

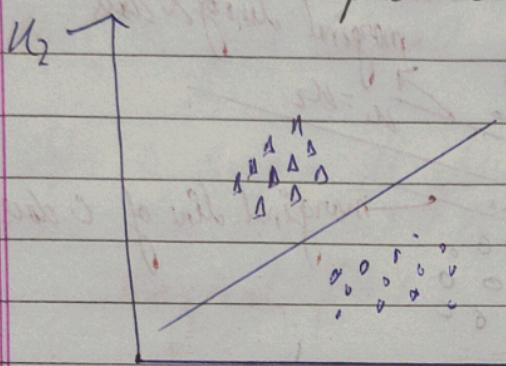
SVM & NAIVE BAYES

- Introduction to Support Vector Machine

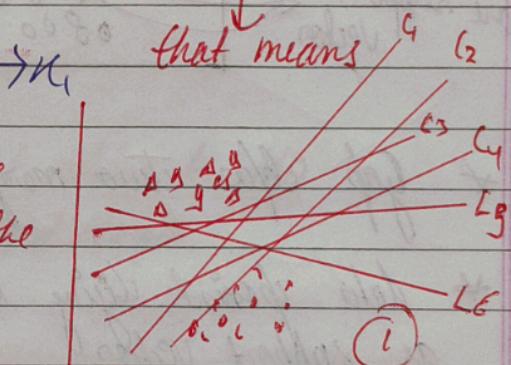
① Support Vector classifier - classification

② Support Vector Regressor - Regression

① what is the problem with logistic regression that we need to come up with SVM → SVC

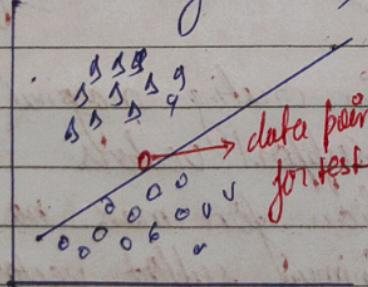


In logistic regression we find linear decision boundary which separates two class



Here all ~~the~~ lines were separating two classes the all ~~are~~ are possible line.

then let's say we got a datapoint of 0 class just above the decision boundary like ; then it is to classify as 1 class



but if some line passes through the center of gap b/w two classes datapoints. it would have classify it correctly.

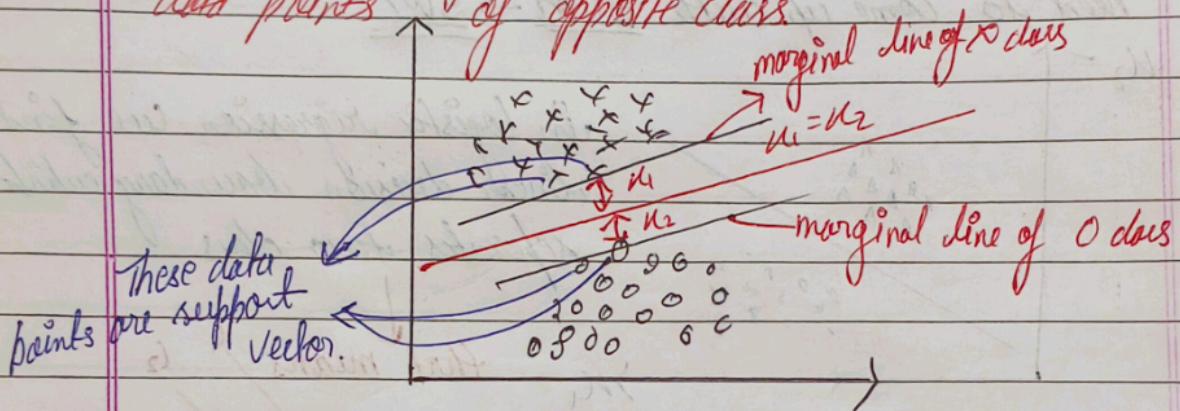
Logistic regression model doesn't care about the margin / space / distribution across two classes as in img ① we don't all lines are possible.

in logistic regression, when data points comes slightly left/right to decision line, it is misclassified

→ in order to solve this we got SVC

* Support Vector Classifier

→ make a decision boundary such that it is equidistant from two or more classes to all data points of opposite class



* Gap b/w two marginal line is margin

* data points lying on marginal line are called as support vector.

Note → Support vector are nearest datapoints of opposite class.

→ as they helps us in creating marginal line so they called as support vectors

→ There are no limitation for max no. of support vectors but minimum support vectors are required are one, one $\rightarrow 1 + 1 = 2$

help with less help of support vector \rightarrow we create marginal line/ plane.

and marginal plane/line ~~but~~ helps in deciding best fit line as it will be equidistant.

This is why SVM is also called as margin classifier.

Steps to get best fit boundary.

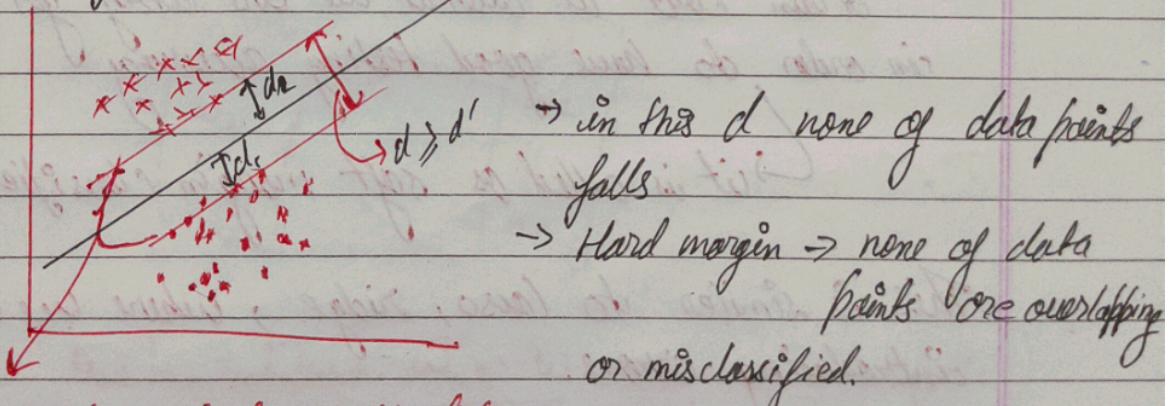
Step 1 \Rightarrow Find out the marginal line or plane using sum.

Step 2 \Rightarrow find out all possible decision boundary

Step 3 \Rightarrow Among all possible decision boundary find out the one that is equidistant from both the marginal plane.

We have two more concept related to this

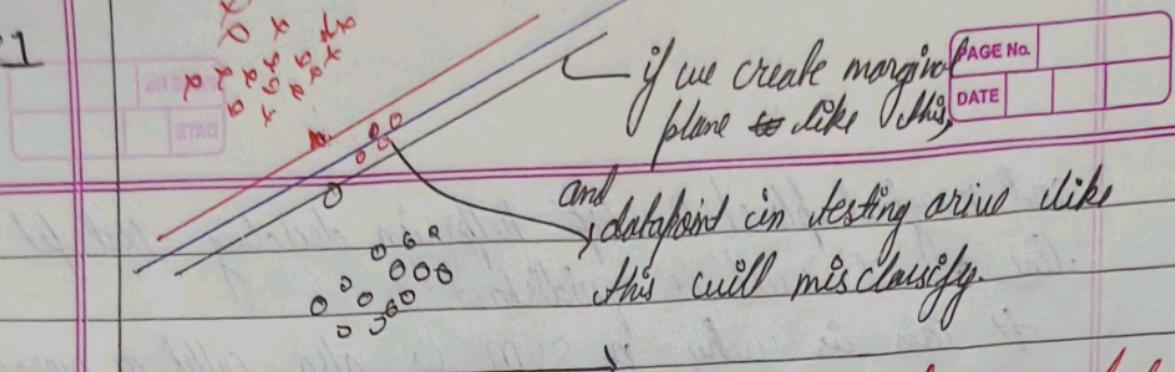
① Soft margin ② Hard margin



What we want is that in the distance d , none of our data points lie.

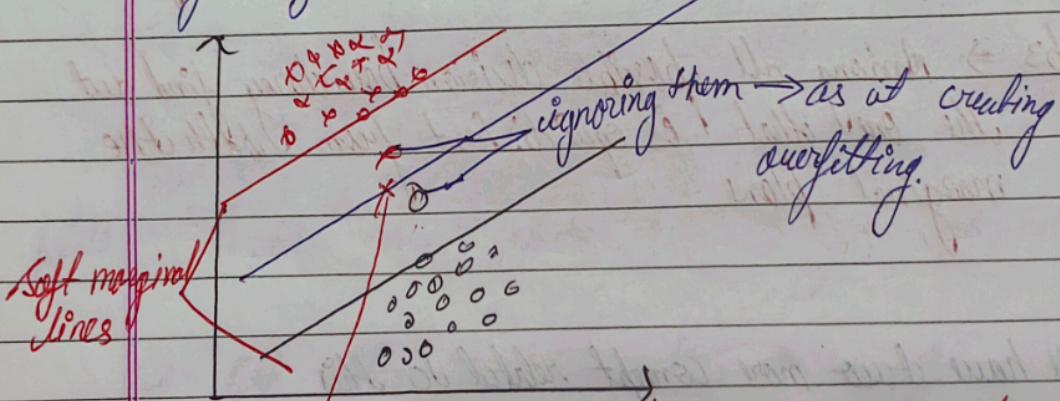
But in practical it is not possible

Case 1



As in this case, a lot of data points are clustered away from marginal plane, because only few or single data points come near to each other, which is causing a problem as the distance between the marginal planes/line is very short.

→ That's where we create a soft marginal line by ignoring some data points



This data point is misclassifying in train but at overall we are ready for it in order to have good testing accuracy

It is called as soft margin classifier

This is similar to Lasso, ridge, where we are introducing errors.

Maths behind Support Vector Machine

eqⁿ of line / plane / hyper plane

① eqⁿ of line

$$y = mx + c$$

$$y = \alpha_0 + \alpha_1 u_1$$

$$\underline{y = \alpha_0 u_0 + \alpha_1 u_1 + \dots + \alpha_n u_n} \Rightarrow \underline{y = -\alpha_0 - \alpha_1 u_0 - \dots - \alpha_n u_n}$$

but when it is 2D or more than 2D: $y = -\frac{a}{b}u_0 - \frac{c}{b}$
 then $y = \alpha_0 + \alpha_1 u_1 + \alpha_2 u_2 + \dots + \alpha_n u_n$

this is also a eqⁿ of
 hyperplane 2
 also written as $y = b + w_1 u_1 + w_2 u_2 + \dots + w_n u_n$

b = bias w = weights

$$w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \end{bmatrix} \quad u = \begin{bmatrix} u_1 \\ u_2 \\ u_3 \\ u_4 \end{bmatrix}$$

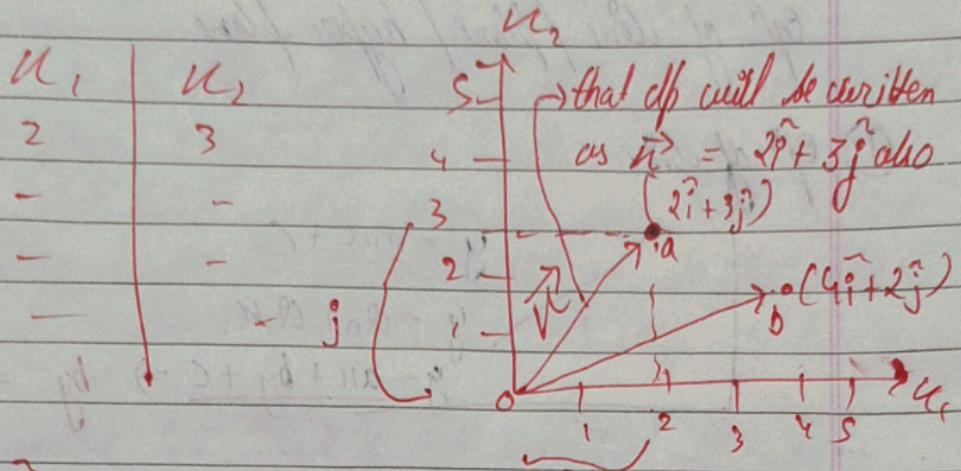
$$w^T u = [w_1 \ w_2 \ w_3 \ w_4] \begin{bmatrix} u_1 \\ u_2 \\ u_3 \\ u_4 \end{bmatrix}$$

$$w^T u = w_1 u_1 + w_2 u_2 + w_3 u_3 + w_4 u_4$$

$$w^T u + b = b + w_1 u_1 + w_2 u_2 + w_3 u_3 + w_4 u_4$$

$$y = w^T u + b$$

② data points as vectors



$$\vec{a} = \vec{u} = 2\hat{i} + 3\hat{j}$$

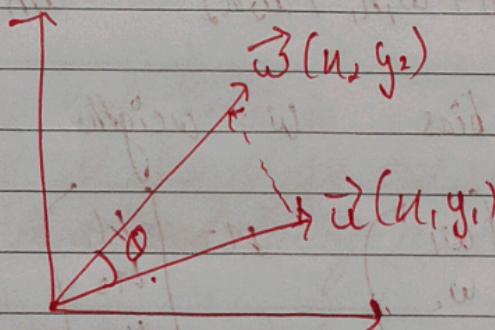
$$\vec{b} = \vec{v} = 4\hat{i} + 2\hat{j}$$

$$\vec{a} + \vec{b} = (2+4)\hat{i} + (3+2)\hat{j} \quad \left. \right\} \rightarrow \text{Vector addition.} \quad ③$$

$$\vec{a} + \vec{b} = 6\hat{i} + 5\hat{j}$$

$$\vec{a} - \vec{b} = (2-4)\hat{i} + (3-2)\hat{j} = 2\hat{i} + 1\hat{j} \quad \left. \right\} \rightarrow \text{Vector subtraction.} \quad ④$$

③ Dot product



$$\vec{w} \cdot \vec{u} = |\vec{w}| |\vec{u}| \cos \theta$$

\downarrow magnitude of w \downarrow magnitude of u

$\vec{w} \cdot \vec{u} \Rightarrow$ ~~projection of \vec{w} on \vec{u}~~ , θ is angle between \vec{w} and \vec{u}

~~Q~~ Gross Product

Now will we know (7, 8) of (-2, -3) lie on left or right of plane? below or above of plane?

$$3x + 5y + 3 = 0$$

$$\underline{3 \times 7 + 5 \times 8 + 3 = 0}$$

$$21 + 40 + 3 = 0$$

$$\boxed{64 > 0}$$

$$3x + 5y + 3 = 0$$

$$3x + 5y + 3 = 0$$

$$3(-2) + 5(-3) + 3 = 0$$

$$-6 + (-15) + 3 = 0$$

$$-24 + 3 = 0$$

$$-18 < 0$$

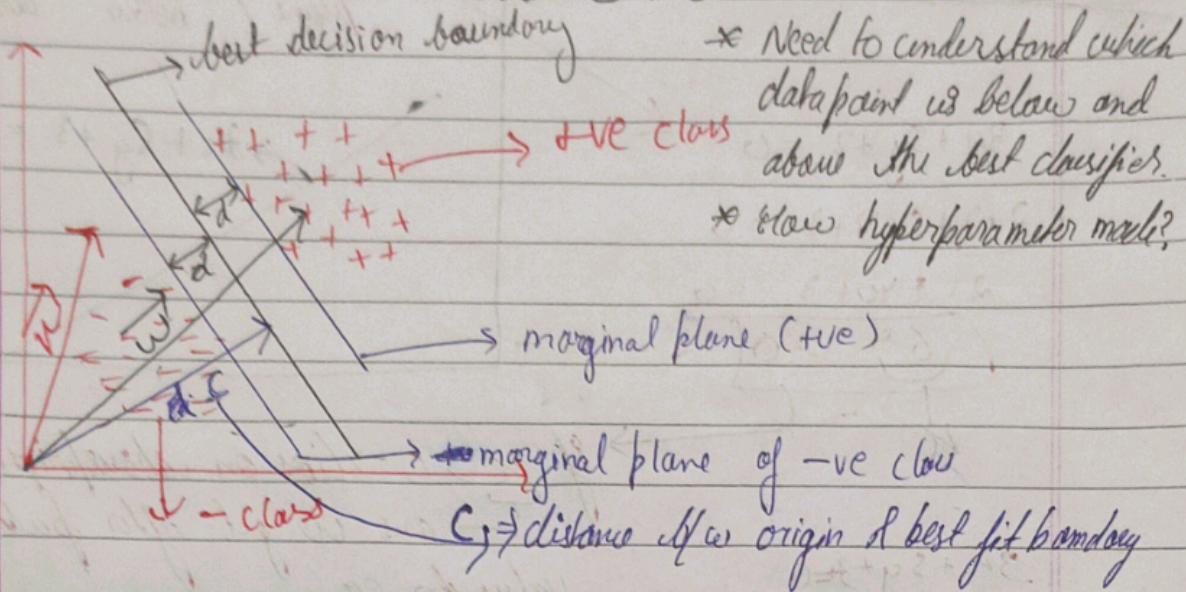
\rightarrow if $n > 0 \rightarrow$ lies on above/right to plane.
ans we get after putting value to eq.

\rightarrow if $n < 0 \rightarrow$ lies on below/left to plane / line.

Now we understood the basic mathematics behind sum,

then now lets dive deep into deep mathematics that How actually these concepts applied on sum.

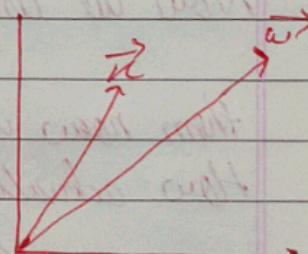
Mathematics behind S.V.C



- First we need a \vec{w} , which is ~~to~~ perpendicular to hyperplane (best decision fit line)
- also ~~to~~ to both marginal planes
- C is distance of hyperplane/optimal line from origin
- now we try to project \vec{n} on \vec{w}

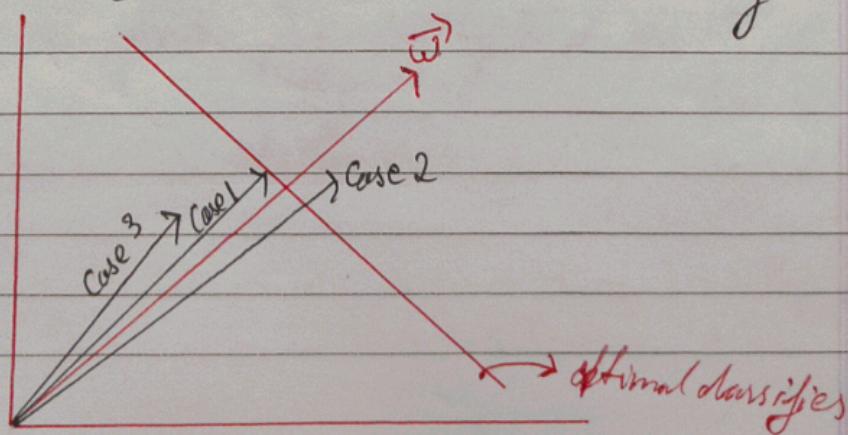
$$\vec{n} \cdot \vec{w} = |\vec{n}| \cos \theta \cdot |\vec{w}|$$

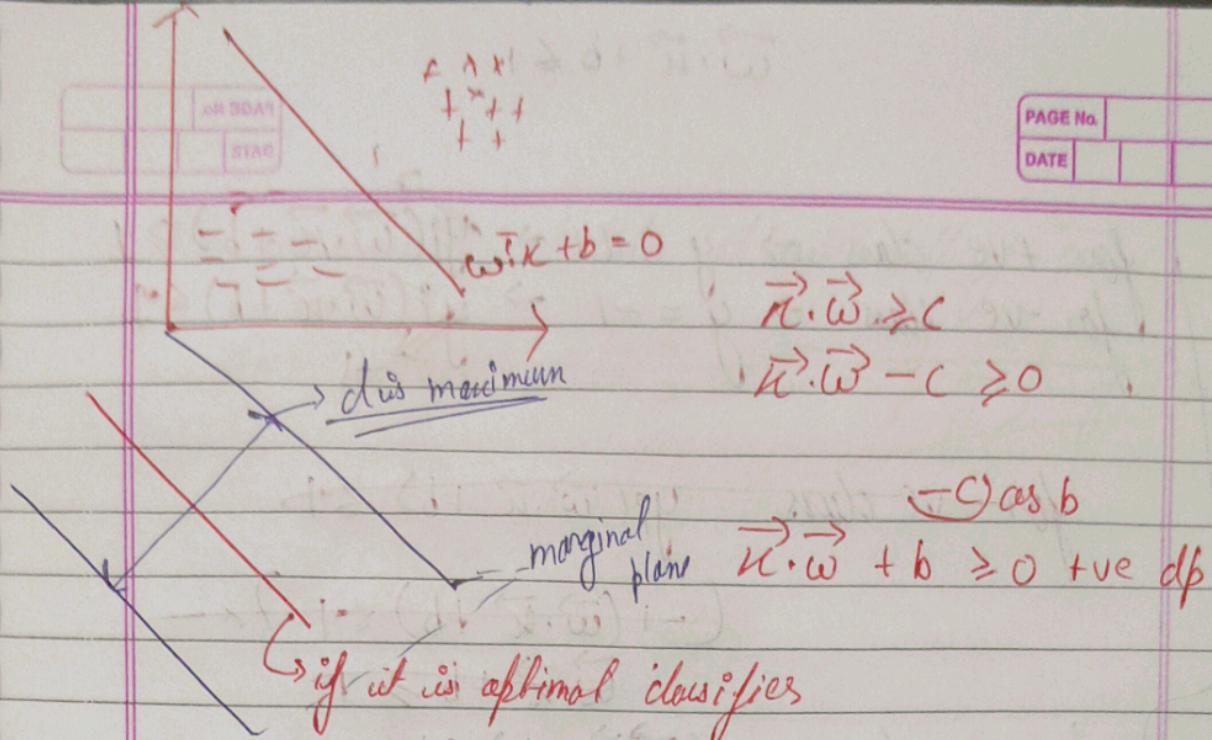
Case 1 if $\vec{n} \cdot \vec{w} = C$ {lie on best decision boundary}



Case 2 if $\vec{n} \cdot \vec{w} > C$ {lie above the decision boundary}

Case 3 if $\vec{n} \cdot \vec{w} < C$ {lie below the decision boundary}





y is +ve if $\vec{w} \cdot \vec{x} + b \geq 0$
 y is -ve if $\vec{w} \cdot \vec{x} + b < 0$

why equal $\{ \text{both}(+1, -1) \}$
 why only $\{ (+1, -1) \}$

because both marginal lines are
 equidistant from best decision
 boundary. ($d_1 = d_2$)

why only 1?

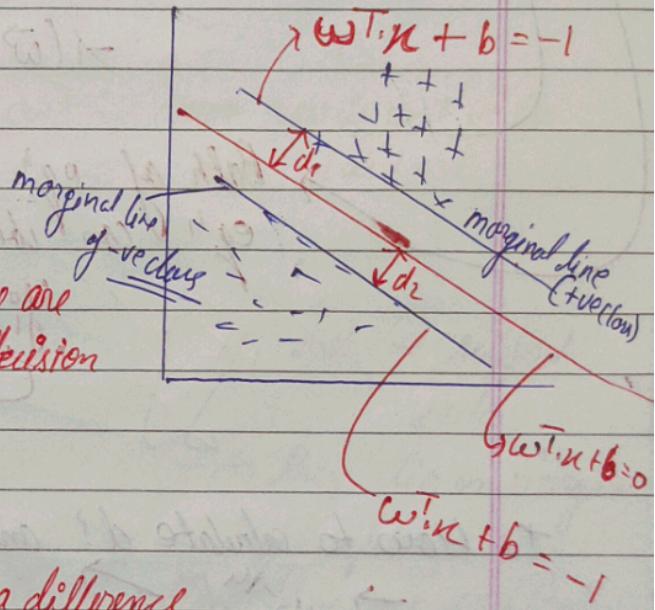
it doesn't make a difference

$$k + g = 1$$

$2k + 2g = 2 \rightarrow$ even if we multiply whole g with other no. line doesn't change

Now for distance $d \rightarrow$ we want to calculate distance(d) in such a way that no positive or negative point can cross marginal line

for +ve class df's $\vec{w} \cdot \vec{x} + b \geq 1$ we want to maximize d
 such that this constant holds true.
 for -ve class df's $\vec{w} \cdot \vec{x} + b \leq -1$



$$\vec{w} \cdot \vec{n} + b \leq 1$$

PAGE No.	
DATE	

for +ve class $\rightarrow y = +1 \rightarrow y_i(\vec{w} \cdot \vec{n} + b) \geq 1$
 for -ve class $\rightarrow y = -1 \rightarrow y_i(\vec{w} \cdot \vec{n} + b) \leq -1$

for +ve class $y_i(\vec{w} \cdot \vec{n} + b) \geq 1$

$$(\rightarrow (\vec{w} \cdot \vec{n} + b) \geq 1) \times$$

$$\vec{w} \cdot \vec{n} + b \geq 1$$

for -ve class $y_i(\vec{w} \cdot \vec{n} + b) \leq -1$

$$-1(\vec{w} \cdot \vec{n} + b) \leq 1$$

changing \leq to \geq

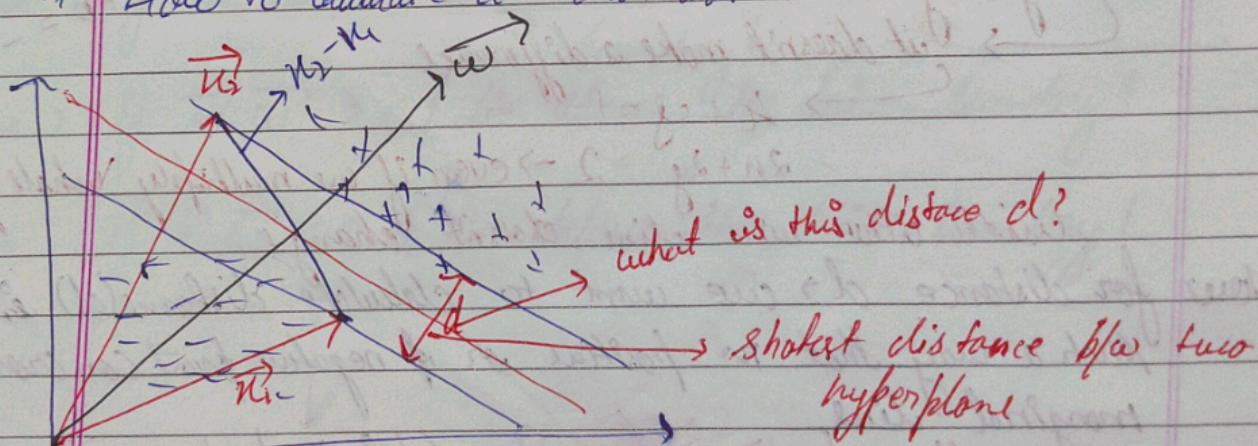
$$1(\vec{w} \cdot \vec{n} + b) \geq -1$$

Both of eqⁿ due to y taking (+1, -1)
 eqⁿ can be re-written as:

$$y_i(\vec{w} \cdot \vec{n} + b) \geq 1$$

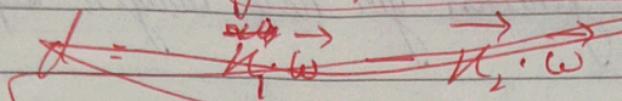
\hookrightarrow this will take care both class

+ How to calculate d ? and what is w ?



Now we will calculate d with help of \vec{w}

~~$d = \text{project of } \vec{w_2} \text{ on } \vec{w} - \text{projection of } \vec{w_1} \text{ on } \vec{w}$~~



$$d = \vec{w} \cdot (\vec{w}_2 - \vec{w}_1)$$

To get shortest distance we ~~will~~ need a unit vector $\frac{\vec{w}}{|\vec{w}|}$ to all marginal plane/hyperplane and that unit vector is \vec{w} .

$$d = \vec{w}_2 - \vec{w}_1 \cdot \frac{\vec{w}}{|\vec{w}|}$$

$|\vec{w}| \leq 1$ since \vec{w} is unit vector

$$|\vec{w}| = 1$$

$$d = \vec{w}_2 - \vec{w}_1 \cdot \frac{\vec{w}}{|\vec{w}|}$$

$$(3) \quad d = \frac{\vec{w}_2 \cdot \vec{w} - \vec{w}_1 \cdot \vec{w}}{|\vec{w}|} \quad \left\{ \begin{array}{l} \vec{w}_1, \vec{w}_2 \rightarrow \text{support} \\ \text{vector} \end{array} \right.$$

so they lie on marginal hyperplane.

for vector $\vec{w}_1 \quad y_1 = 1$

$w^T \vec{w}_1 = 1 - b \rightarrow ①$

Since they are support vector they should follow $y_1(\vec{w} \cdot \vec{w}_1 + b) = 1$

for vector $\vec{w}_2 \quad y_2 = -1$

$w^T \vec{w}_2 = -b - 1 \rightarrow ②$

Putting all in eq (3)

This is $CF = \frac{2}{|\vec{w}|}$

$$d = \frac{-b - 1 - 1 + b}{|\vec{w}|} = \frac{2}{|\vec{w}|}$$

Cost function $CF = \frac{1}{2} \|w\|^2$ by changing (w, b)

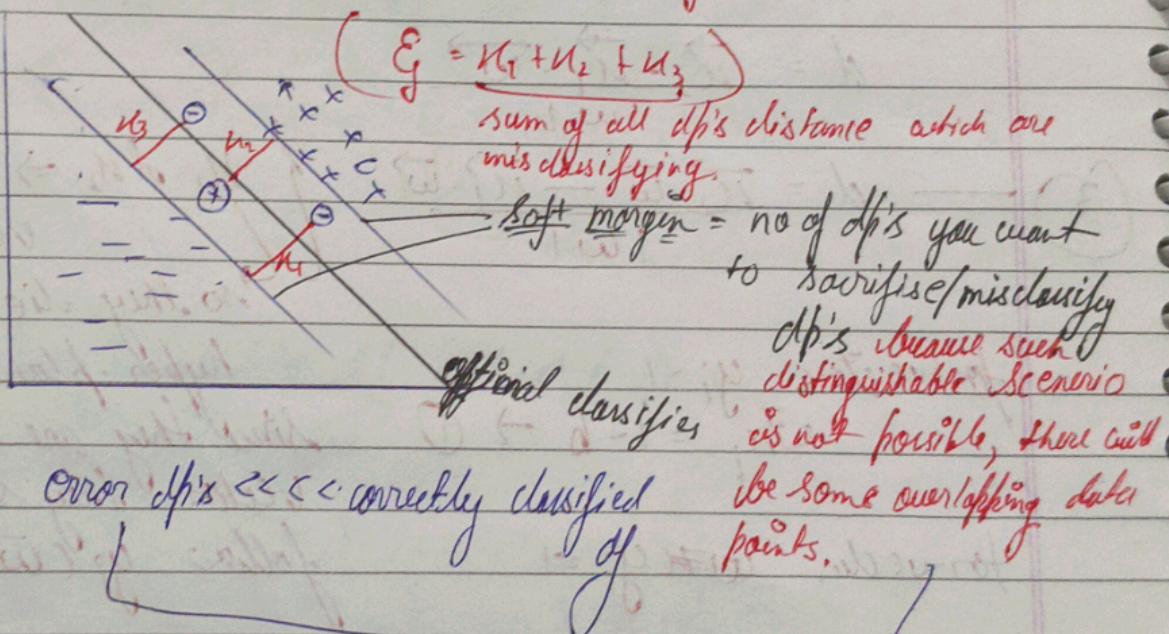
* modified cost function for hard margin support vector classifier

minimise $\frac{1}{2} \|w\|^2$ by varying w & b

Constraints such that $y_i(w^T x_i + b) \geq 1$

$\max f(u) = \min \frac{1}{2} \|u\|^2$ max converts to minimise
in case of inverse

$CF = \min_{w, b} \frac{1}{2} \|w\|^2$ constraint $y_i(w^T x_i + b) \geq 1$



due to this our cost function changes

* Soft margin SVC cost function: \rightarrow zeta

$$\text{minimise } \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n E_i \quad \text{such that } y_i(w^T x_i + b) \geq 1$$

$C = \text{no of misclassified dp's}$
(Hyperparameter)

$E(zeta) \rightarrow$ is the distance of all misclassified dp's to correct margin planes.

* higher the ' δ' \rightarrow distance b/w hyperplane of
 ↓
 two classes

Cover the error \rightarrow in train set, sacrifice some dp's but
 while testing it performs good

$$CF = \frac{||w||^2}{2} + C \sum_{i=1}^n \xi_i$$

also called as hinge loss
such that

soft margin Support vector Classifier CF $y_i(w^T n + b) \geq 1 - \xi_i$

Sum errors

$$3 + 3 + 3 + 3 = 12$$

$$\text{margin} = 3 - 1 = 2$$

$$\frac{3+3+3+3}{4} = \frac{12}{4} = 3$$

$$3 \geq 1 + \xi_i$$

hand p.c.

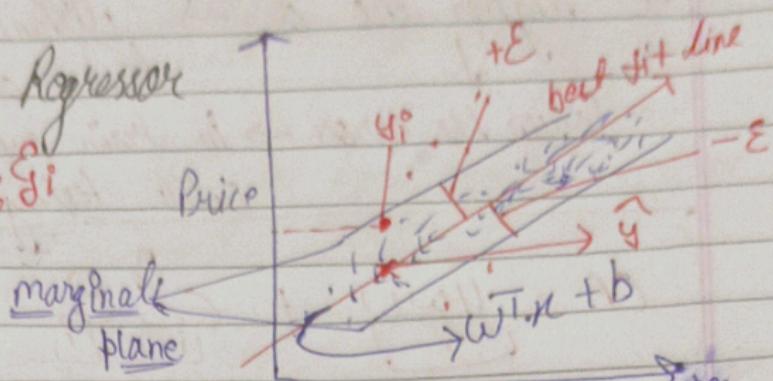
Maximizing the margin

$$(3 + 3) \cdot 3 + 10 \cdot 1 = 32$$

$$(3 + 3) \cdot 3 + 10 \cdot 1 = 32$$

Support Vector Regressor

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



we want all our datapoint surrounding to best line.

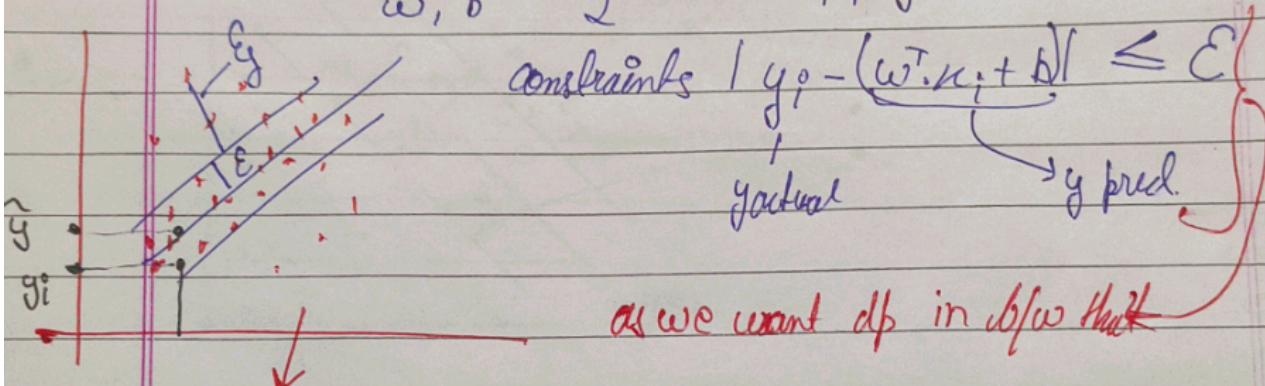
All datapoint should be on the boundary of marginal plane

all dp's should be in between

$$w^T x + b - \epsilon \quad \text{---} \quad w^T x + b + \epsilon$$

$\epsilon = \text{tolerance}$

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



But few of dp's will be away from marginal line

So we want ξ_i & ϵ to be minimum.
as we want all dp's to be close to best fit line.

$$\text{CF} = \min_{w, b} \frac{\|w\|}{2} + C \sum_{i=1}^n (\xi_i + \epsilon)$$

Constraint $|y_i - w^T x_i| \leq \epsilon + \xi_i$ distance b/w nearest
marginal plane and datapoint

a. → change dimension to higher dimension by mathematical transformation to distinguish b/w classes

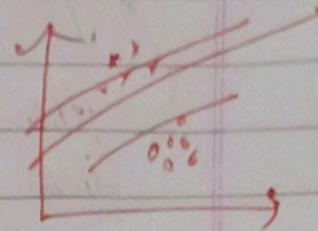
PAGE NO.	
DATE	

Now ($\S 18$) → ~~for~~ Hyperparameters do we want

SVM Kernel :-

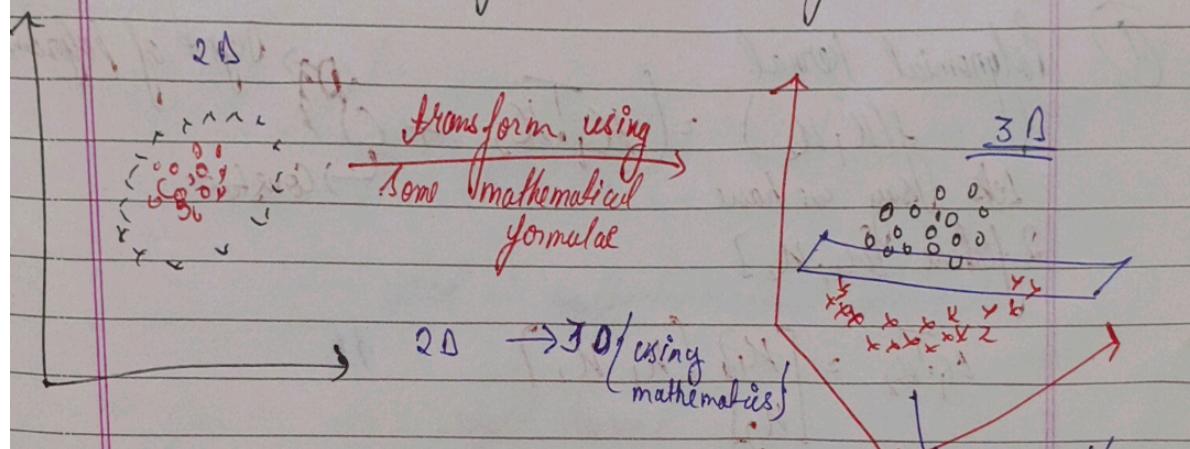
SVM-Kernel-trick tell now our data is like

what if it looks like \rightarrow non-linear?
if data is not linear

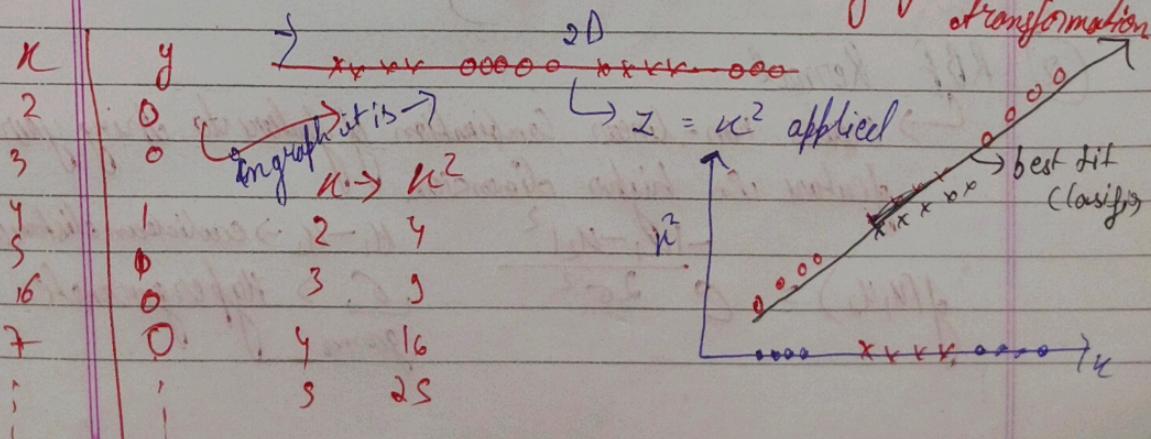


→ if we try to solve it \rightarrow 50% of data will be misclassified

→ for that we have got SVM kernel trick



$f(u) \Rightarrow$ Kernel \Rightarrow kernel transformation
there can be some overlapping after kernel transformation



why called as kernel trick

↪ lets say we are not sending data in higher dimension
but just using mathematical transformation to achieve it

↪ then using SVC classifier you can classify non linear data

Types of kernel function →

- ① Polynomial
- ② RBF (Radial Basis Function)
- ③ Sigmoid

① Polynomial kernel

$$f(u_1, u_2) = (u_1^T u_2 + c)^p \quad \begin{matrix} \text{Degree of polynomial} \\ \text{Constant} \end{matrix}$$

lets say we have
2 feature $[u_1, u_2]$

$$u_1^T u_2 = \begin{bmatrix} u_1^T \\ u_2^T \end{bmatrix} \begin{bmatrix} u_1 & u_2 \end{bmatrix}$$

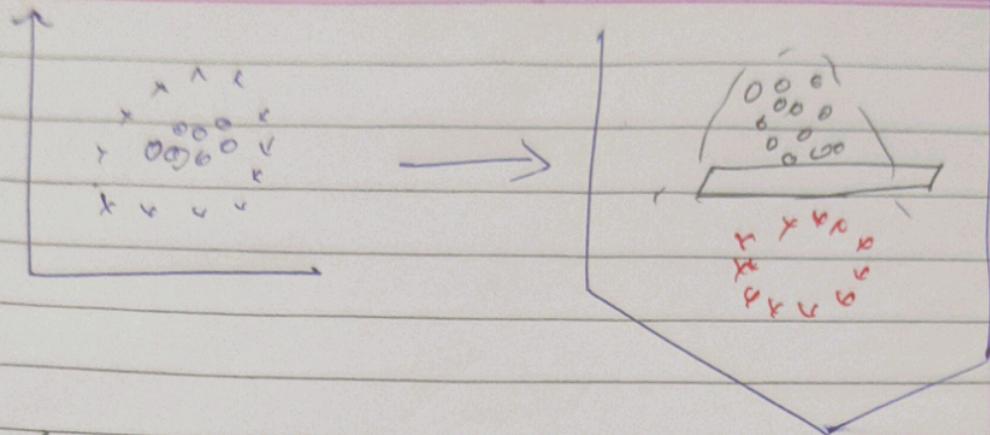
$$= \begin{bmatrix} u_1^2 & u_1 u_2 \\ u_1 u_2 & u_2^2 \end{bmatrix} \quad u_1, u_2 \rightarrow u_1 u_2, u_1^2, u_2^2$$

② RBF Kernel

↪ creates non-linear combination of features to bring your features in higher dimension

$$f(u_1, u_2) = C \frac{-|u_1 - u_2|^2}{2\sigma^2}$$

$u_1 - u_2 \rightarrow$ euclidean distance
 $\sigma \rightarrow$ Hypergeometric gamma?



(3) Sigmoid Kernel

$$f(u, g) = \frac{1}{1 + e^{-u}}$$

* Bessel kernel \Rightarrow

$$f(u, g) = \frac{J_{\nu+1}(\sigma \|u - g\|)}{\|u - g\|^{\nu+1}}$$

* Anova Kernel

$$f(u, g) = \sum_{k=1}^n \exp(-\sigma (u^k - g^k)^2)$$

→ used in multidimensional regression.

* How to choose right kernel

→ build on all

→ By Hyperparameters
→ get best one