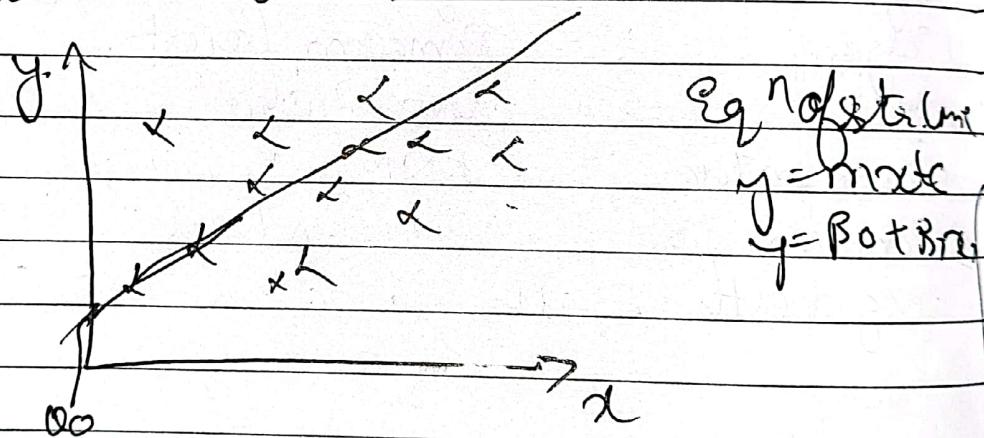


→ Independent
↓ dependent

Simple Linear Regression.

Find the best fit line ~~so~~ in such a way that the difference between the errors and after doing summation should be minimum.

whichever is the best fit line and normal distance is there we can say that model is trained well.



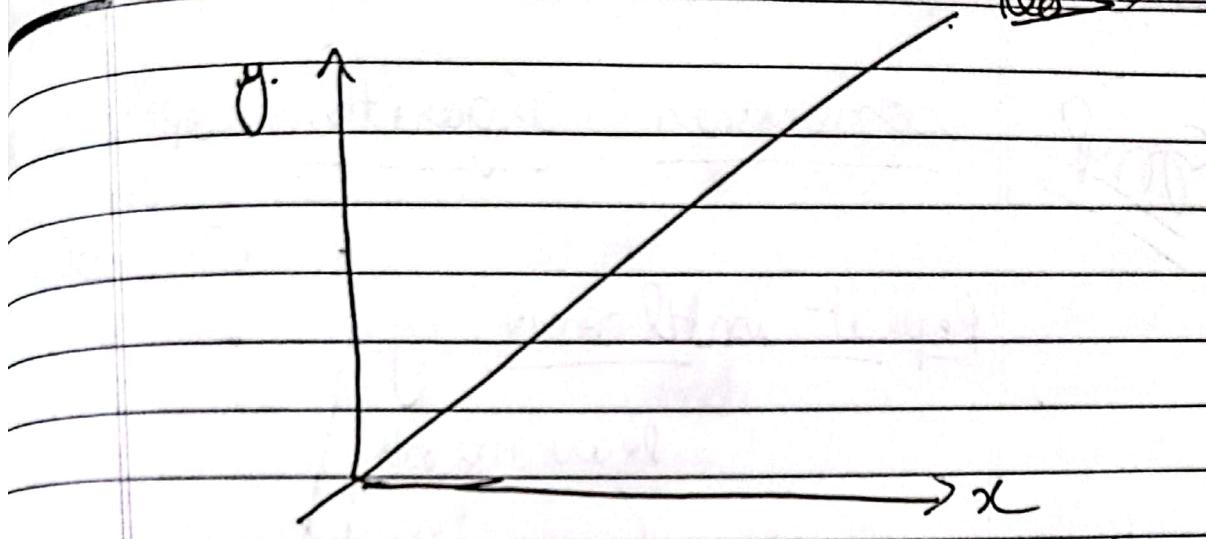
$$\boxed{h_0(x) = \theta_0 + \theta_1 x}$$

$\theta_0 \rightarrow$ Intercept
 $\theta_1 \rightarrow$ Slope

Intercept ≥ 0

Page.....

$$\boxed{\alpha_0 = 0}$$



Cost function

predict actual
↓ ↓

JMF
$$J(\alpha_0, \alpha_1) = \frac{1}{m} \sum_{i=1}^m (h_\alpha(x^{(i)}) - y^{(i)})^2$$

↓
mean squared error.

$m \Rightarrow$ all the datapoints

final aim

$$\text{minimize } J(\alpha_0, \alpha_1) \Rightarrow \frac{1}{2m} \sum_{i=1}^m (h_\alpha(x^{(i)}) - y^{(i)})^2$$

$\frac{1}{2}$

~~QPF~~ conversion algorithm. [optimize the changes of α_j value]

Repeat until conv.

$$\alpha_j' = \alpha_j - \alpha \frac{\partial J(\alpha_j)}{\partial \alpha_j}$$

learning rate



$$\alpha_j' = \alpha_j - \alpha (-\nabla)$$

$$= \alpha_j + \alpha$$

learning rate decides the speed
of convergence

Q1 what is linear regression?

Ans It is a statistical method that is used for predictive analysis. Linear reg alg shows a linear relations between a dependent (y) and one or more dependent (y) variables hence called as linear regression.

Q2 how we can calculate error in linear regression?

Ans mean squared error
error b/w predicted value and real error

Q3 diff b/w loss and cost function?

Ans loss is for one

Linear Reg Alg.

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_m x_m$$

~~Def~~convergence algorithm

$$J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

mean squared error

cost function

loss function vs

cost function

In cost function we have to
check every observation.

In loss function

$$= (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

↓

$$(y^{(i)} - \hat{y}^{(i)})^2$$

predicted value

actual value

$$\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0} = \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)}) - y^{(i)}$$

$$= \frac{1}{2m} \sum_{i=1}^m [(\theta_0 + \theta_1 x^{(i)}) - y^{(i)}]$$

$$\Rightarrow \frac{1}{2m} \sum_{i=1}^m [(1 + x^{(i)}) - y^{(i)}]^2$$

=

$$j=1. \quad \frac{\partial}{\partial \theta_1} \left[\sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)}) - y^{(i)} \right]^2$$

$$= 2 \sum_{i=1}^m [(\theta_0 + \theta_1 x^{(i)}) - y^{(i)}] x$$

cost functions.

① MSE = mean \times y error.

② RMSE = $\sqrt{\text{MSE}}$ "absolute"

③. MAE = mean \times ~~absolute~~ error

MSE

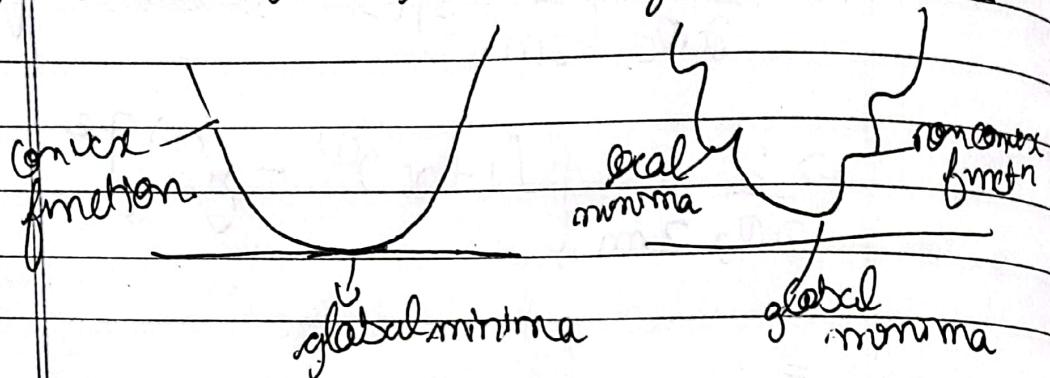
$$\sum_{i=1}^n (y - \hat{y})^2$$

Date..... /..... /.....

Page.....

① Advantages

- (i) This y^n is differentiable
- (ii) " " also has one global minima



Disadvantage

- (i) This is not robust to outliers
- (ii) Penalizing the error

② MAE [mean absolute error]

$$= \frac{1}{n} \sum_{i=1}^n |y - \hat{y}|$$

A advantage

- (i) Robust to outliers
- (ii) It will also lie in the same envt.

Disadvantage

- 1 convergence takes more time
- 2 optimization is a complex task.

H-W

Huber Loss

RMSE

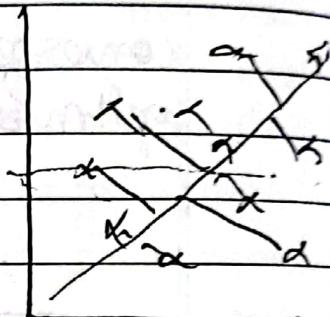
Performance Matrix

- ① R squared
- ② Adjusted R squared.

~~Graph~~

R squared

$$R^2 = 1 - \frac{SS_{Res}}{SS_{Total}}$$



SSRes - Sum of square residual

SSTotal - Sum of square average.

$$= 1 - \frac{\sum_{i=1}^n (y_i - \bar{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

$$\text{avg of } \bar{y}$$

$$R^2_{\text{adjusted}} = 1 - \frac{\{\text{small number}\}}{\{\text{bigger number}\}} \rightarrow \text{small no.}$$

R squared is used to measure the performance of the model that you have actually created.

ADJUSTED R SQUARED

(overfitting & underfitting)

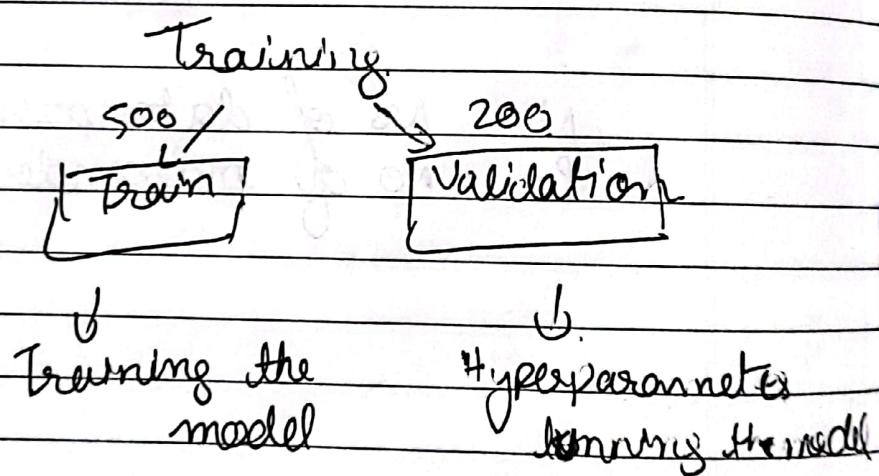
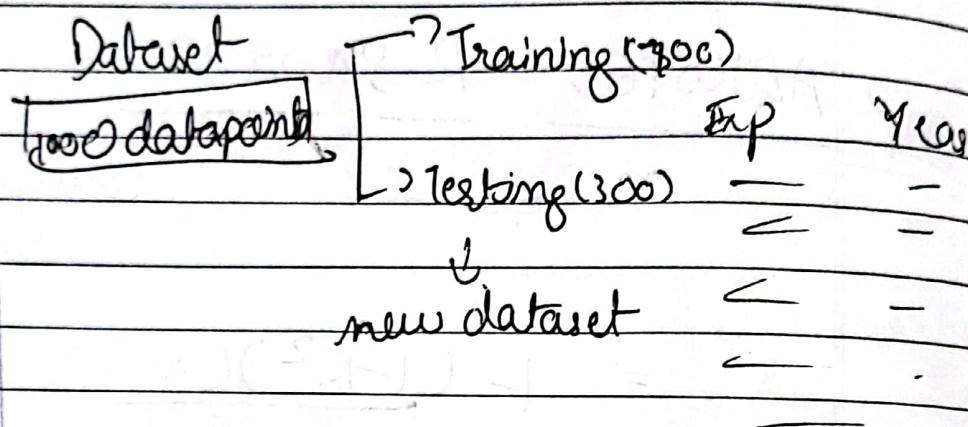
$$= \frac{1 - ((1-R^2)(N-1))}{N-P-1}$$

N = no of data points

P = no of independent features

most important

Overshooting and Underfitting
(Bias Variance)



MODEL

TRAIN DATA	Very good Accuracy (90%) [Bad]
TEST DATA	" " " (85%) [Var.]
	↳ generalised model

Train	Very good accuracy [90%] (low bias)
Test	Bad Accuracy [50%] [high var.]

Test Bad Accuracy [50%] [High var.]

16 10/21/2018

overfitting

hyperparameter tuning

model laws. (Hugh Bibas)

accuracy low / high (low or high bias)

11

underfitting.

~~Qmp~~

RIDGE REGRESSION , LASSO " " ELASTIC NET. "

① Ridge regression. (L2 regularization)

↓ aim

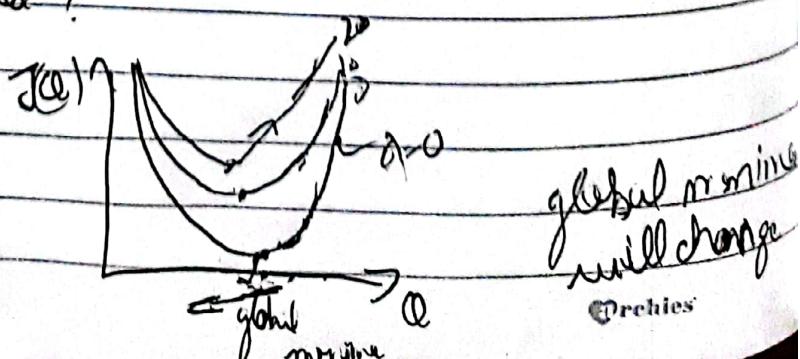
[reduce overfitting]

we will use same cost functn

$$= \frac{1}{n} \sum_{i=1}^m (\hat{y}_i - y_i)^2 + \lambda (\text{slope})^2$$

λ = hyperparameter.

→ ~~Qmp~~ what is the relationship b/w slope & lambda?



λ is inversely proportional
to slope

when using ridge regression it
will never overfit any line.

In ridge regression λ will never
become zero. because when
 $\lambda = 0$ you will see one
feature gets deleted and this
will effect our model.

LASSO REGRESSION {L1 norm} {L1 regularization}

(i) aim

to reduce the feature

(j)
feature selection

$$\text{Cost function} = \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2 + C \sum_{i=1}^m |\beta_i|$$

combination of ridge & loss

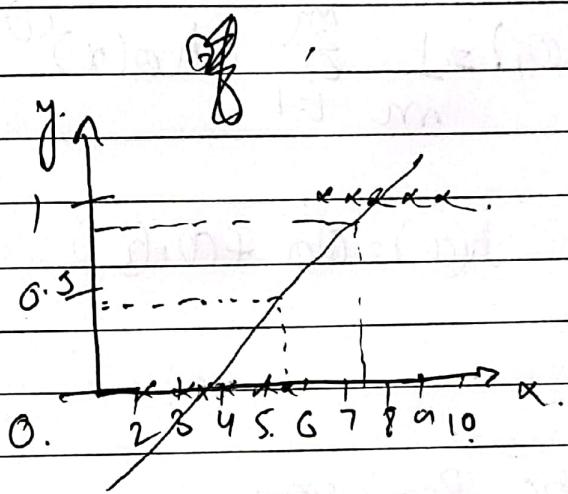
ELASTIC NET [L1 & L2 norm]

$$\text{cost fn} = \frac{1}{m} \sum_{i=1}^m (h_0 x^{(i)} - y^{(i)}) + \lambda_1 \sum_{i=1}^m |x_i| + \lambda_2 \sum_{i=1}^m |\text{slope}|$$

LOGISTIC REGRESSION

→ It is used to solve classification problem.

with the help of regression we try to make best fit line



$$\text{if } y \leq 0.5 = 0 \\ y > 0.5 = 1$$

It means 0.5 is my threshold

this best fit line $f(x) = \theta_0 + \theta_1 x$

sigmoid activation

output = 0 to 1

② Sigmoid fn = $\frac{1}{1+e^{-x}} \Rightarrow 0 \text{ to } 1$

Linear regression cost function.

(C) $J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$

gradient descent
convergence

$$h(\theta) = \theta_0 + \theta_1 x$$

Logistic Regression

- ① create a best fit line
- ② Squashing. \rightarrow Sigmoid fn.

$$J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$h_\theta(x) = \frac{1}{1+e^{-(\theta_0 + \theta_1 x)}}$$

Sigmoid. activation.

$$\text{let } z = \theta_0 + \theta_1 x$$

$$= \sigma(z)$$

$$= \frac{1}{1+e^{-z}}$$

$$h_0(x) = \frac{1}{1+e^{-(\alpha_0 + \alpha_1 x)}}$$

* log loss cost function.

$$\text{cost}(h_0(x)^{(i)}, y^{(i)}) = \begin{cases} -\log(h_0(x)) & \text{if } y=1 \\ -\log(1-h_0(x)) & \text{if } y=0 \end{cases}$$

y = truth value

$$\text{cost}(h_0(x)^{(i)}, y^{(i)}) =$$

$$-y \log(h_0(x)) - (1-y) \log(1-h_0(x))$$

↓

This will create a convex
function

↓

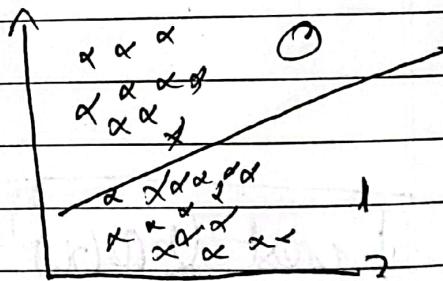
never local minima.

minimize cost function $J(\theta_0, \theta_1)$ by changing θ_0, θ_1 convergence alg.

Repeat convergence

$$\left\{ \begin{array}{l} j=0 \text{ and } \\ \theta_j := \theta_{0j} - \alpha_2 J(\theta_0, \theta_1) \\ \theta_{0j} \end{array} \right.$$

Performance metrics



- ① confusion matrix
- ② Accuracy
- ③ precision
- ④ Recall
- ⑤ f-Beta score

DIRECT	BY	OP	+
f1	f2	OP	+
-	-	0	1
-	-	1	1
-	-	0	0
-	-	1	1
-	-	1	1
-	-	0	1

Confusion matrix

		Actual value	
		1	0
predicted	1	3	2
	0	1	1
predicted	1	TP	FP
	0	FN	TN

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

$$= \frac{3+1}{3+2+1+1} = \frac{4}{7}$$

$\approx 57\%$.

* DATASET \rightarrow BINARY CLASSIFICATION

\hookrightarrow 1000 datapoints. \hookrightarrow 900 \rightarrow ?
 \hookrightarrow 100 \rightarrow 0 } Imbalanced dataset.

* Precision : $\frac{TP}{TP+FP}$

out of all the actual values how many are pred correctly predicted

Act.	
1	0
TP	FP
FN	TN

* RECALL $\frac{TP}{TP+FN}$ \rightarrow out of all the predicted values how many are correctly predicted

* F-Beta Score : $\frac{(1+\beta)^2}{(\beta^2 \text{ precision} + \text{recall})}$ Precision & Recall

① If FP and FN are both important.
 $\beta = 1$

$$\text{F1 score} = 2 \frac{P \times R}{P+R}$$

② If FP is more important than FN
 $\beta = 0.5$

$$\text{F.S. score.} = \frac{1+0.25}{(0.25)} \frac{P \times R}{P+R}$$

$$\frac{1+0.25}{(0.25)} \frac{P \times R}{P+R}$$

③ If FN \gg FP

$$\text{F2 score} = \frac{(1+4)(P \times R)}{4(P+R)}$$

SUPPORT VECTOR MACHINE (SVM)

vector classifier

SVC ① solve both classification and

SVR ② regression problem

vector reg.

$$\Rightarrow y = w \cdot x + c$$

$$y = \beta_0 + \beta_1 x$$

$$w \cdot x + b \cdot y + c = 0$$

$$y = \frac{-a}{b}x - \frac{c}{b}$$

coeff intercept

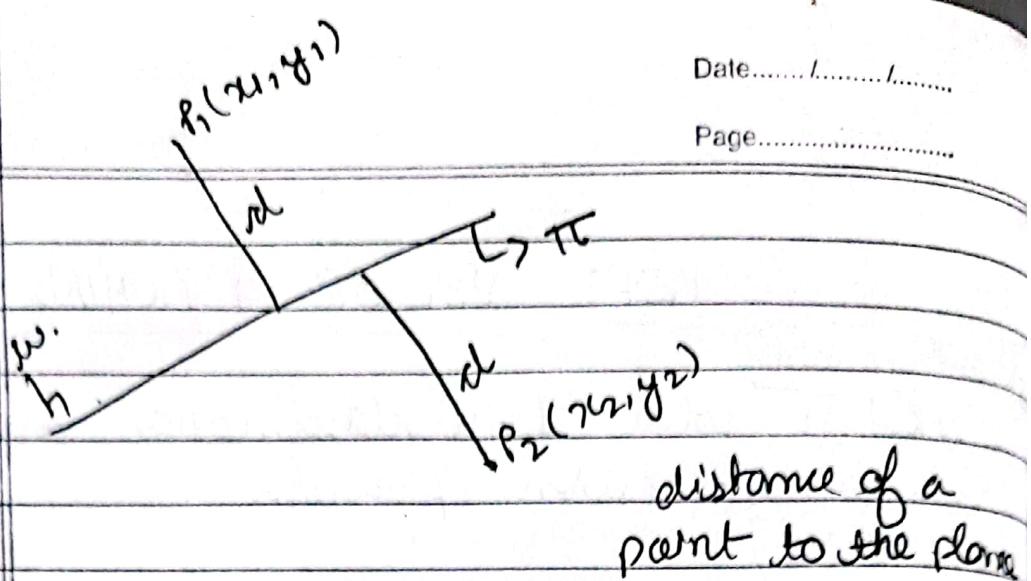
$$a_1x_1 + b_1x_2 + c = 0$$

$$w_1x_1 + w_2x_2 + b = 0$$

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

$$w^T x = 0$$

equation of line passing
through origin

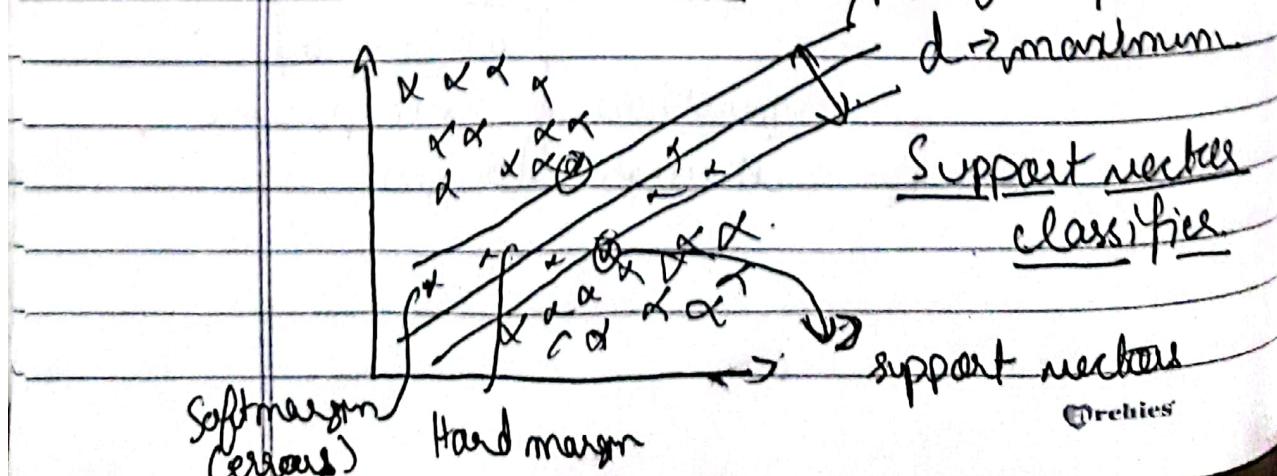
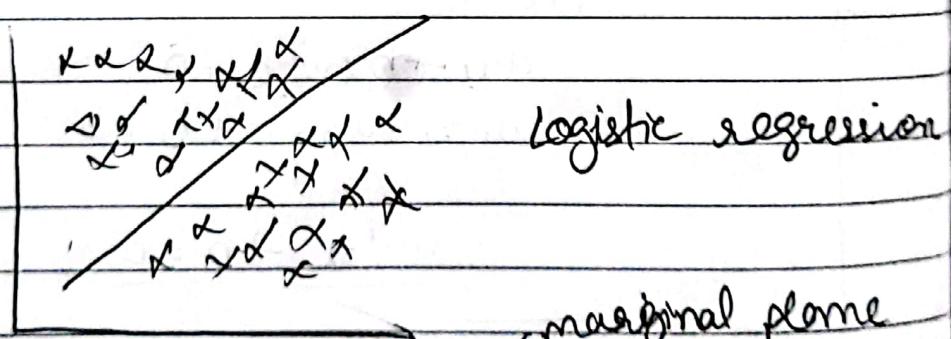


$$d = \frac{w^T p}{\|w\|} \Rightarrow \|w\| \|p\| \cos\theta$$

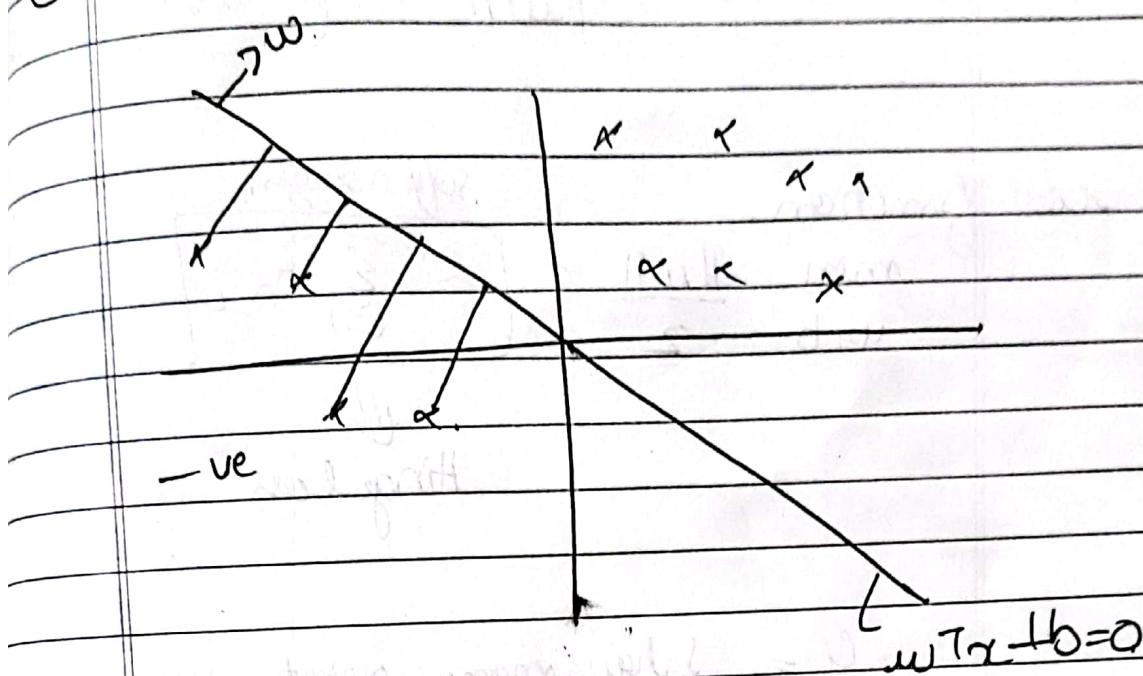
unit vector \Rightarrow vector which has a magnitude of 1

~~$$\frac{w^T p}{\|w\|} \Rightarrow \text{distance}$$~~

Geometric Intuition behind SVM



② SVM Mathematical Intuition



Cost function.

$$\text{maximise}_{w,b} \frac{2}{\|w\|} \Rightarrow \text{distance b/w marginal planes}$$

$$\text{constraint } y_i \begin{cases} 1 & w^T x_i + b \geq 1 \\ -1 & w^T x_i + b \leq -1 \end{cases}$$

for all current point

$$\text{constraints} \rightarrow y_i * (w^T x_i + b) \geq 1$$

maximize $\frac{2}{\|w\|} \Rightarrow \boxed{\text{minimize } \frac{\|w\|}{2}}$

cost function.

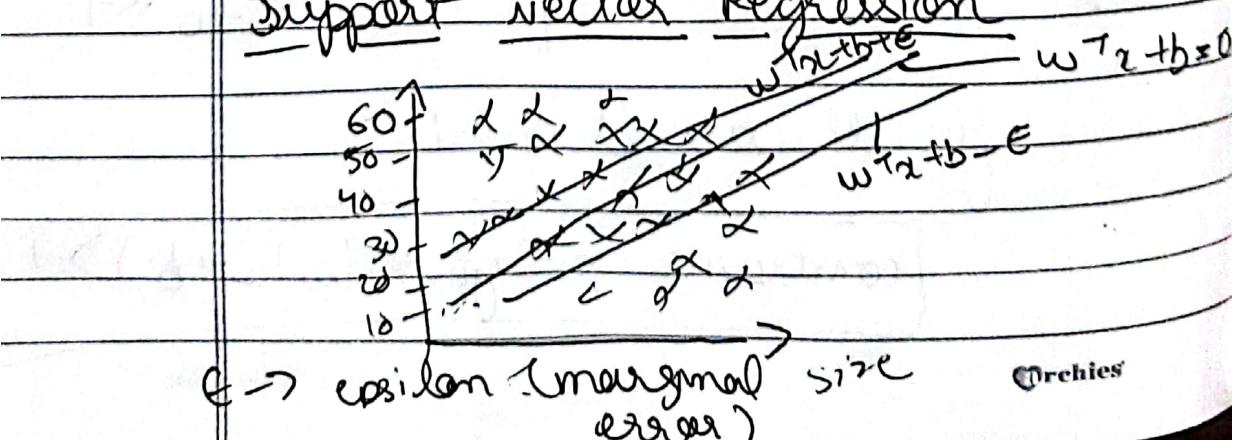
$$\min_{w,b} \frac{\|w\|}{2} + \boxed{c_i \sum_{i=1}^m \xi_i}$$

\Downarrow
Hinge loss

$c_i = \{ \text{how many points are even}\}$
 $\text{ignore for misclassification?}$

$\left\{ \begin{array}{l} \text{data} \\ \rightarrow \text{submission of the distance} \\ \text{of the incorrect data points} \\ \text{from S points} \end{array} \right.$

Support vector Regression



constraint

$$|y_i - w_i x_i| \leq \epsilon + \xi_i$$

~~non-convex function~~

→ Will SVM impacted by outliers?
Yes.

→ SVM kernel