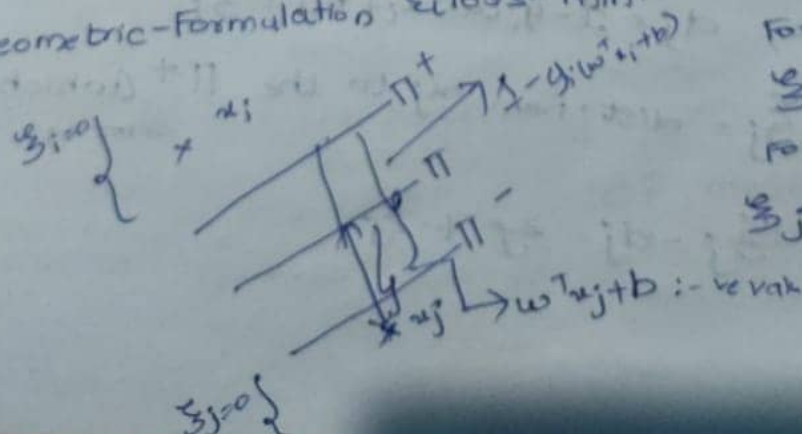
Case 2:-  $z_i \leq 1$

$$1 - z_i = +ve (> 0)$$

So,  $\max(0, +ve \text{ value}) = +ve \text{ value}$   
 $\Rightarrow 1 - z_i$

These two cases are exactly same as ①

→ Geometric-Formulation  $\mathcal{E}_{\text{loss}} = \min$



For  $x_i$

$$z_i = 0$$

For  $x_j$ ,

$$z_j = 1$$