



# Linear and Logistic Regression HandWritten Notes

- Darshan R M

"Learn, Share,  
and  
Collaborate"

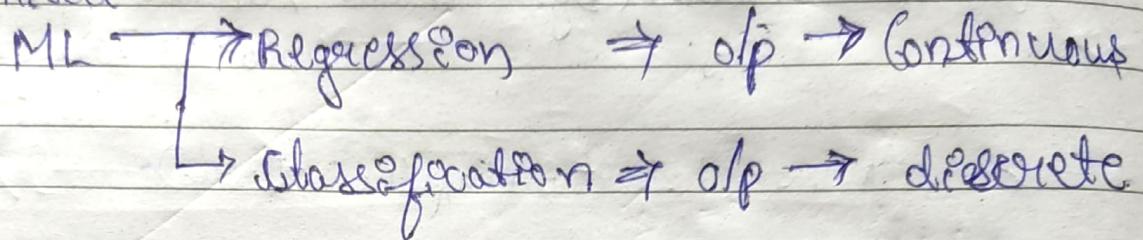
15

## Week -15

### Linear and logistic regression

\* 1

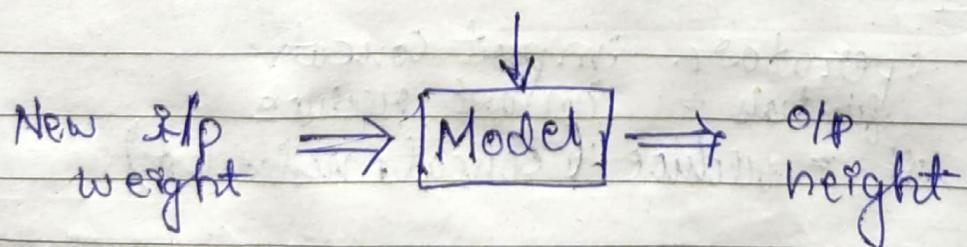
Supervised



### Simple linear regression

Dataset's <del>source</del>	(Independent)		(Continuous) o/p (dependent)
	o/p weight	height	
	74	170	
	80	180	
	75	175.5	
	-	-	
	=	=	

Train



\* Simple linear regression

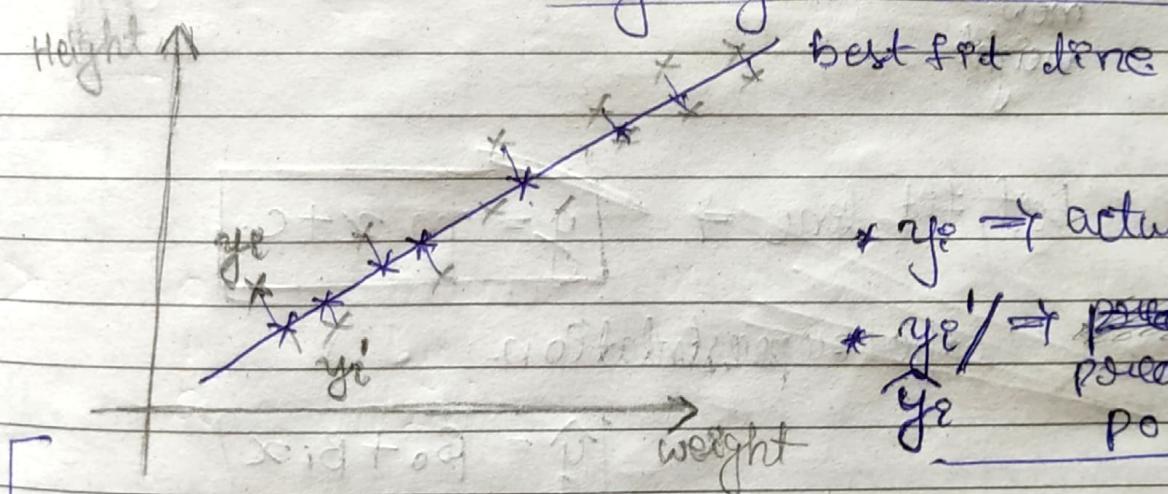
↳ we have 1 input and 1 output

\* Multiple linear regression

↳ we have multiple input and 1 output

\* Geometrical intuition

↳ gives main aim what we are going to do and how.



→ we ~~will~~ do find out best fit line.

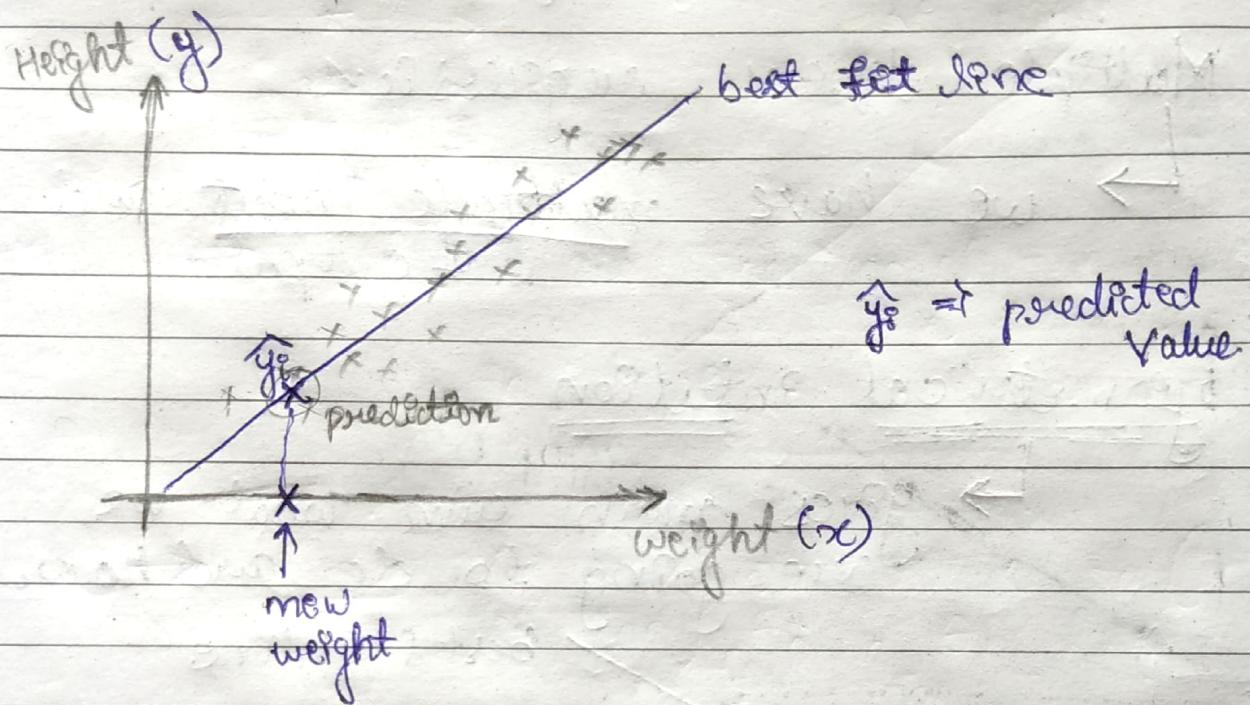
~~such that the summation of the errors is least.~~

↓  
between actual and predicted point

$$\sum \text{Error} \rightarrow \text{minimally}$$

\* why to find best fit line.

Because whenever the new weight comes we have to predict height.



best fit line  $\Rightarrow$   $\boxed{\hat{y} = m x + c}$ .

other representation.

$$\hat{y} = \beta_0 + \beta_1 x$$

substituting  $\rightarrow$  
$$h_0(x) = \theta_0 + \theta_1 x$$

because  $m$  &  $c$  also treated  
as hypothesis  
testing

Error  $\{ y_i - \hat{y}_i \}$  actual - predicted

\*

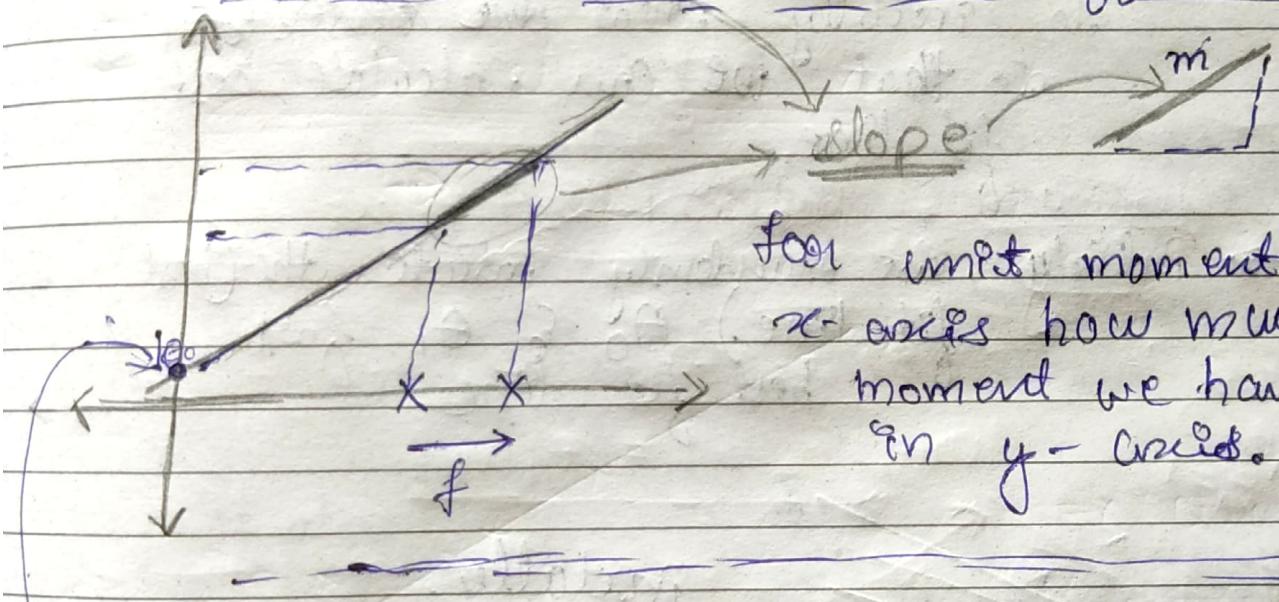
$$h_{\theta}(x) = \theta_0 + \theta_1 \cdot x$$

where,

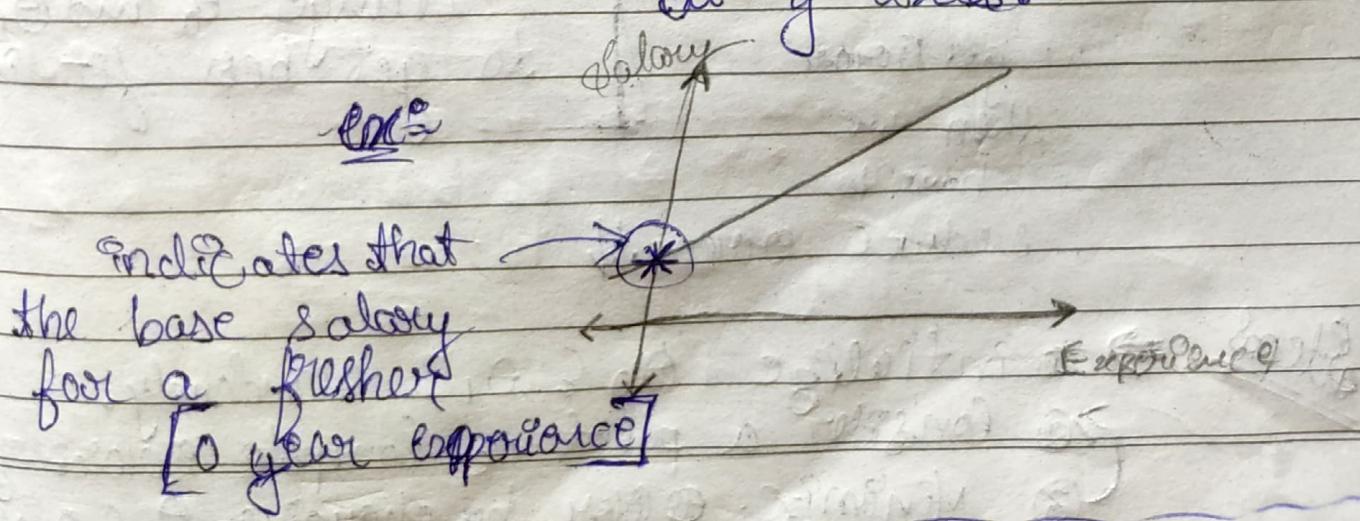
$x_i$  = data point

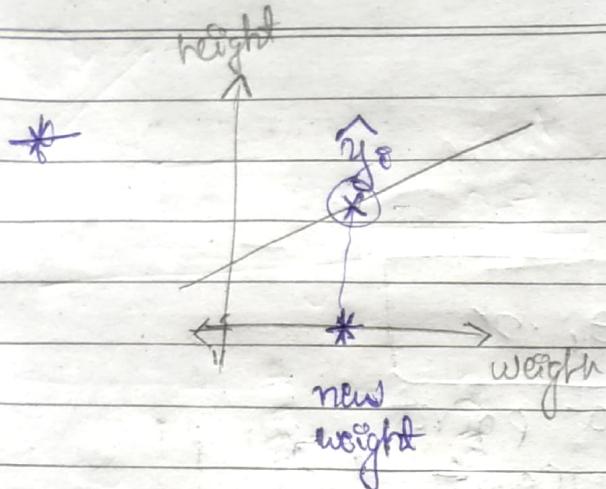
$\theta_0$  = Intercept. (c)

$\theta_1$  = Slope (or) Coefficient (m)



$\theta_0 \Rightarrow$  intercept  $\Rightarrow$  when  $x=0$  where do our line intercept at y-axis.

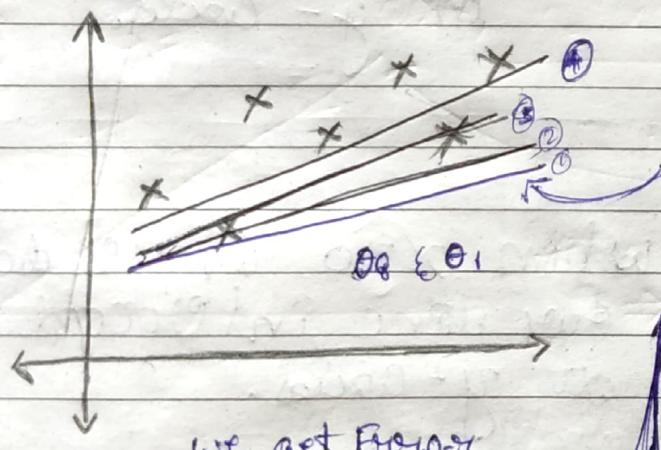




For trained model  
for any new weight  
it is calculated  
like this

\* But, to train the model we have to initially treat the  ~~$\theta_0, \theta_1$~~ , so that we can calculate best fit line.

→ Initially, we randomly indicate the (best fit line)  $\theta_0 \in \theta_1$  for best fit line



we get Error  
then we  
have to  
reduce error

Initially  
do now our aim  
is to continuously  
change  $\theta_0$  and  $\theta_1$   
to get best fit line.

Step

- ① Initialize  $\theta_0 \in \theta_1$  randomly.
- ② Consider a point & Compute error.
- ③ Minimize error by changing  $\theta_0 \in \theta_1$ .
- ④ Again

## 2 important things about Optimization

### ① Cost Function [Mean squared error]:

→ To find the sum summation of error we use cost function.

reduce ↓↓

$$J(\theta_0, \theta_1) = \frac{1}{n} \sum_{i=1}^n [y_i - h_{\theta}(x_i)]^2$$

actual - predicted

because we  
change values  
of  $\theta_0$  and  $\theta_1$ .

mean square  
error

\* we need to "minimize":  $J(\theta_0, \theta_1)$   
[cost function] to get the best fit line

where,  $n$  = number of data points

$y_i$  = actual value

$[y_i]$   $h_{\theta}(x)$  = predicted value

Final aim: To get best fit line.

by reducing value of Cost function

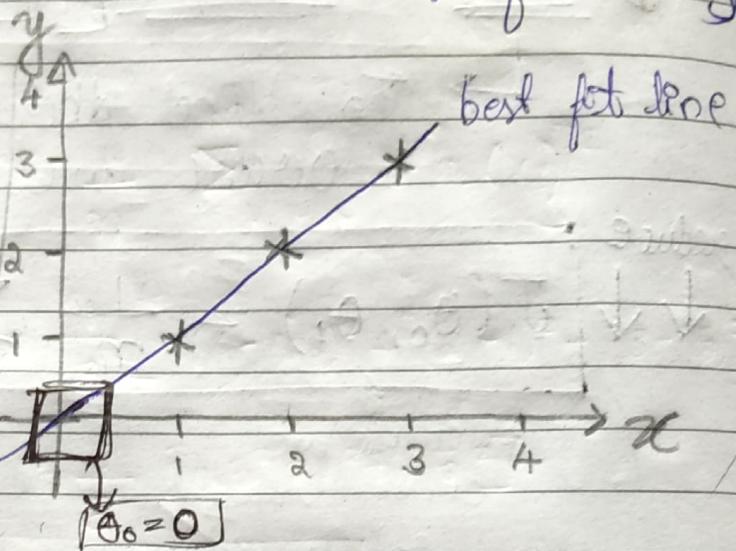
→ This entire process is  
an "optimization  
process."

## ~~Given~~ Dataset

$x$	$y$	$h_{\theta}(x)$
1	1	1
2	2	2
3	3	3

(independent) (dependent)

Optimization { Minimize the cost function }



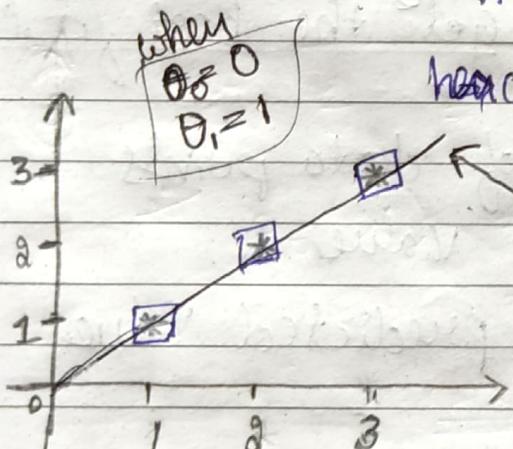
$\therefore$  as it touch origin

①

Equation of line :  $h_{\theta}(x) = \theta_0 + \theta_1 x_i$

As,  $\theta_0 = 0$

hence,  $h_{\theta}(x) = \theta_1 x_i$



let  $\theta_1 = 1$

\* for  $x=1$

$$\Rightarrow h_{\theta}(x) = 0 + 1(1)$$

$$\underline{h_{\theta}(x) = 1}$$

\* For  $x=2$

$$\Rightarrow h_{\theta}(x) = 0 + 1(2)$$

$$\underline{h_{\theta}(x) = 2}$$

\* For  $x=3$

$$\Rightarrow h_{\theta}(x) = 0 + 1(3)$$

$$\underline{h_{\theta}(x) = 3}$$

& let's calculate cost function

$$\theta_0 = 0 \Rightarrow [h_0(x) = \theta_0 x_0]$$

$\theta_1 = 1$

(cost Function)

$$J(\theta) = \frac{1}{n} \sum_{i=1}^m [y_i - h_\theta(x_i)]^2$$

$\underbrace{\qquad\qquad\qquad}_{\text{Error}}$

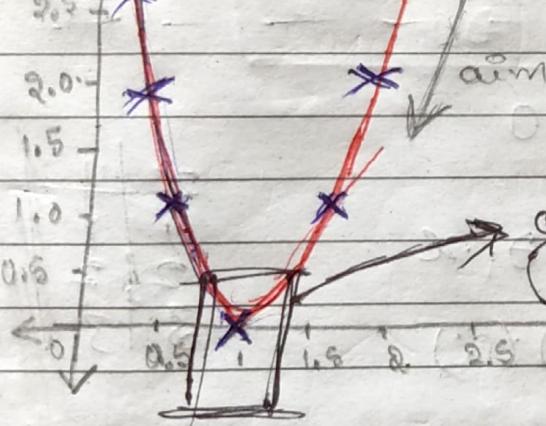
w.k.t  $m = 3$

$$= \frac{1}{3} [(1-1)^2 + (2-2)^2 + (3-3)^2]$$

$$J(\theta_1) = 0$$

[cost function]

Graph  
of  
cost  
function  
with  
 $\theta_1$



→ cost value  
with  
slope

$\theta_1$	$J(\theta_1)$
0	0
0.5	1.16
1	0.66

global minima  
[cost function ↓]  
(slope)

③ start  $\theta_1 = 0.5$

$$x_0 = 1, h_0(x_0) = 0.5(1)$$

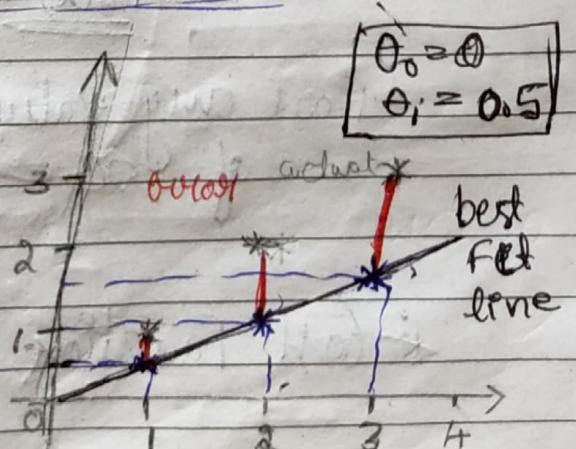
$$= 0.5$$

$$x_0 = 2, h_0(x_0) = 0.5(2)$$

$$= 1.0$$

$$x_0 = 3, h_0(x_0) = 0.5(3)$$

$$= 1.5$$



let's find error

$$\theta_0 = 0, \theta_1 = 0.5$$

Cost function,

$$[\theta_1 = 0.5]$$

~~Best fit~~

$$\Rightarrow [h_{\theta}(x) = 0.5(x)]$$

x	y	$h_{\theta}(x)$ $\theta_1 = 0.5$
1	1	0.5
2	2	1
3	3	1.5

$$J(\theta_1) = \frac{1}{n} \sum_{i=1}^n [y_i - h_{\theta}(x_i)]^2$$

$$= \frac{1}{3} [(1-0.5)^2 + (2-1)^2 + (3-1.5)^2]$$

$$= \frac{1}{3} [0.5^2 + 1^2 + 1.5^2]$$

$$[J(\theta_1) = 1.16]$$

$$\textcircled{3} \quad \text{let } \theta_1 = 0, \theta_0 = 0$$

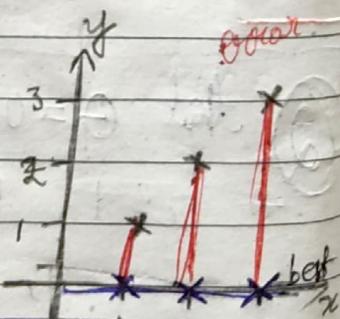
$$\begin{aligned} * h_{\theta}(x) &= \theta_0 + \theta_1 x \\ &= 0 + 0(x) \end{aligned}$$

$$\underline{h_{\theta}(x) = 0}$$

x	y	$h_{\theta}(x)$
1	1	0
2	2	0
3	3	0

For any value  
of  $x \in \{1, 2, 3\}$

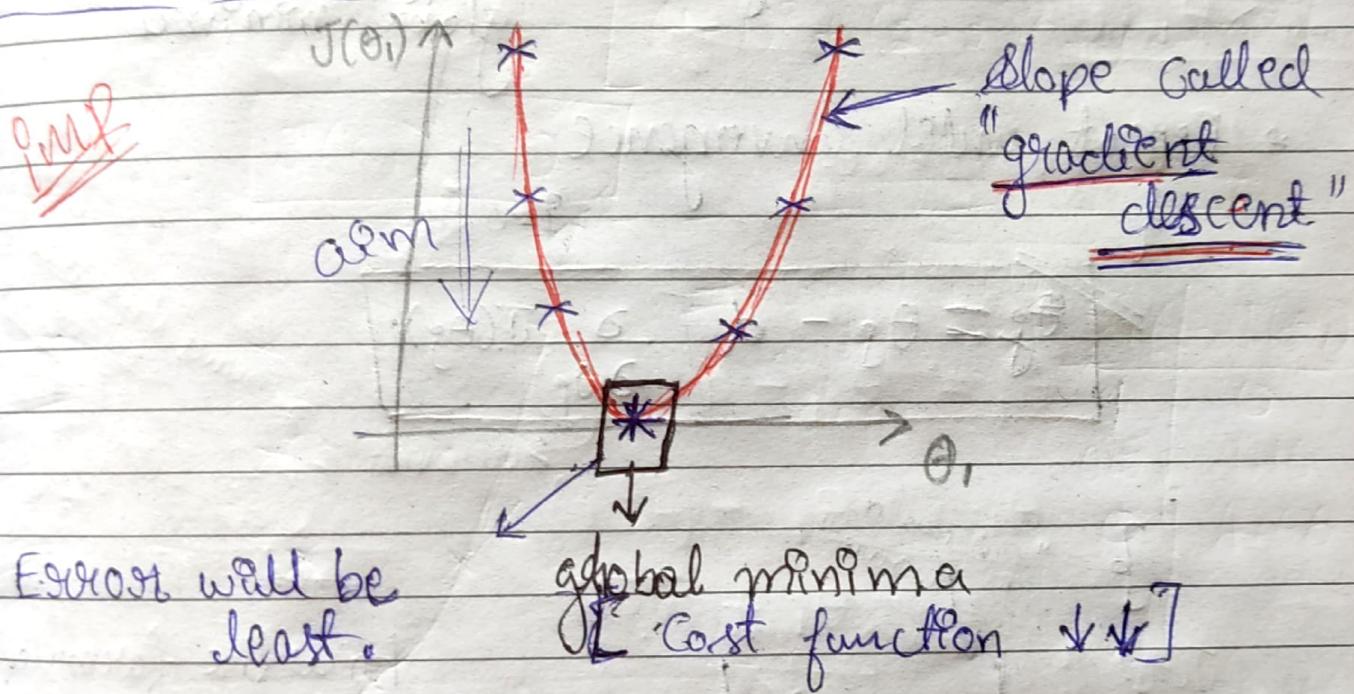
$$[h_{\theta}(x) = 0]$$



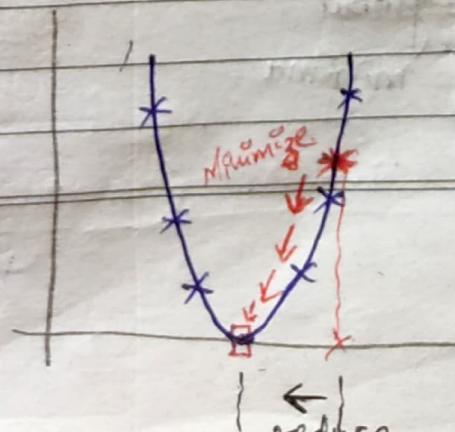
$$* \text{ Cost function, } J(\theta_0) = \frac{1}{m} \sum_{i=1}^m [y_i - h_{\theta}(x)]^2$$

$\approx 0$

$$\begin{aligned}
 J(\theta_0) &= \frac{1}{3} [(\theta_1 - 0)^2 + (\theta_2 - 0)^2 + (\theta_3 - 0)^2] \\
 &\geq \frac{1}{3} [1^2 + 2^2 + 3^2] \\
 &= \frac{14}{3} \\
 \boxed{J(\theta_0) \approx 4.66}
 \end{aligned}$$



\* we ~~can't~~ can't always initialize a ' $\theta$ ' value.  
so we have to find a technique  
to reduce error.



we should converge to reach global minima

## (2) Converge algorithm:

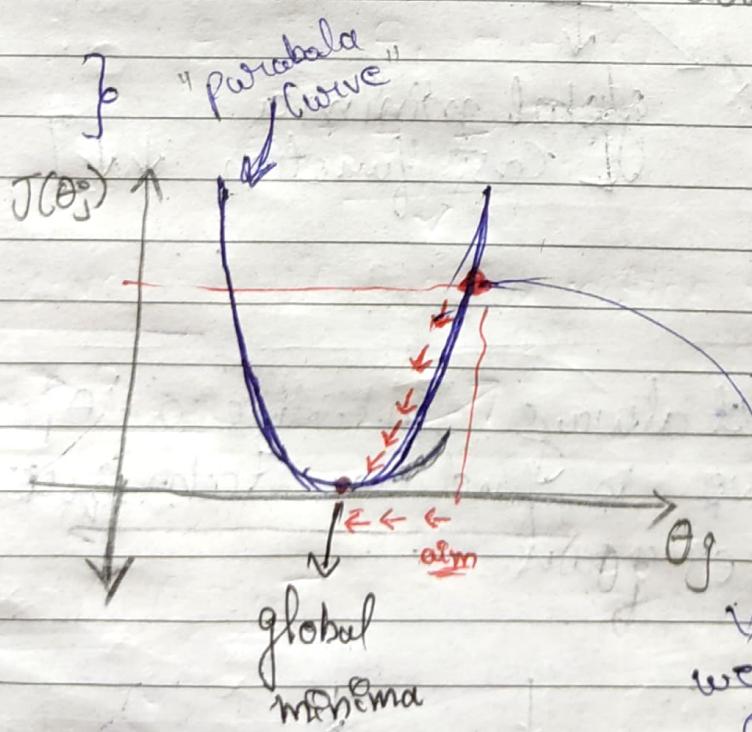
- we need to optimize the change of  $\theta_i$ ,
- we cannot always initialize  $\theta_i$ ,  
so we have to initialize once  
then try to converge towards  
global minima

\* Repeat until convergence

$$\theta_j = \theta_j - d \left[ \frac{\partial J(\theta_j)}{\partial \theta_j} \right]$$

means "slope at a specific point"

derivative  $\Rightarrow d \left[ J(\theta_j) \right] / d \theta_j$

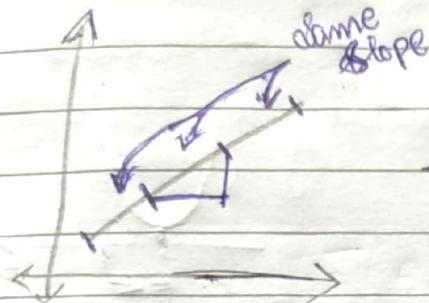


how the formula  
getting applied?

how to reduce  
 $\theta_j$ ?

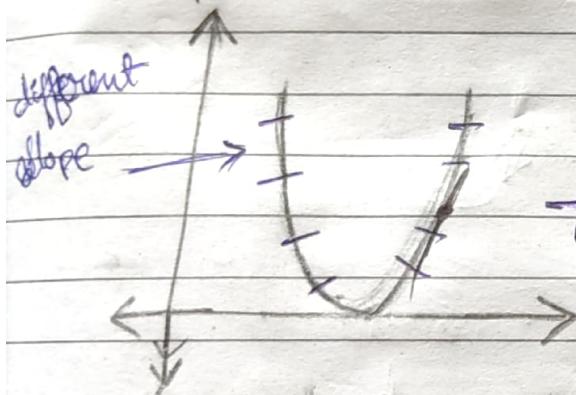
we need to  
calculate the  
derivative of slope  
to reduce ' $\theta_j$ '

\*  $\alpha$  = learning rate

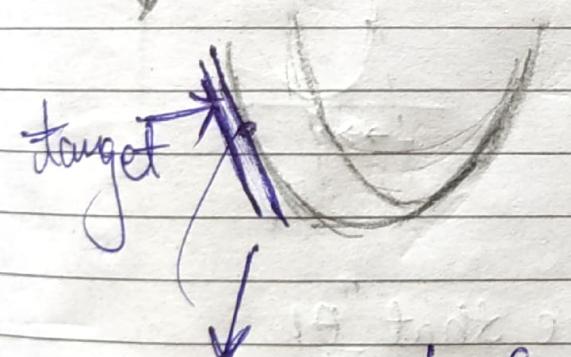


For a best fit line we can easily find out the slope of curve line

but



In parabola curve,  
we need derivative of slope of a point.



~~slope can be +ve/-ve~~

To find slope  
we have to calculate at point

we draw a tangent at the point

$$\frac{d}{d\theta_j} J(\theta_j)$$

slope at a point

direction is up

negative slope

positive slope

~~slope~~

$\cup$  (+ve) Slope

Now, if slope is +ve

$$\Theta_{\text{true}} = \Theta_{1,p} - d \quad (+ve)$$

$$\Theta_{\text{new}} = \Theta_{\text{old}} - (+ve)$$

when

"slope is +ve"

$$\theta_{\text{new}} < \theta_{\text{ratch}}$$

so, we can notice

that " $\theta$ , will reduce".

specie

10

\* if slope is -ve

(-ve)  
slope

$$\Theta_1 = \Theta_1 - d(-ve)$$

$$\theta_1 = \underline{\theta_1} - (-ve)$$

$$\theta_{\text{new}} = \theta_{\text{old}} + \text{value}$$

*when*

only  
"slope is -ve"

$$\theta_{\text{new}} \rightarrow \theta_{\text{old}}$$

so, we can notice that  $\theta$

will inoculate

# learning rate

V<sub>n</sub> V<sub>n</sub> PMP

## Endportant interview

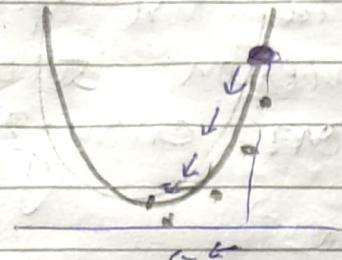
→ it tells about the "speed at which convergence happens"

which convergence happens?

Cyril

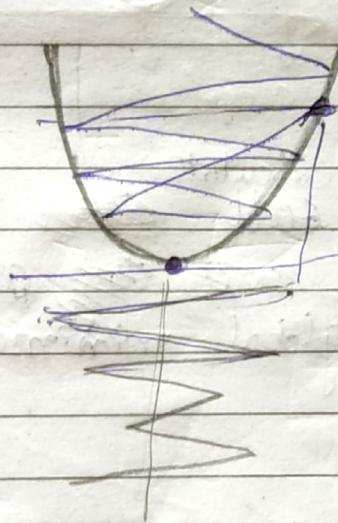
$\lambda = 0.01$  tells speed

if  $\alpha = 0.01$  [less.]



$\theta$ , change  
slowly. reach the  
global minima.

if  $\alpha = 2$  [more / high]



$\theta$ , change  
fast and.  
"absently move"  
here and there  
and near search  
the global minima

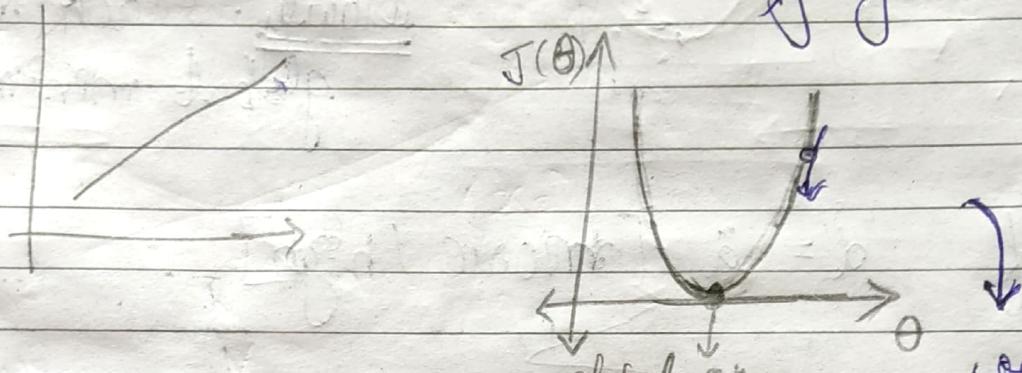
\* So, we should choose learning rate  
( $\alpha$ ) that helps in Converge.

### ~~Conclusion~~

→ "we need to repeat converge  
algorithm until we reach  
the global minima where  
cost function is very less. So,  
at that point we get best  
fit line!"

## \* gradient descent

→ it is an optimiser function that keep on changing slope



we have  
calculate point we try  
to fall ~~slope~~ ~~slope~~  
and ~~slope~~ aim to reach  
global minima

## \* Convergence algorithm

repeat until convergence

$$\theta_j^n = \theta_j^0 - \alpha \left[ \frac{d}{d\theta_j} J(\theta_j) \right]$$

## \* Cost function

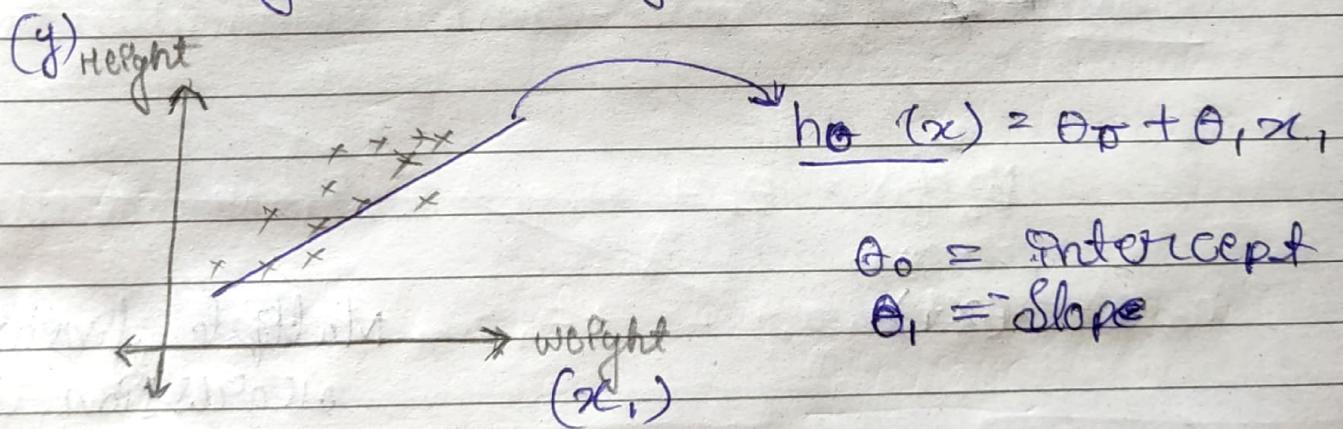
$$J(\theta_0, \theta_1) = \frac{1}{n} \sum_{i=1}^n [y_i - h_\theta(x)]^2$$

## \* Multiple linear regression

Ex: Dotatset

independent  
 $\downarrow$   $x_1$ .  
weight

Simple linear

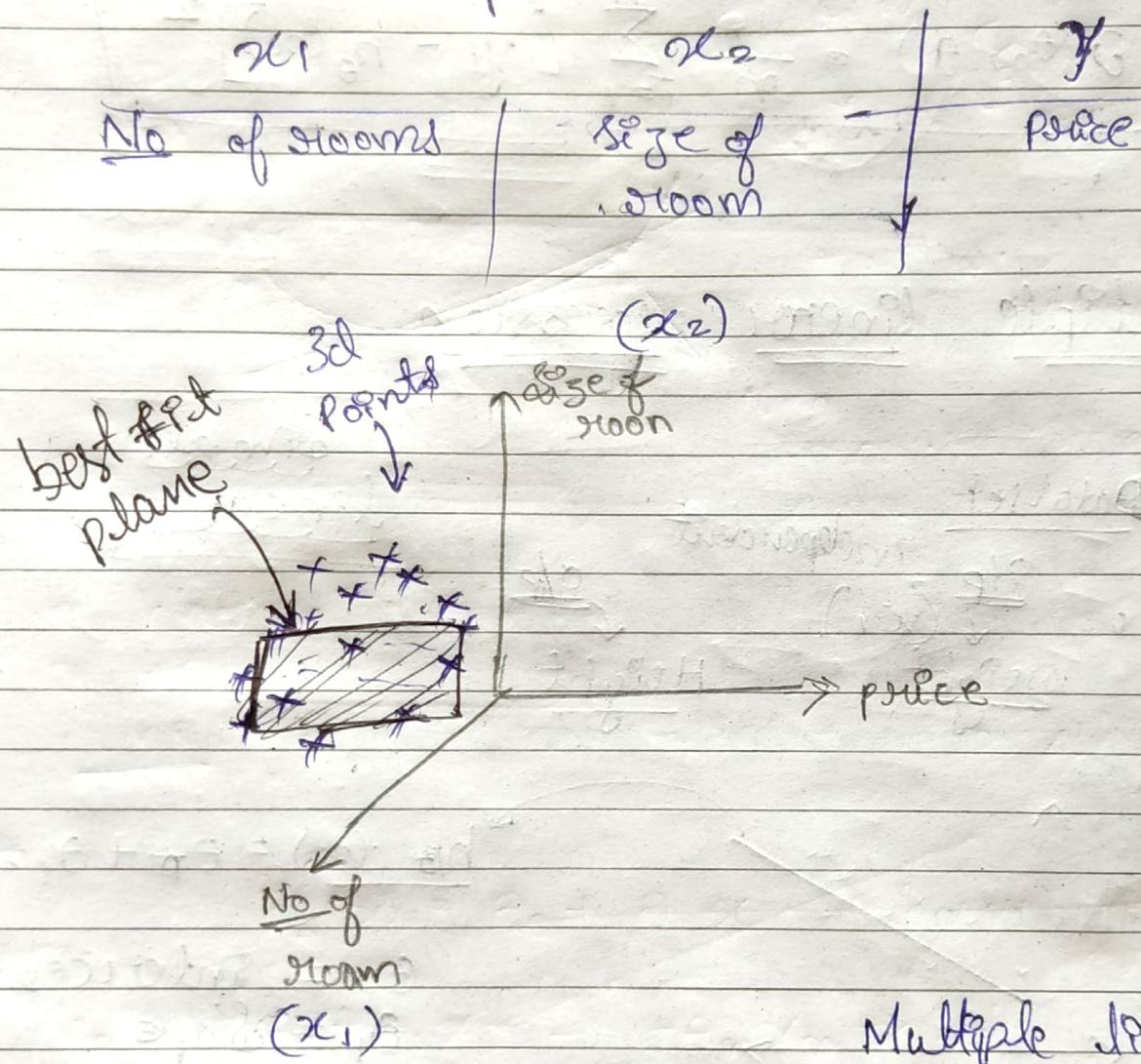


$\theta_0$  = intercept  
 $\theta_1$  = slope

P.T.C

## House price Dataset

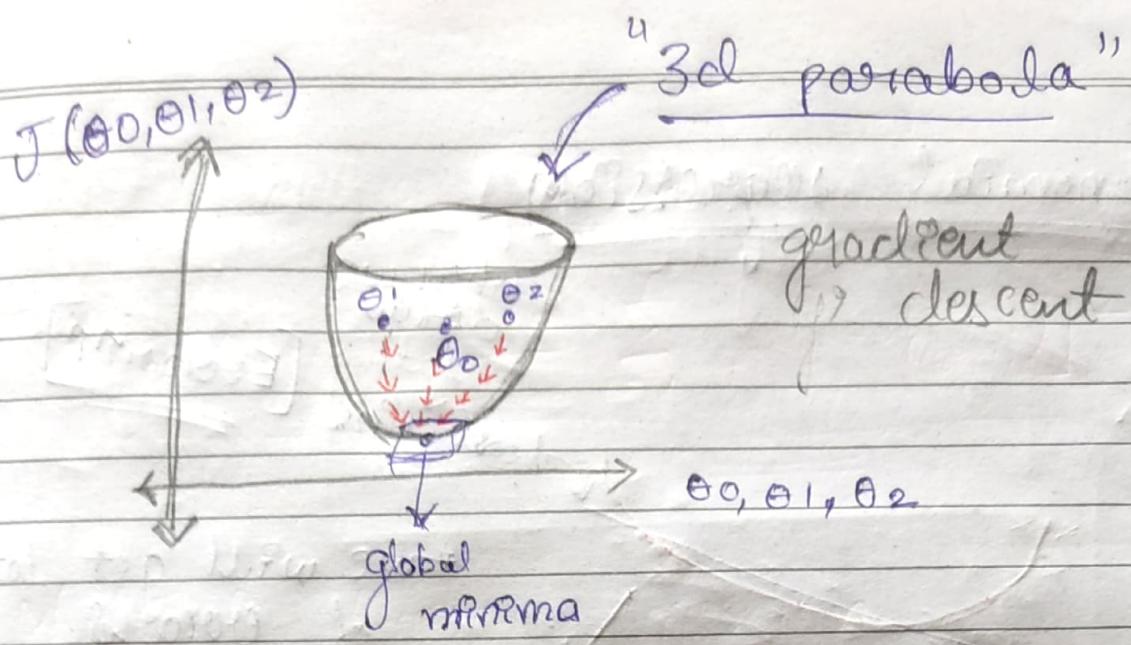
independent



$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

$\theta_0 \rightarrow$  intercept  
 $\theta_1, \theta_2 \rightarrow$  slopes /

No of Feature  $\uparrow\uparrow$  then  $\frac{\text{No of}}{\text{slope}} \uparrow\uparrow$

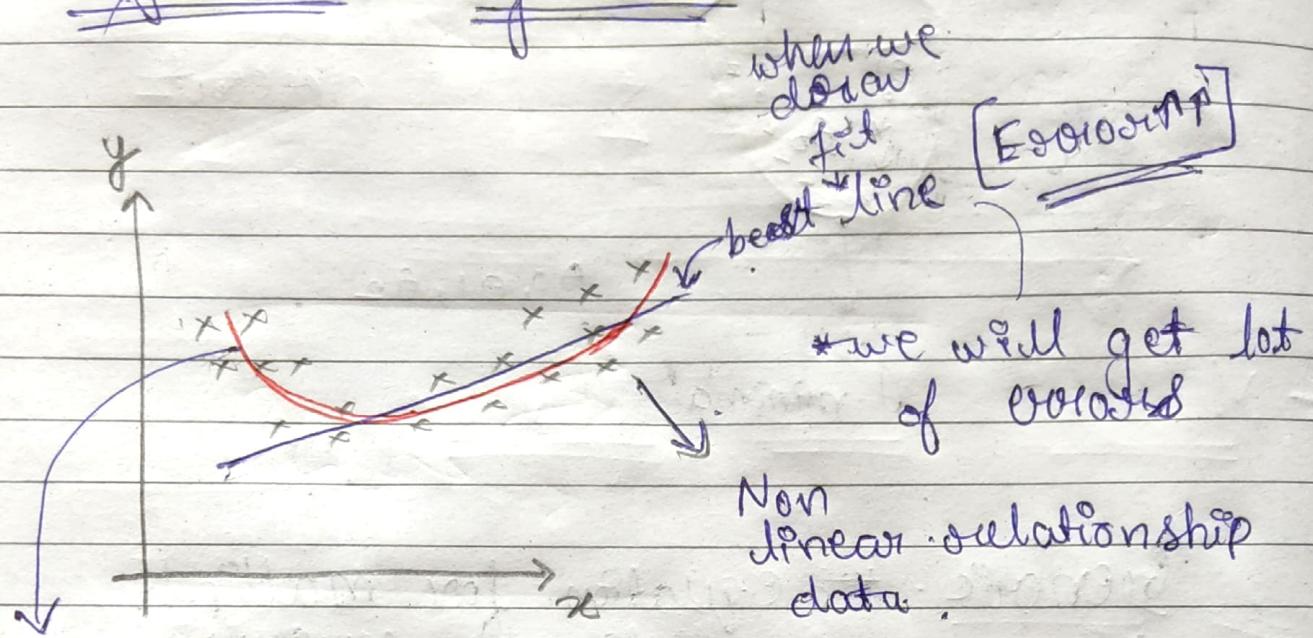


Generic equation for multiple linear regression

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

% Maximum  $\rightarrow$  Human can see 3D figure

## \* Polynomial Regression



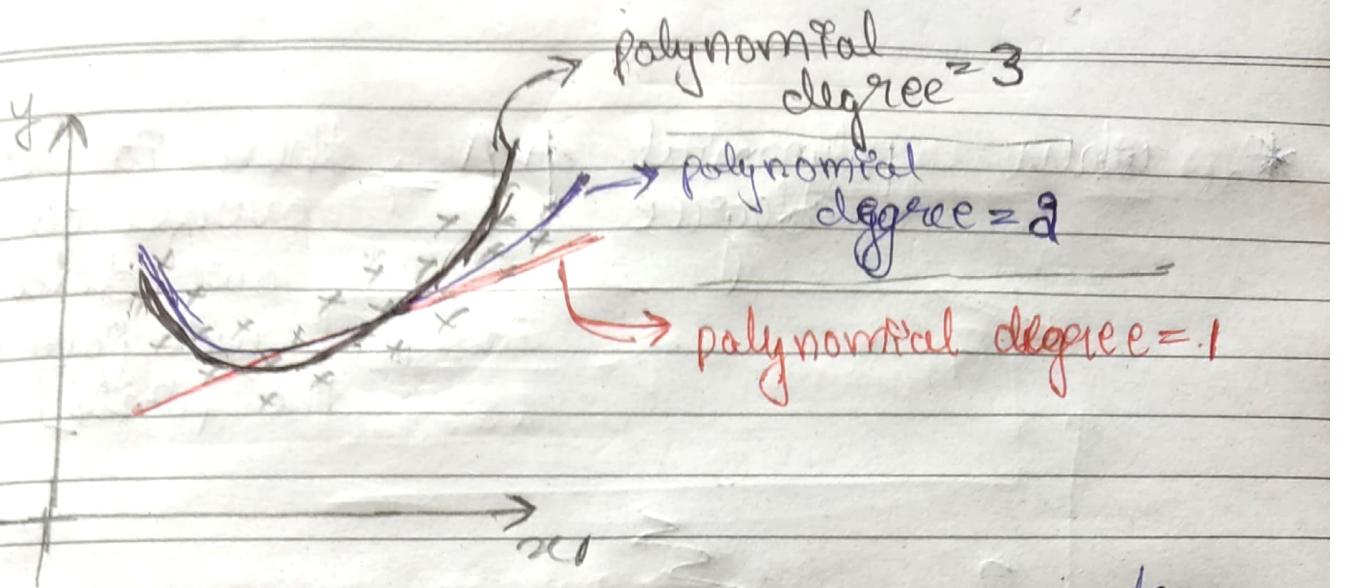
\* we have created a best fit curved line for a non-linear relationship data.

$$h_0(x) = \theta_0 + \theta_1 x \rightarrow \text{linear regression}$$

$$h_0(x) = \theta_0 + \theta_1 x + \dots + \theta_n x_n \rightarrow \text{multiple regression}$$

used when data is having non-linear relationship. ← Polynomial regression

↓  
polynomial degrees.



\* Polynomial degree = 0.  $\Rightarrow h_0(x) = \theta_0 + \underline{\theta_1 x_1}$

$$h_0(x) = \theta_0 * 1$$

$$h_0(x) = \underline{\theta_0} * x_1^{(0)}$$

\* Polynomial degree = 1  $\Rightarrow h_0(x) = \theta_0 * x_1^{(0)} + \underline{\theta_1 * x_1^{(1)}}$

[Simple linear regression]

\* Polynomial degree = 2  $\Rightarrow h_0(x) = \theta_0 x_1^{(0)} + \theta_1 x_1^{(1)} + \underline{\theta_2 x_1^{(2)}}$

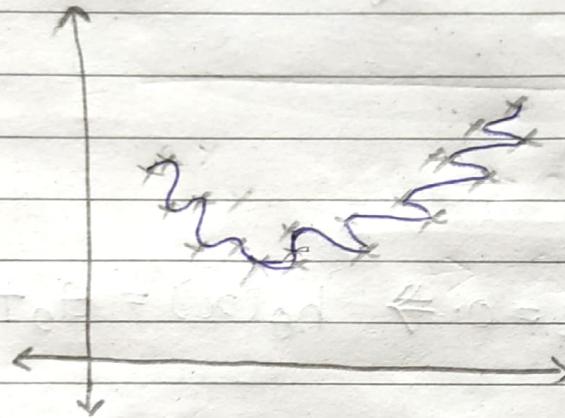
Curve fit even better for our data

\* Polynomial degree = 3,

$$h_0(x) = \underline{\theta_0 x_1^{(0)}} + \theta_1 x_1^{(1)} + \underline{\theta_2 x_1^{(2)}} + \theta_3 x_1^{(3)}$$

\* when polynomial degree ↑↑  $\Rightarrow$  then it will overfit to our data.

[it will join all the points]



∴ we need to find a suitable polynomial degree so that it does not overfit to data.

\* polynomial  $\approx n$  degree

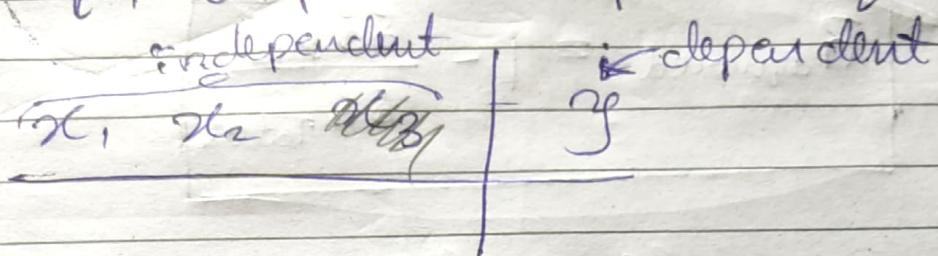
$$h_0(x) = \theta_0 x_i^0 + \theta_1 x_i^1 + \theta_2 x_i^2 + \dots + \theta_n x_i^n$$

Simple polynomial regression

[just for 1 independent feature & 1 dependent feature]

## \* Multiple polynomial Regression

{ Multiple independent features?



then we have equation for polynomial  
degree = 2

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2$$

## \* Performance metrics used in linear regression

- \* what is the accuracy?
- \* How do we calculate it?
- \* when do we say that we got better model?

We have 2 different metrics

- ① R squared.
- ② Adjust R squared.

introduction

## ① R Squared:

Say error getting from best fit line should be less than error getting from average line.

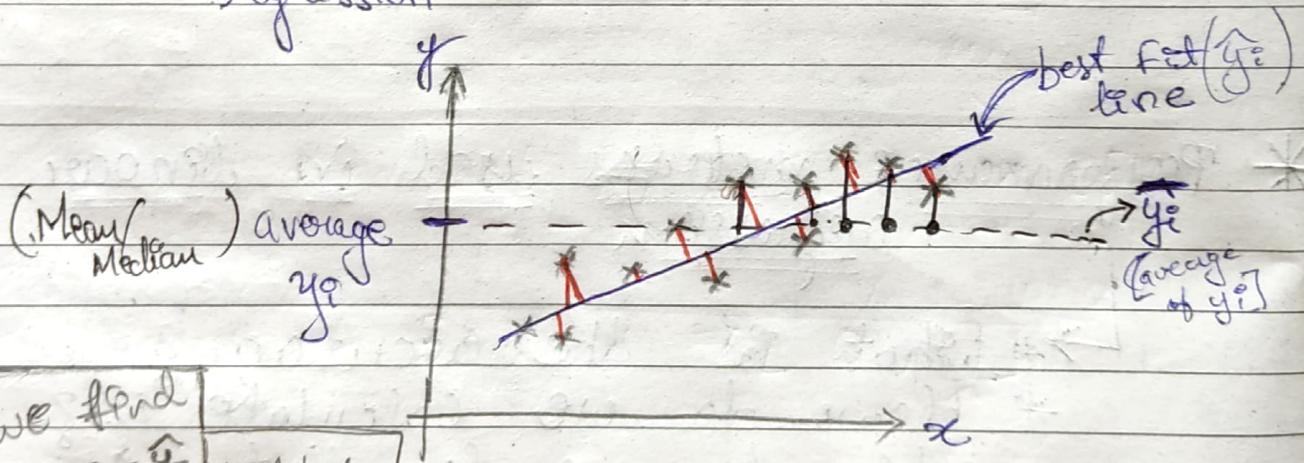
$$R_{\text{Squared}} = 1 - \frac{\text{SS Res}}{\text{SS Total}}$$

\*  $\text{SS} \rightarrow \text{sum of square}$

$\text{Res} \rightarrow \text{residual } \{ \text{difference of } y_i \text{ & } \hat{y}_i \}$

(or) error

Aim of sample linear regression  $\Rightarrow$  To get Best Fit line



We find  $y - \hat{y}$  from best fit line

\*  $\text{SS Res} \Rightarrow \text{sum of square residuals/error}$

$\text{SS Total} \Rightarrow \text{sum of square total.}$

[difference b/w actual point & predicted point w.r.t to mean line]

error w.r.t  
best fit line

$$R^2 = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

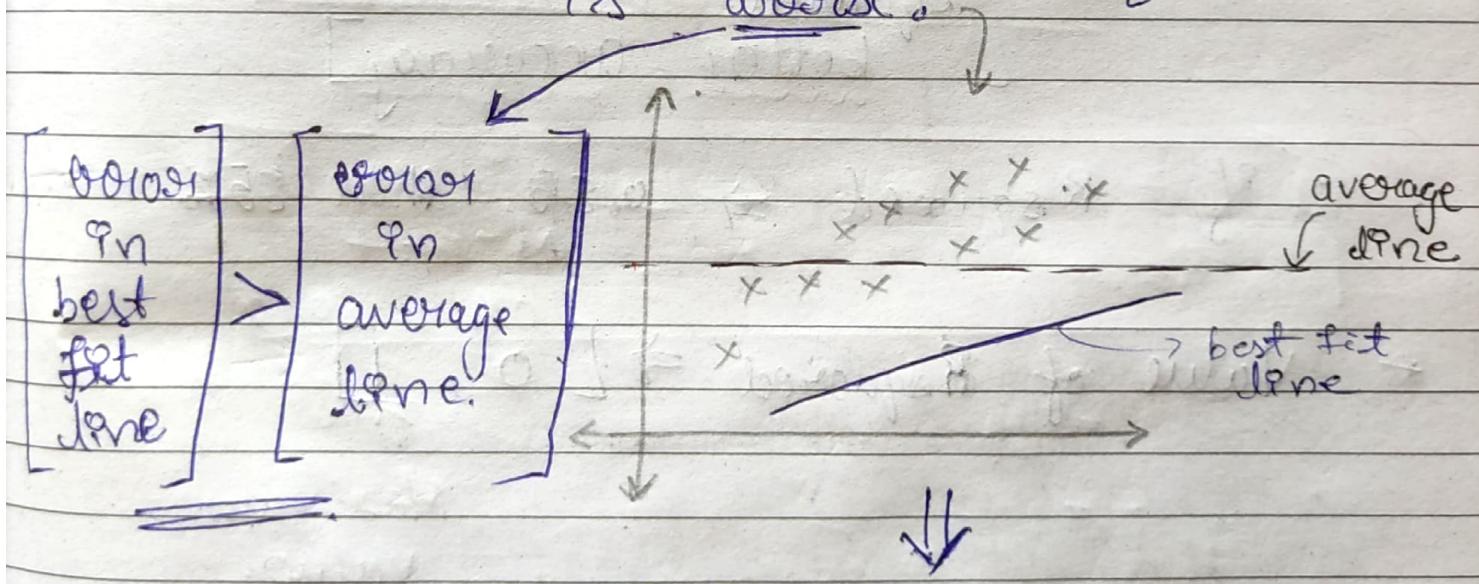
↓      ↓  
smaller value      larger w.r.t.  
the value      average line  
of  $y_i$

better the accuracy.

interview asked

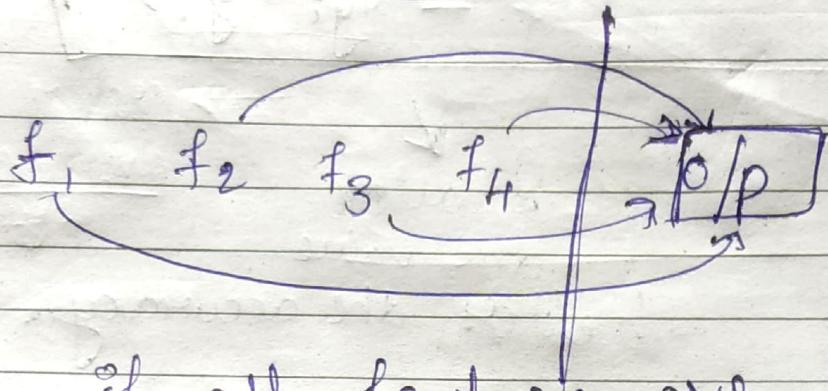
Q) Can ~~best fit~~ "  $R^2$  have -ve value " ?

⇒ Yes, only if the best fit line is worse.



$R^2$  will be "-ve"

\* Suppose



if all feature are highly correlated with output value  
then ↓

\*  $R^2$  will increase.  
[ better accuracy ]

\* if  $R^2$  = 0.85 then 85%.

\* Value of  $R^2$   $\Rightarrow [0, \frac{1}{\text{↑}}]$   
0% to 100%.

Maximum value.

\* we should not get  $[R^2 = 1]$  because it tell model is "overfitted"

## ② Adjusted R squared:

<u>Explanatory</u>	Independent	Dependent
Size of house	No of bed rooms	Price
$x_1$	$x_2$	$y$

Here,  $x_1$  &  $x_2$  are directly correlated with price

$$\text{suppose } \rightarrow R_{\text{squared}} = 85\%$$

\* Now, if we had "location" feature

↙  
it is also highly correlated  
with output

$$\rightarrow R_{\text{squared}} = 90\% \quad \begin{matrix} \text{drastically} \\ \text{increase} \end{matrix}$$

\* Now, if we add "gender"  
less correlation ↴

$$\Rightarrow R_{\text{squared}} = 91\% \quad \begin{matrix} \text{less} \\ \text{increase} \end{matrix}$$

\* So, when we add on feature our R square will keep on increasing.

" tells how model is behaving with addition of new feature "

This will not happen with "adjusted R square".

To prevent this we use.

$$\text{Adjusted R squared} = 1 - \frac{(1-R^2)(N-1)}{N-p-1}$$

N = No of data points

p = No of independent features

Suppose,  $R^2 = 0.8 \Rightarrow 80\%$ ,  $N = 11$ ,  $p = 2$

$$\begin{aligned} \text{Adjusted R square} &= 1 - \frac{(0.8)(10)}{11-2-1} = 1 - \frac{8}{8} \\ &= 1 - \frac{1}{4} \\ &\Rightarrow \frac{3}{4} = \underline{\underline{0.75}} \end{aligned}$$

when  $R^2 = 80\%$  then adjusted  $R^2 = 75\%$

~~$R^2 > \text{adjusted } R^2$~~  (always)

\* if we add ~~1~~ independent feature  
~~(imp)~~

$$\cancel{R^2 = 85\%} \rightarrow \text{adjusted } R^2 = 78\% \cancel{\uparrow}$$

\* if we add ~~1~~ independent feature  
(not imp)

$$\cancel{R^2 = 86\%} \rightarrow \text{adjusted } R^2 = 76\% \downarrow$$

better and  
to measure accuracy

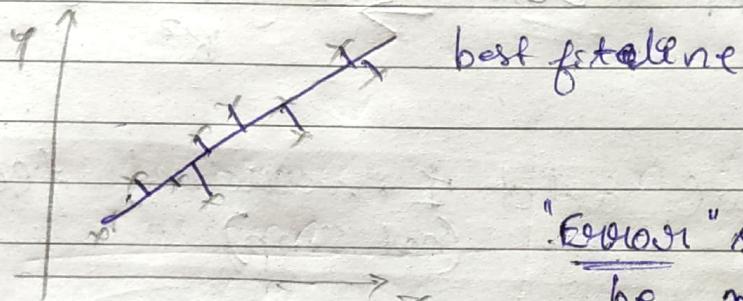
	R-square	Adjusted R-square
* For any new <del>feature</del> , <u>imp</u> feature added	Increase drastically $\uparrow\uparrow$	Increase $\uparrow\uparrow$
* <del>not imp</del> feature added	Increase by small amount $\uparrow$	Decrease by small amount $\downarrow$

## \* MSE, MAE and RMS E

### [Cost function]

~~Explained clearly~~

\* To reduce error we consider



"Error" should be minimum

To measure <sup>error</sup> we use  
[Cost function].

we say

$$J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m [y_i - f_{\theta}(x_i)]^2$$

This ~~is~~ "mean-squared error"

Cost function is called as

Our aim is to reduce this error

We reduce this

error by using  $\rightarrow$  "Convergence algorithm"

## ① Mean Squared Error (MSE):

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

Quadratic equation.

we get  
parabola curve.

we can have

$$a\theta_0 + b\theta_1 + c = 0$$

$$(a-b)^2 = a^2 + b^2 - 2ab$$

we reduce

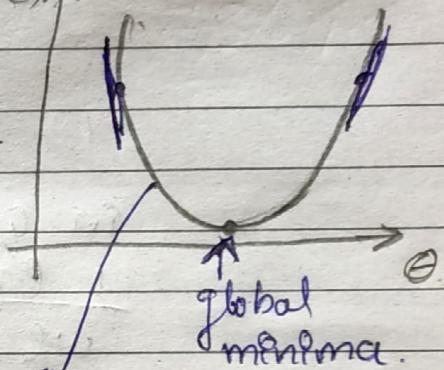
error

using "convergence algorithm"

this function

is called

"convex function"



$J(\theta)$

global  
minima

$\theta_1$

until we reach  
optimal solution  
by using gradient  
descent

used to minimize

(or) maximize  
the function

## \* Advantages

① It is differentiable.

So at any point we are able to find slope of the point.

② It has one local and one global minima.



### In quadratic equation

we know our global minima

we know where we need to converge

### In Non-quadratic equation

In graph we have many local minima.

"Non Convex Function"



So, we can't move

Value doesn't change

[ $\text{Let } \frac{d}{dx} \text{slope} = 0$ ]

$$\text{Convergence} \Rightarrow \theta_0 = \theta_0 - \alpha \left[ \frac{\partial}{\partial \theta_0} h_0(x) \right]$$

local minima  
↓  
global minima

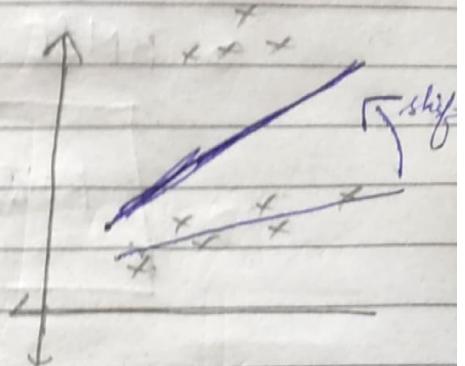
### disadvantage

pt will stuck in that region as it cannot calculate slope.

## \* DLS advantages

① Not Robust to outliers.

best fit line shift more from actual position.



② It is not in the same unit.

Price

TNR

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

$$(TNR - \hat{TNR})^2$$

Profit

③ Mean absolute error. [MAE]

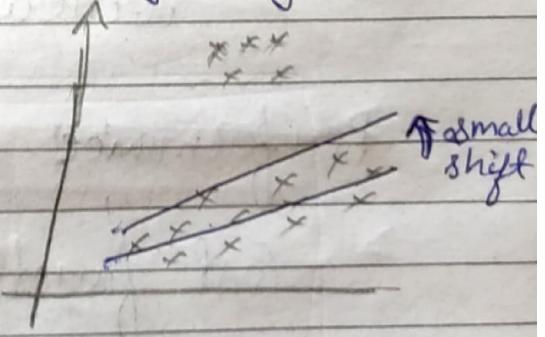
$$MAE = \frac{\sum_{i=1}^n |y_i - \hat{y}_i|}{n}$$

only magnitude

## Advantage

① Robust to outliers.

② It will be in same unit.



## Disadvantage

① Convergence usually takes time.

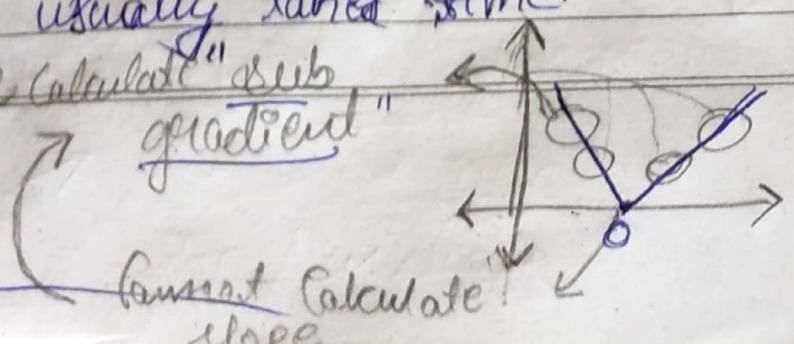
"optimization"

② Is "Complex"

③ Time

we calculate "sub gradient"

cannot calculate slope



### ③ RMSE [Root mean square error]

$$RMSE = \sqrt{MSE}$$

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

#### \* Advantages

- ① Same unit.
- ② Differentiable

#### \* Disadvantage

- ③ Not robust to outliers.

### Notes

linear regression.

Performance check  $\rightarrow R^2$  and adjusted  $R^2$

Cost function  $\rightarrow MSE, MAE, RMSE$

# \* Simple Linear Regression in Python

> plt.scatter(—, —)

- ① Dataset division (Independent & dependent)
  - ② Train, test split of dataset.
  - ③ Standardization.
  - ④ Model training (⑤ testing & final accuracy).
- > from sklearn.model\_selection import train\_test\_split

> X-train, x-test, y-train, y-test = train\_test\_split(x, y, test\_size=0.2, random\_state=42)

> from sklearn.preprocessing import StandardScaler

> scalar = StandardScaler()

> { scalar.fit, transform(x-train)}  
{ scalar.transform(x-test)}

train model > from sklearn.linear\_model import LinearRegression

> regressor.fit(x-train, y-train)  
finds  $\theta_0$  &  $\theta_1$

- > regressor.coef\_  $\theta_i$  give slope
- > regressor.intercept\_  $\theta_0$  intercept
- > plt.plot(x\_train, regressor.predict(x\_train))
- >  $y_{\text{pred}} = \text{regressor.predict}(x_{\text{test}})$   
 $\text{plt.plot}(x_{\text{test}},$   
 $y_{\text{pred}})$   
 test data  
 predicted data

### Performance metrics

- > from sklearn.metrics import ~~mse~~
  - mean\_squared\_error,
  - mean\_absolute\_error,
- >  $\text{mse} = \frac{1}{n} \sum (y_{\text{test}} - y_{\text{pred}})^2$   
 $\frac{(y_{\text{test}}, y_{\text{pred}})}{\text{actual} \quad \text{predicted}}$
- >  $rmsse = \sqrt{mse}$

accuracy  $\rightarrow$  from sklearn.metrics import ~~accuracy~~  $\text{r2\_score}$

$$\text{r2\_score}(x_{\text{test}}, y_{\text{pred}}) = 1 - \frac{\sum (y_{\text{pred}} - y_{\text{test}})^2}{\sum (y_{\text{test}} - \bar{y})^2}$$

$$\text{adj\_r2} = 1 - \frac{(1 - r^2)(n-1)}{(n-p-1)}$$

## \* Multiple regression with python

> from sklearn.datasets import fetch\_california\_housing

> California = fetch\_california\_housing()

> from sklearn.model\_selection import train\_test\_split.

(2) > x\_train, x\_test, y\_train, y\_test =  
train\_test\_split(x, y, test\_size=0.33,  
random\_state=10)

(3) > from sklearn.preprocessing import StandardScaler

(4) > from sklearn.linear\_model import LinearRegression.

> regressor.fit(x\_train, y\_train)

> regressor.coef\_

> regressor.intercept\_

5 y\_pred = regressor.predict(x-test)

⑥

~~now!~~

from sklearn.metrics import mean\_squared\_error  
mean\_squared\_error

mse\_error = mean\_squared\_error(y-test, y-pred)

⑦

accuracy → from sklearn.metrics import r2\_score

R-square = r2\_score(y-test, y-pred)

percent ( $r^2 \times 100$ )

$$R^2_{\text{adj}} = 1 - \frac{(1-R^2)(n-1)}{(n-p-1)}$$

⑧

new  
data

data

scaled\_data = scaler.fit\_transform(data)

pred = regressor.predict(scaleddata)

⑨

Assumptions

- ① scatter plot (y-test, y-pred) → linear graph
- ② test plot (residuals) → normal distribution
- ③ scatter plot (y-pred, residuals) → uniform distribution

## ⑩ Pickling

> Import pickle

> pickle.dump (regression, open('model.pkl', 'wb'))

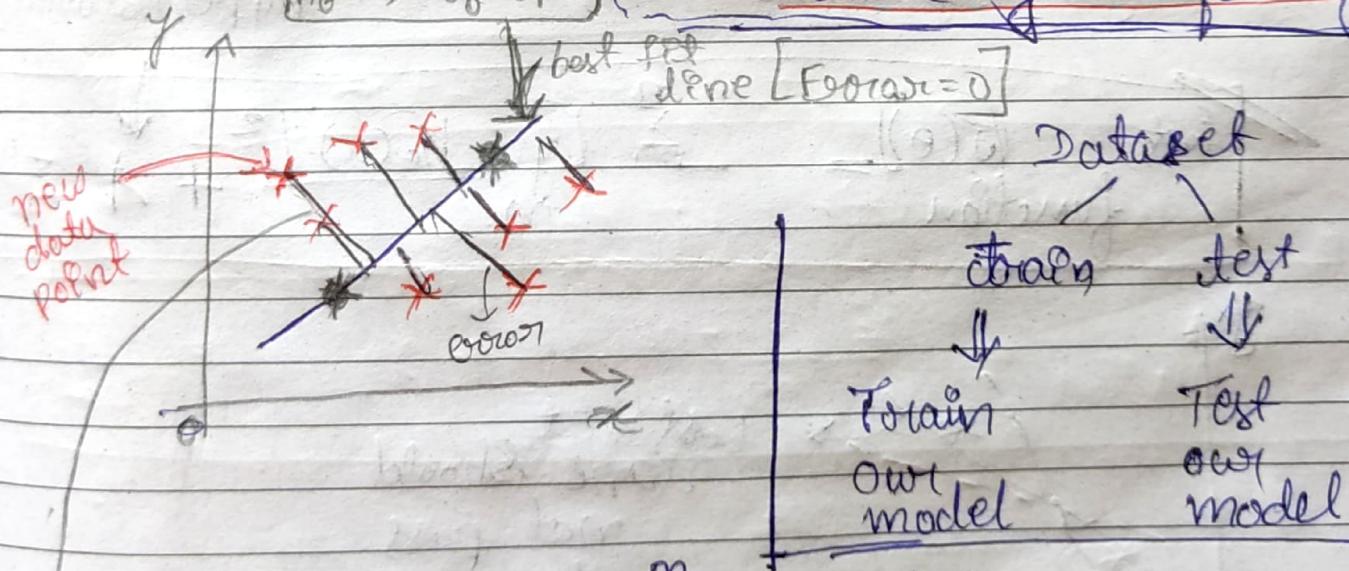
> model = pickle.load(open('model.pkl', 'rb'))

> model.predict(new\_scaled\_data)

## \* Ridge regression [L2 regularization]

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

Used to "reducing overfitting"

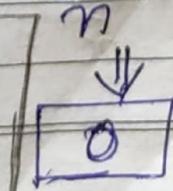


$$\text{Cost function} = \frac{1}{n} \sum_{i=1}^n [y_i - h_{\theta}(x_i)]^2$$

Overfit is more.  $\theta_0$

Accuracy decreases

This is called overfitting



as variance is 0

low bias  
high variance

model is trained properly.

## \* overfitting

Train dataset  $\rightarrow$  Accuracy  $\uparrow\uparrow$

Test dataset  $\rightarrow$  Accuracy  $\downarrow\downarrow$

Training show bias and Testing High variance.

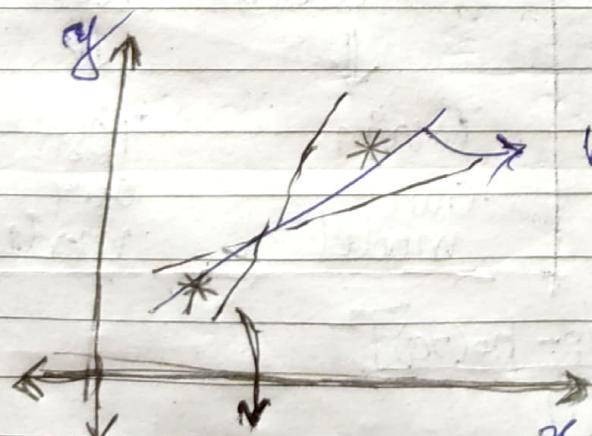
$\downarrow$  error is more as data spread is more so, accuracy is less.

\* In order to overcome the above problem we use ridge regression

## \* Ridge regression

by parameter

~~Cost function~~  $J(\theta) = \frac{1}{n} \sum_{i=1}^n [y_i - h_\theta(x)]^2 + \lambda \sum_{j=1}^m \theta_j^2$



best fit line should not pass through 2 points

because error will be 0.

best fit line may be looked like this.

\* Let  $\lambda = 1$

$$h_{\theta}(x) + \theta_0 + \theta_1 x_1$$

$\theta_1 \rightarrow \text{slope}$

then

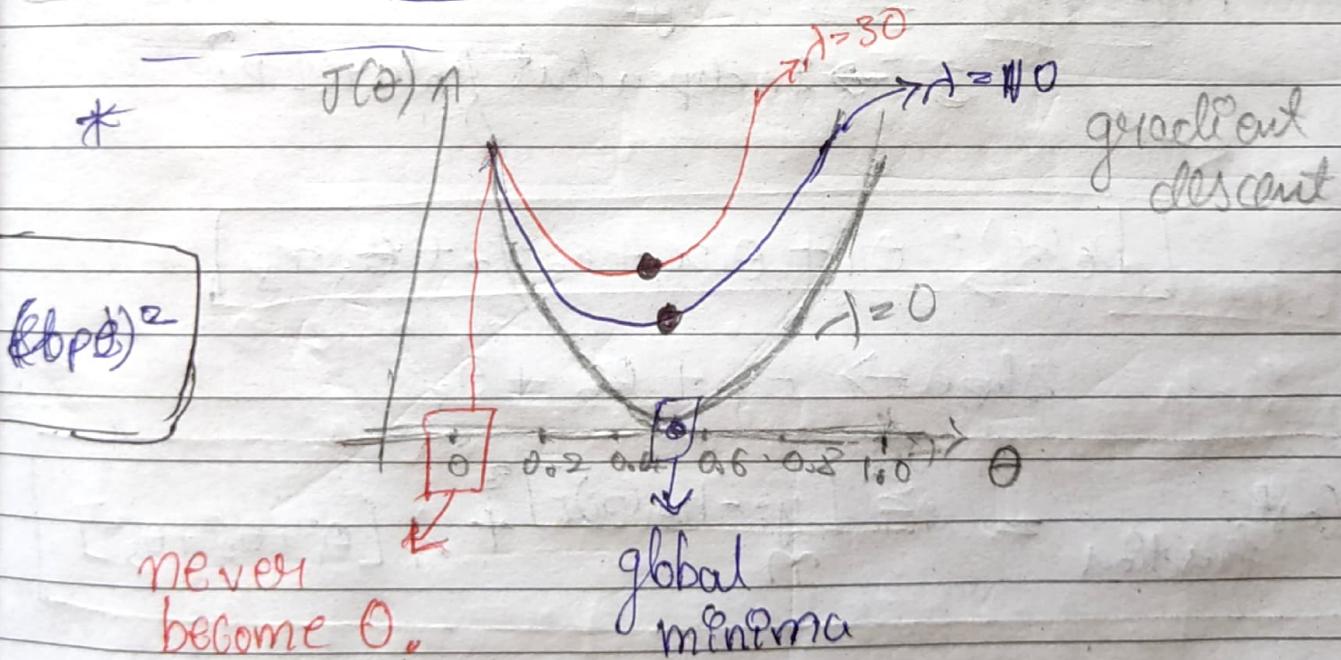
Cost function =  $\frac{\sum [y_i - h_{\theta}(x_i)]^2}{n} + \lambda \sum (\text{slope})^2$

by penalizing cost function  $= 0 + \lambda [(\theta_1)^2]$  [penalizing cost function]

reduce overfitting

$\rightarrow$  error will not be 0.

~~Always~~: Cost function  $> 0$  and never 0.



\* As  $\lambda \uparrow \uparrow \rightarrow$  global minima move upwards.  
 $\rightarrow$  slope  $\downarrow$

## Simple linear regression

$$h_{\theta}(x) = \theta_0 + \theta_1 x \quad \theta_1 \rightarrow \text{slope}$$

Cost function =  $\frac{\sum_{i=1}^n [y_i - h_{\theta}(x)_i]^2}{n} + \lambda \left( \sum_{i=1}^n (\theta_j)^2 \right)$

$\downarrow \theta_1^2$   
d<sub>2</sub> norm

## Multiple linear regression

3  $\rightarrow$  independent feature [p]

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

slope  $\Rightarrow \theta_1 + \theta_2 + \theta_3$

Cost function =  $\frac{\sum_{i=1}^n [y_i - h_{\theta}(x)_i]^2}{n} + \lambda [\theta_1^2 + \theta_2^2 + \theta_3^2]$

## \* Lasso Regression [L1 regularization]

<sup>if</sup> Feature selection

$$\text{Cost function} = \frac{1}{n} \sum [y_i - h(\theta)x_i]^2 + \lambda \sum_{i=1}^n |\text{slope}|$$

\* lastly, consider

$$\text{MLR}, h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

tells that  
unit change  
in  $x_i$ , how much  
change in  $y$

$$= 0.52 + 0.65 x_1 + 0.72 x_2 + 0.12 x_3$$

tells that  
unit change in  
 $x_3$ , there is  
only 0.12  
change in  $y$

we can conclude that  
 $x_3$  has less correlation  
with output feature ( $y$ )

so, we can remove this  
 $x_3$  feature.

\* For the above process we  
use lasso regression.

Cost function.

$$\frac{1}{n} \sum [y_i - h_{\theta}(x)]^2$$

+  $\lambda \sum_{j=1}^m | \theta_j |$  if sloped

if  $\lambda = 1$ ,

$$Error + 1 [|\theta_1| + |\theta_2| + |\theta_3|]$$

$\lambda = 10$

$\lambda = 20$

$\lambda = 10$

gradient  
descend

$\lambda = 0$

at one point it  
become 0.

but in ridge regression it  
will never become 0

$(\theta)$  which is less correlated become 0.

\* Now,

$$h_{\theta}(x) = 0.52 + 0.65x_1 + 0.072x_2 + 0.12x_3$$

Here, feature  
selection has  
happened

This less  
Correlated feature  
will actually become  
0 at 1 point.

\*  $\lambda \uparrow$  slope  $\downarrow \downarrow$

At one point,  $\theta \approx 0$

\* Elastic Net regression  $\rightarrow$  lasso + Ridge  
Combination of

① Reduce overfitting

[Ridge]

② Feature selection

[Lasso]

Cost function =  $\frac{1}{n} \sum_{i=1}^n [y_i - h_\theta(x)_i]^2 + \lambda_1 \sum_{j=1}^m (\text{slope})_j^2 + \lambda_2 \sum_{j=1}^m |(\text{slope})_j|$

MSE

Reduce overfitting

[L1 regularization]

Feature Selection

[L2 regularization]

$\lambda_1, \lambda_2 \Rightarrow$  { hyperparameter tuning }

$\lambda_1, \lambda_2 = 0.5$   
remaining given to  
L2 regularization

## \* Ridge, Lasso, ElasticNet with python

Ridge  $\Rightarrow$  from sklearn.linear\_model import Ridge

Lasso  $\Rightarrow$  from sklearn.linear\_model import Lasso

ElasticNet  $\Rightarrow$  from sklearn.linear\_model import ElasticNet

\* All the supervised learning model have →

fit and predict method  
 $\rightarrow$  regressor.fit(x\_train, y\_train)  
 $\rightarrow$  regressor.predict(x\_test)

## \* Algorithm forest

① ~~# cleaning~~ <sup>new</sup> New feature  
    > df.loc[6122, "Region"] = 1  
        df.loc[1228, "Region"] = 2.  
            <sup>now</sup> <sup>Column</sup>

> df = df.dropna().reset\_index(drop=True)

> df.iloc[[122]]

> df.columns = df.columns.str.strip()

> df[[-1]] = df[[-1, -2]].astype(float)

>

## ② EDA

> df["class"] = ~~df["class"]~~  
    np.where(df["class"] str.contains("fire"),  
              1, 0)

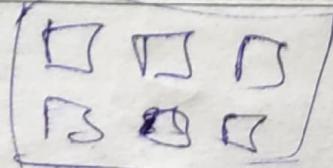
> drop → day, month, year.

— plot

> plt.style.use("seaborn")

df.hist(bins=50, figsize(20, 15))

— plt.show()



>  $\text{percent} = \text{df\_copy['classes'].value\_counts(normalize)} * 100$

`pct = pd.DataFrame([percent, labels = class_labels, autopct = '%.1f%%'])`

## Training

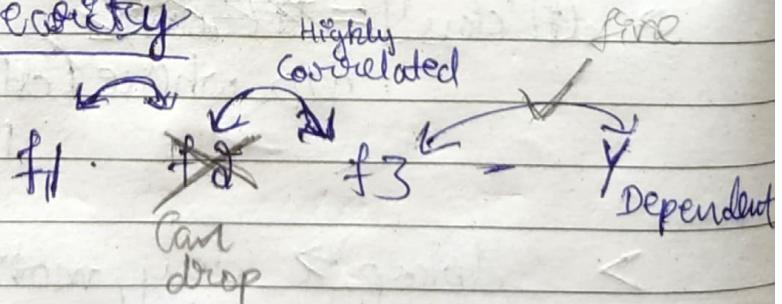
① Divide dependent and independent feature.

② Train - Test split.

→ ③ Feature selection.

domain expert → > 85% correlation drop.  
gives threshold because same  
S/P most needed  
[Similar feature  
should be dropped]

## Multicollinearity



→ def. covarlet(df, threshold):

col\_coor = set()

col\_matrix = df.cov()

for i in range(df.columns):

for j in range(i)

if abs(col\_matrix.iat[i, j]) > thres:

-ve  
value.

col\_name = ~~colname~~

col\_matrix.col[i]

col\_coor.add(col\_name)

return col\_coor

→ Xtrain, Xtest → drop col\_coor.

④ Standardize. (Xtrain, Xtest)

understand imp of scaling

⑤.1 linear regression

> from linear\_model import LinearRegression

⑤.2 lasso regression (feature selection)

> import Lasso

⑤.3 cross validation Lasso.



Dataset

Test

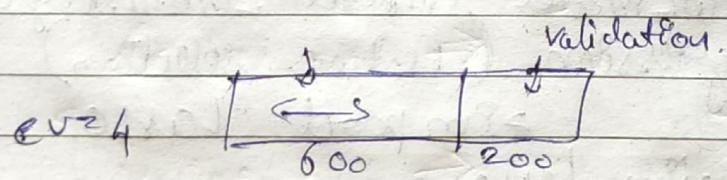
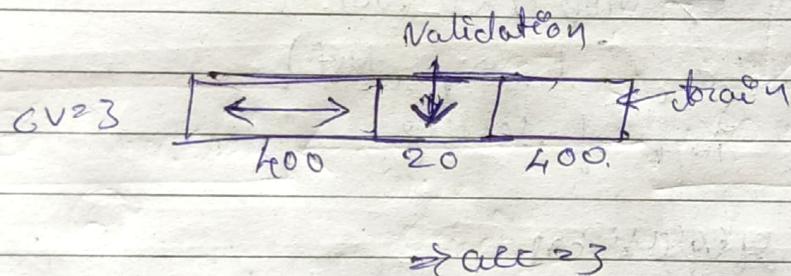
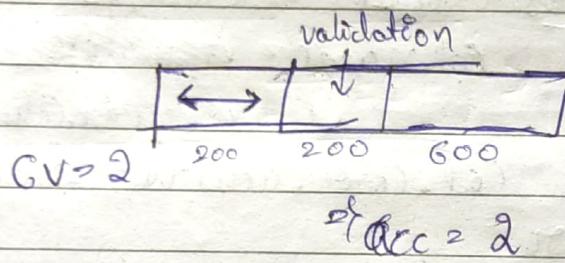
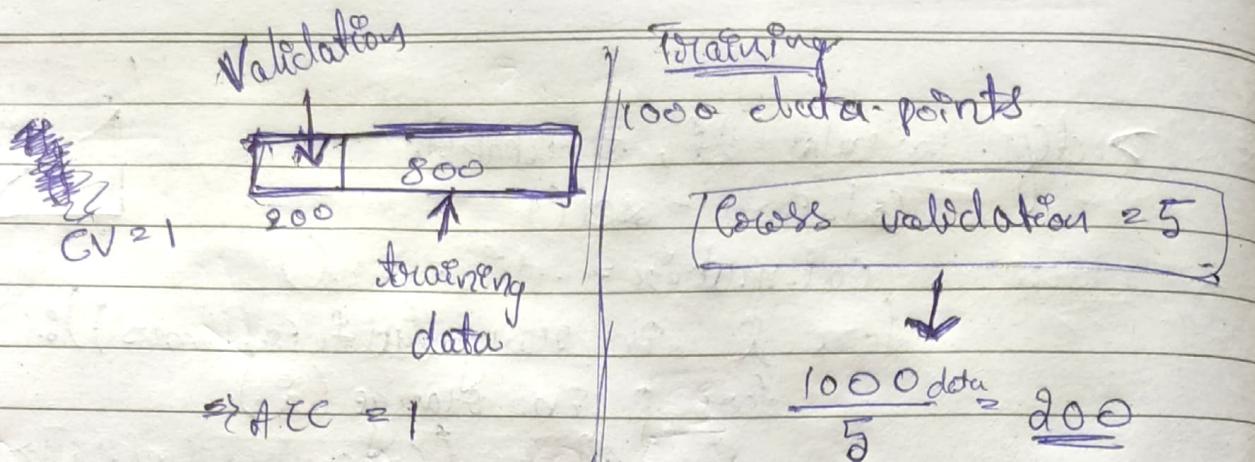
Train

Train

Validation

Training  
our model

Hyperparameter  
tuning.



Validation data  
↓  
used for hyperparameter tuning.

\* Final accuracy = average of all accuracy

\* Here, in every cross validation we change our training test data & validation data.

↓  
To check model

So, this called Cross-validation.

\* > from sklearn.linear\_model import LassoCV

↳ LassoCV = LassoCV(CV=5)

# rest all same

④ Ridge regression - reduce overfitting

> from sklearn.linear\_model import Ridge

⑤ Ridge  
cross validation  
> import RidgeCV

⑥ Elastic net regression

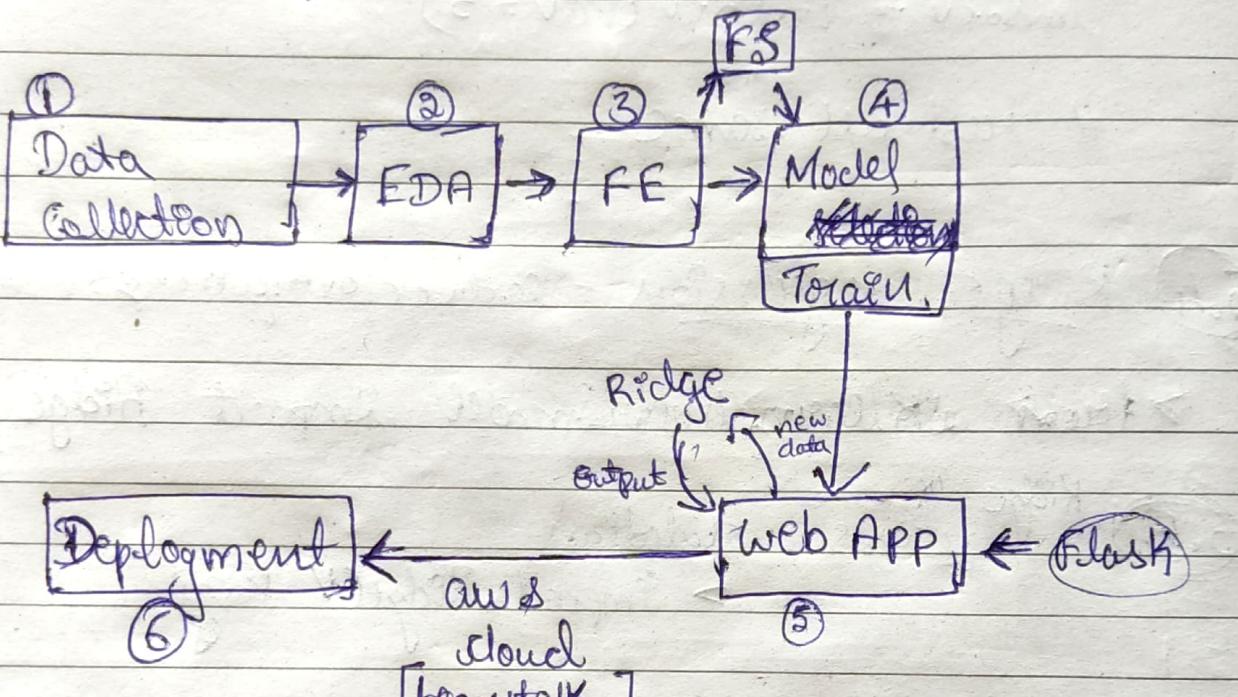
> import ElasticNet

⑦ > import ElasticNetCV

⑧

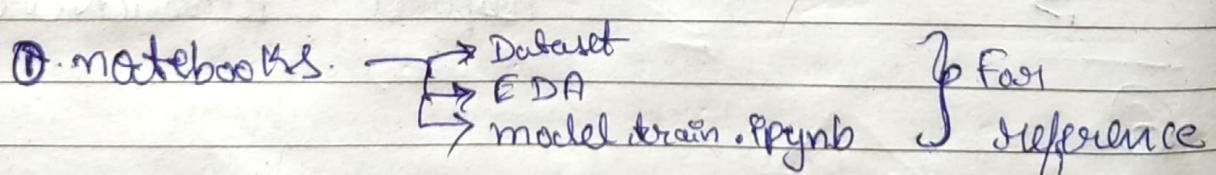
# \* Machine learning Projects

## ML Project life cycle



{ Scaler.pkl → Standardization } → Feature Scaling  
 { Ridge.pkl → Data Prediction } → Model.

## \* Folder



③ templates ↴  
 index.html ←  
 home.html ←

add all the library required for project

→ pip install -r requirements.txt

④ application.py



It will install all library to environment when run Command.

→ to automatically add required lib  
→ python -m pipreqs -p requirements

→ from flask import Flask, request, jsonify, render\_template

→ import pickle

→ import numpy as np  
pandas

→ from sklearn.preprocessing import StandardScaler

→ ridge\_model = pickle.load(open('model/ridge.pkl'))  
scaler = StandardScaler

→ app = Flask(\_\_name\_\_)

② app.route('/')

def index():

render\_template('index.html')

goes to template folder and renders "index.html" file

\* GET → getting some a static page.

POST → <sup>requesting</sup> passing some info and getting response to that particular info.

@app.route ('/predict-data', methods=['GET', 'POST'])

def predict\_datapoint():

if request.method == 'POST':

temp = float(request.form.get('Temp'))

}

new\_data = [ ]

new\_scaled\_data = scaler.transform([new\_data])

result = ridge\_model.predict(new\_scaled\_data)

return render\_template('home.html', resu

else:

return ('home.html')

if name == 'main':  
app.run(host='0.0.0.0')

⑥

home.html

def predict\_data():

<form action="{} url\_for('predict\_datapoint')>

method="post">

<input type="text", name="Temp.",  
placeholder="—" required="true" />

<button type="submit".> Predict </button>

</form>

<h2> {{ results }} </h2>

# Basic Library + Final

→ home.html

<input type="text" name="Temp"/>

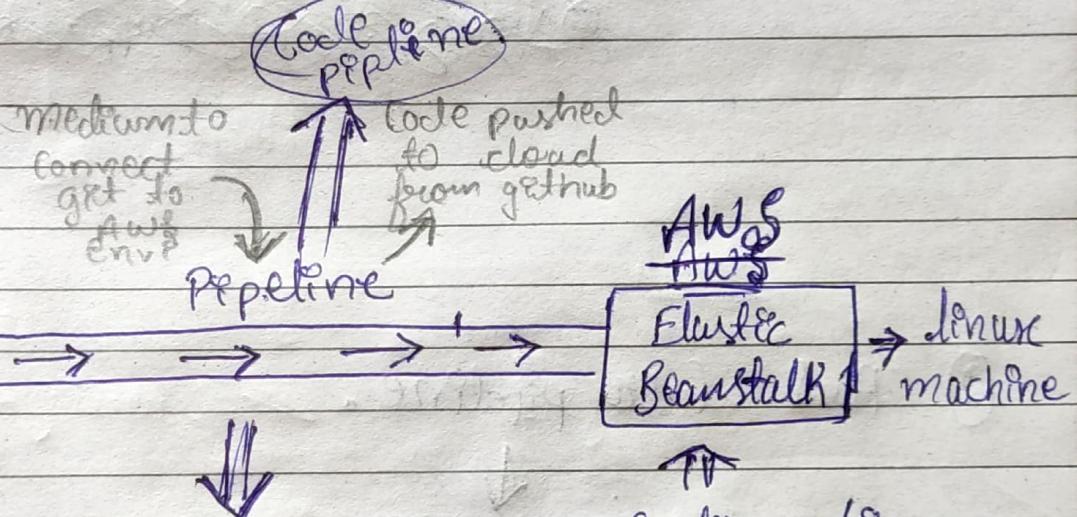
## Aws deployment

= result[0])

Code base

github  
repository

Any  
change

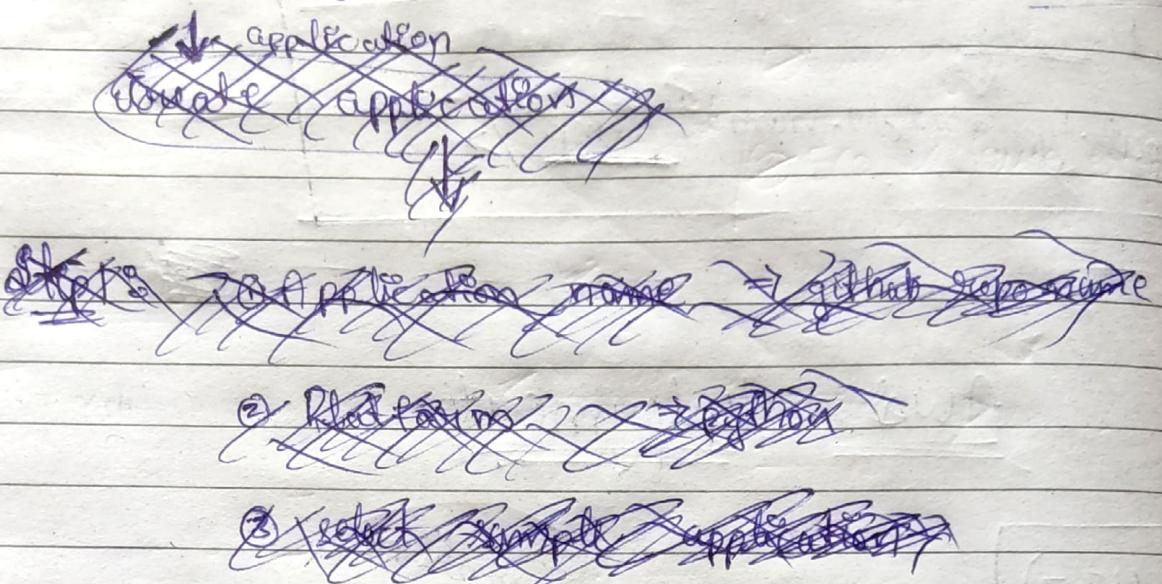


Folders → <sup>+ add</sup> ① .ebextensions  
↳ python.config

"link + predict data"

In AWS

## ① AWS Elastic Beanstalk.



## ② AWS CodePipeline.

- ↓
- Create new pipeline
- Step 1: ① pipeline name → forest-forepipe
- Step 2: ① source provider → github <sup>select</sup> version+  
② Connect github.  
③ select repo.  
④ select branch.

Step 3: skip (for now)

Step 4: ~~skip~~ select Beanstalk.

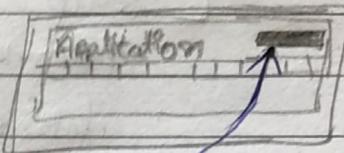
+ name + environment

## Elastic Beanstalk



### ① Create application [Folder]

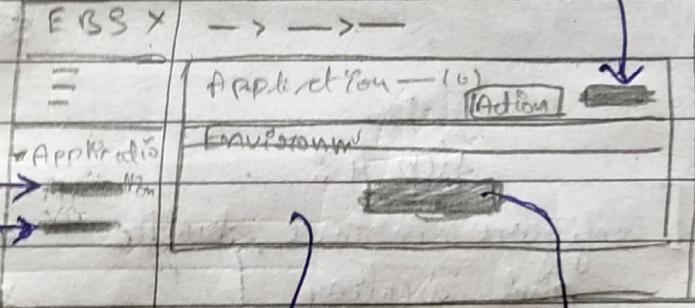
- ① enter name
- ② ~~Create new application~~



### ② \* we can see

Application versions

Saved configurations



Environment

Create Environment

### ③ Create environment

step 1: \* web server environment  
Configure Application Name [already created]  
environment ~~platform~~ ~~python [select]~~  
~~\* Sample application [selected]~~

For first three we have  
to create it

Step 1:

[Configure service access]

\* Use an existing service role.

(~~existing service role~~)

\* other default

(~~EC2 instance profile~~)

Step 2:

[Setup networking, database & tags]

\* VPC [select default]

instance settings

\* select 2 different subnets

database

\* - - - - -

~~Setup~~ defaults

Step 3:

[Configure instance traffic and scaling]

\* ~~default~~ default

~~Security group~~

Step 4:

Enhanced

\* Monitoring

\* System [select ~~Basic~~]

\* Managed platform updates

\* "deactivate"

\* defaults

Step 5:

Subnets

Need to create IAM role

AWS Console → IAM



roles

① Create role

Step 1:

④ AWS service.

② Use Case  select Elastic Beanstalk

Step 2:

③ policy



~~AWS Elastic Beanstalk WebTier~~

WebTier

Step 3:

"AWS Elastic Beanstalk WebTier"

Step 3:

① Role name

Submit

## \* Logistic Regression

→ Used to solve classification problem

Binary classification

o/p: 2 categories

Multiclass classification

o/p: > 2 categories

\* Consider,

Dataset

[Independent]

No of play hours

9  
8  
7  
6  
5  
4  
3  
2

[IP]

[Dependent feature]

Pass / Fail [y]

Fail

[0]

Fail

0

Fail

0

Fail

0

Pass

[1]

Pass

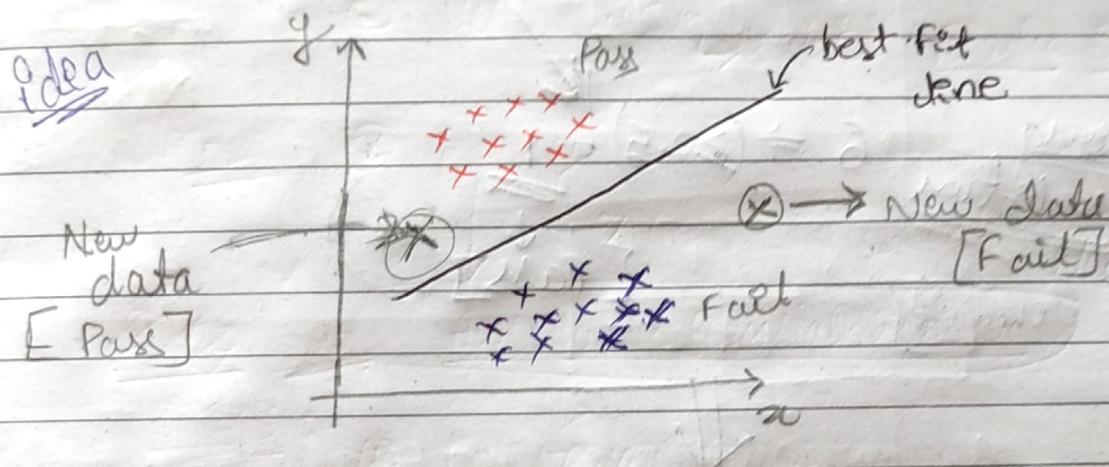
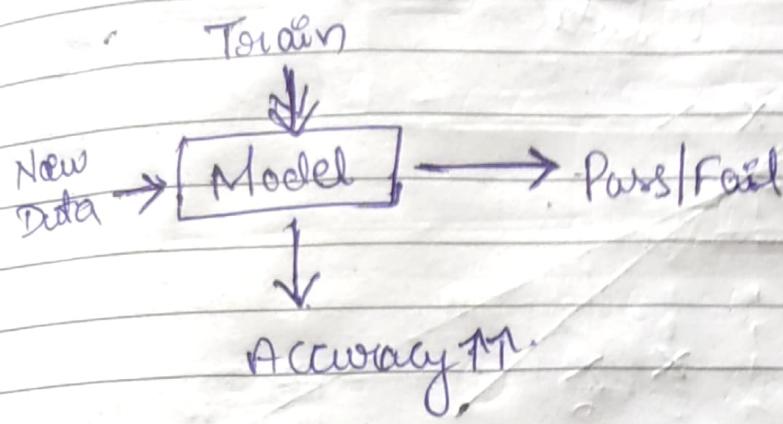
1

Pass

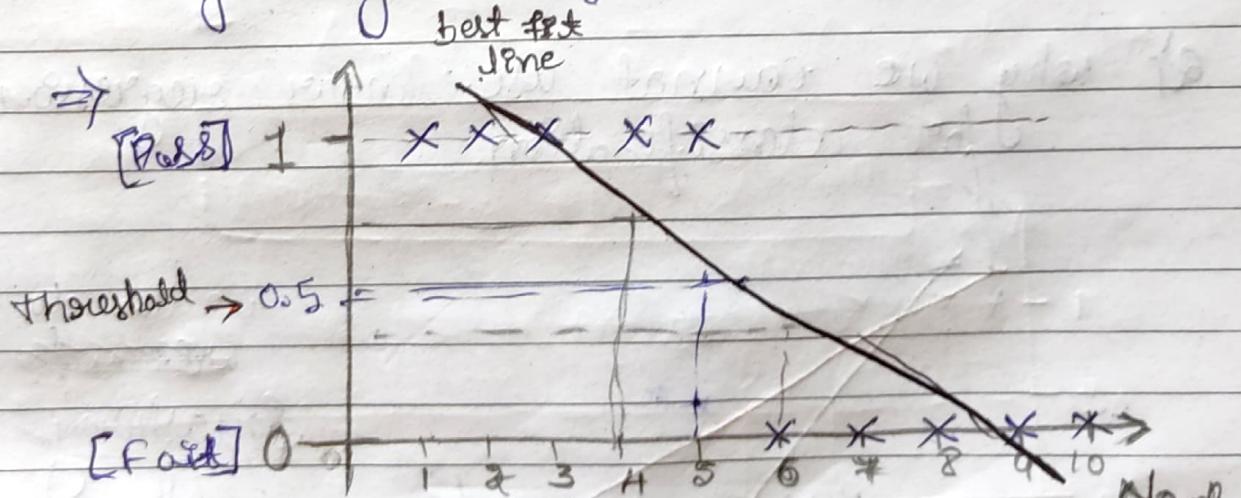
1

Pass

1



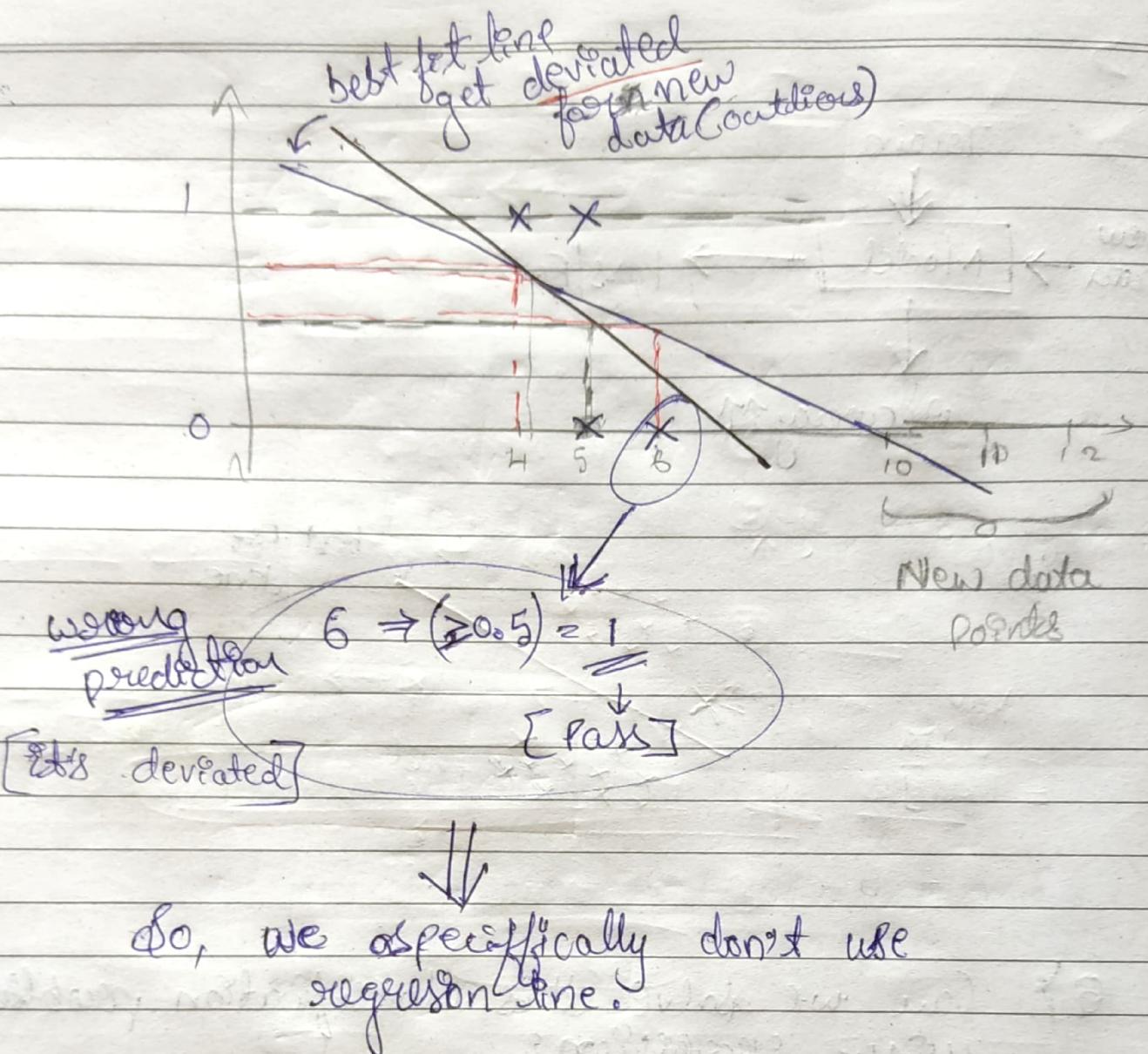
Q) Can we solve this classification problem using regression?



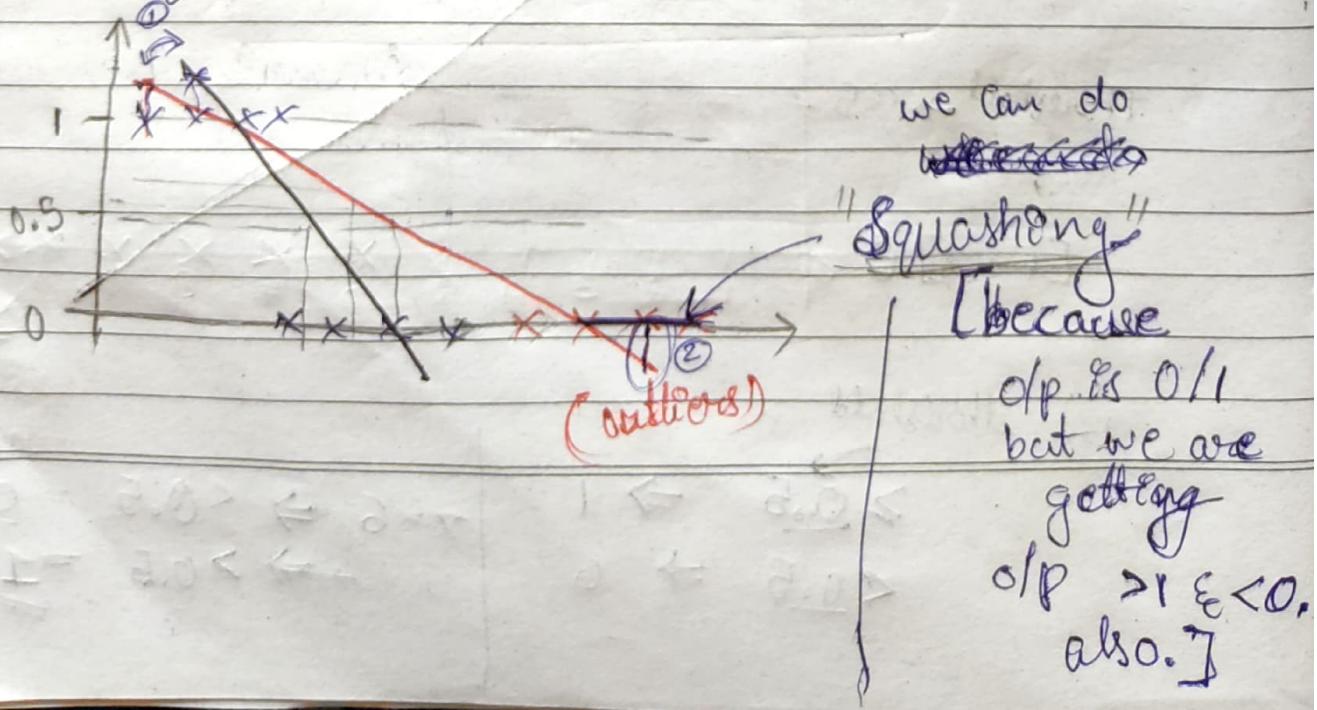
threshold

$$\begin{aligned} \geq 0.5 &\Rightarrow 1 \\ < 0.5 &\Rightarrow 0 \end{aligned}$$

$x=6 \Rightarrow < 0.5 = 0$	$x=4 \Rightarrow > 0.5 = 1$
-----------------------------	-----------------------------



Q) why we cannot use linear regression for classification.





① Best fit line changes because of outliers

↳ problem

Predict goes wrong

② The outcomes comes ( $>1$  and  $<0$ ).

solved by

logistic regression

O/P

0 to 1

has

done by  
squashing

a) How logistic regression solves Aim  
classification problem?

we have to  
squash best fit  
line

best fit line  
+ squashing



Squashing Technique

best fit line,  $\rightarrow h_{\theta}(x) = \theta_0 + \theta_1 x$

$z$

↓ "squashing"

"sigmoid activation function"

$[0, 1]$

↳ This will transform  
the line to give  
output b/w 0 & 1

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

where,  $z = \theta_0 + \theta_1 x$

Op:  $[0 \text{ to } 1]$

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

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

logistic regression

for

independent {

| dependent feature ]

activation  
function

responsible  
for

"squashing"  
best fit  
line

\* Linear regression Cost function

Cost  
function,

$$J(\theta_0, \theta_1) =$$

$$\sum_{i=1}^n \frac{[y_i - h_\theta(x)]^2}{n}$$

mean square error

$$h_\theta(x) = \theta_0 + \theta_1 x_1$$

\* Logistic regression Cost function

$$J(\theta)$$

more  
 $\downarrow$   
error  $\downarrow$   
 $\downarrow$   
less

Convex  
function

used for  
Convergence  
algorithm

gradient  
descent

global  
minimum  
[slope = 0]

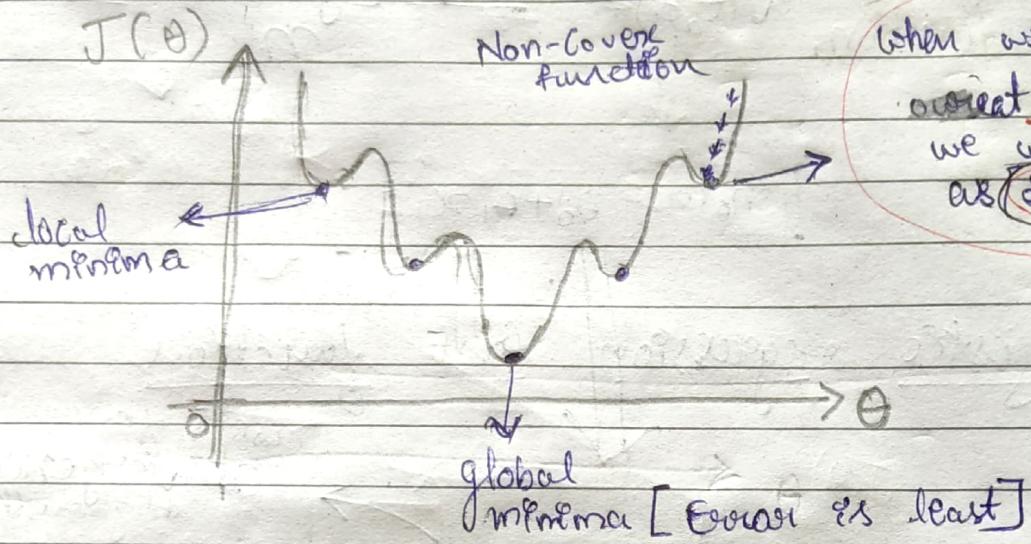
\* Aim is to reach global minima.

\* logistic regression Cost function

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

$$h_\theta(x) = \frac{1}{1 + e^{(\theta_0 + \theta_1 x)}} \quad \text{output 0 to 1}$$

\* we ~~don't~~ use this Cost function of logistic regression



when we do convergence  
at local minima  
we will get stuck  
as ( $\text{slope} = 0$ )

Non-Convex Function → because we have 1 global minima and many local minima

problem: we get "stuck during converging" as "slope become 0 at local minima's".

\* So, Cost function which ~~we~~ we use is called "log loss function".

This give only 1 global minima  
and no local minima

### \* log loss function

actual =  $y_p$ , predicted =  $\hat{y} / h_0(x)$

$$J(\theta_0, \theta_1) = -\frac{1}{m} \sum_{i=1}^m [y_i \log[h_0(x_i)] + (1-y_i) \log[1-h_0(x_i)]]$$

Cost Function

$$J(\theta_0, \theta_1) = -\frac{1}{m} \sum_{i=1}^m [y_i \log[h_0(x_i)] + (1-y_i) \log[1-h_0(x_i)]]$$

where,

$y_i$  = actual data [0/1]

$h_0(x)_i$  = predicted data.

If  $y_i = 1$ ,

$$J(\theta_0, \theta_1) = \begin{cases} -\log[h_0(x_i)], & \text{if } y=1 \\ -\log[1-h_0(x_i)], & \text{if } y=0 \end{cases}$$

Final aim:

Minimize Cost function  $J(\theta_0, \theta_1)$   
by changing  $\theta_0 \in \theta_1$

### \* Convergence algorithm

Repeat until Convergence

$$\theta_j^{(k)} = \theta_j^{(k)} - \alpha \cdot \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

{

\*  $\alpha \rightarrow$  learning rate

## \* logistic regression with regularization

parameter

feature selection

reduce overfitting

→ To include  $L_1$  and  $L_2$  regularization.

[like ridge and lasso with linear regression]

\* Cost function,

$$J(\theta_0, \theta_1) = -y \log[h_\theta(x)] - (1-y) \log[1-h_\theta(x)]$$

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

$$J(\theta_0, \theta_1) = \begin{cases} -\log[h_\theta(x)], & \text{if } y=1 \\ -\log[1-h_\theta(x)], & \text{if } y=0 \end{cases}$$

Regularization

(ridge)

① Reduce overfitting -  $L_2$  regularization

$$J(\theta_0, \theta_1) = -y \log[h_\theta(x)] - (1-y) \log[1-h_\theta(x)] + \frac{\lambda}{2} \sum_{i=1}^n (\text{slope})^2$$

$$\lambda \sum_{i=1}^n (\text{slope})^2$$

② feature selection [Lasso]

$$J(\theta_0, \theta_1) = \text{Cost function} + \lambda_1 \text{regularization}$$

$$\rightarrow \sum_{i=1}^n |\text{slope}|$$

③ Com bination (elastic net)

$$J(\theta_0, \theta_1) = \text{Cost function} + \lambda_2 \text{reg} + \lambda_1 \cdot \text{reg}$$

$$J(\theta_0, \theta_1) = -y_i \log[h_\theta(x)] - (1-y_i) \log[1-h_\theta(x)] \\ + \lambda \sum_{i=1}^n (\text{slope})^2 + \lambda \sum_{i=1}^n |\text{slope}|$$

> From sklearn.linear\_model import LogisticRegression

we have parameters

regularization  $\Rightarrow$  penalty = 'l1' / 'l2' / 'elasticnet' / None

$$A \Rightarrow B(C) = 1.0$$

how this 'C' related to 'λ'?

$$C \propto \frac{1}{\lambda}$$

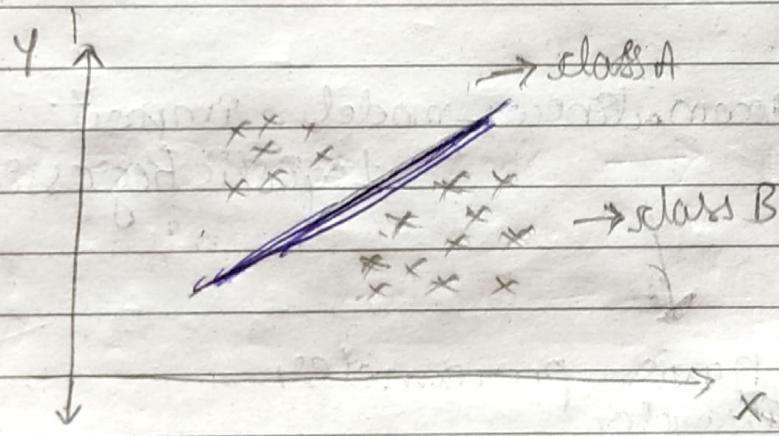
why specifically?

$\Rightarrow$  other algo we use 'C' so here also C

- > regressor. fit (x\_train, y\_train)
- > regressor. predict (x\_test)
- > regressor. score (x, y) → mean accuracy
- > regressor. predict\_proba(x) → probability estimates

## \* Performance Metrics

- ① Confusion matrix.
- ② Accuracy.
- ③ Precision
- ④ Recall
- ⑤ F-Beta Score.



\* Considering Dataset

↖ actual value

	$x_1$	$x_2$	$y$	$y'$	of dependent features	prediction
[W] Wrong prediction	-	-	0	1		
[C] Correct prediction	-	-	1	1		
C	-	-	0	0		
C	-	-	1	1		
C	-	-	1	0		
W	-	-	0	1		
W	-	-	1	0		

## ① Confusion matrix:

		1	0	Actual value [y]
		1	0	
predicted value [ŷ]	1	3	2	
	0	1	1	Correct prediction
				Matrix classification $\Rightarrow (2 \times 2)$

		1	0	Actual value
		True positive.	False positive.	
predicted value.	1	True positive.	False negative.	
	0	False positive.	True negative.	
				TP FP FN TN

actual value is "True" and prediction ~~gives~~ gives "positive" result. tells Value is 1

Model accuracy	$\Rightarrow \frac{TP + TN}{TP + FP + FN + TN}$
----------------	---

$$\frac{3+1}{3+2+1+1} = \frac{4}{7} = 57.1\%$$

## ② Accuracy

$\rightarrow$  M2  
 $\rightarrow$  adjusted  
SMA

## ③ Precision

\* Consider,

Dataset  $\rightarrow$  Imbalanced dataset

1000 datapoints  $\rightarrow$  900  $\rightarrow$  1 off  
 $\rightarrow$  100  $\rightarrow$  0

\* Default mode  $\rightarrow$  1

0	0
900	100
0	0

$$\text{accuracy} = \frac{900}{1000} = 90\%$$

what %/p but off always

baised towards 1 class

even we get good results

so accuracy

is not only

thing to tell

model is performing good

We use when FP [False positive]  
is important

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

	L	D	K
I	(TP)	(FP)	
O	(FN)	(TN)	

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

(A) Recall  $\rightarrow$  when FN is important

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

	1	0
1	TP	FP
0	FN	TN

out of all the predicted values  
how many are correctly predicted  
with ~~as~~ actual value

## ~~Use cases~~

## ① Spam Classification

Text  $\rightarrow$  Model  $\Rightarrow$  Spam/  
Not Spam

		1	0	actual
predicted	1	TP	FP	
	0	FN	TN	

$TP \leftarrow$  Scenario 1: Mail  $\xrightarrow{\text{model}}$  Spam  $\xrightarrow{\text{predict}}$  spam  $\xrightarrow{\text{generate}}$  S. Garcia

$TN \leftarrow$  mail  $\rightarrow$  not spam  $\xrightarrow{\text{model}}$  not spam  $\xrightarrow{\text{predict}}$  not spam  $\xrightarrow{\text{generate}}$  accurate scenario  
 $\rightarrow$  if mail is not spam  $\xrightarrow{\text{if it messes it is a blunder}}$

[Important]  $FP \leftarrow$  mail  $\rightarrow$  not a spam  $\xrightarrow{\text{model}}$  spam  $\xrightarrow{\text{predict}}$  spam  $\xrightarrow{\text{generate}}$  wrong scenario

$FN \leftarrow$  mail  $\rightarrow$  spam  $\xrightarrow{\text{model}}$  not a spam  $\xrightarrow{\text{predict}}$  not a spam  $\xrightarrow{\text{generate}}$  wrong scenario  
 $\rightarrow$  we can ignore spam message  
 $\rightarrow$  (not a blunder)

~~precision~~ "precision" (To calculate accuracy)  $\rightarrow$   $\frac{TP}{TP+FP}$

\* Use Case 2:

To predict whether a person has diabetes or not

Diabetes		No
Diabetes	Yes	1 0
No	0	TP FP
Diabetes		FN TN

$TP \leftarrow$  person  $\rightarrow$  diabetes  $\cap$  Correct  
 model  $\xrightarrow{\text{pred}}$  diabetes  $\cap$  scenario

$TN \leftarrow$  person  $\rightarrow$  No diabetes  $\cap$  Correct  
 model  $\rightarrow$  No diabetes  $\cap$  scenario

$FP \leftarrow$  person  $\rightarrow$  No diabetes  $\cap$  wrong  
 model  $\rightarrow$  Diabetes  $\cap$  scenario  
 (made a blunder)  
 as person taken care more

[imp]  $FN \leftarrow$  person  $\rightarrow$  Diabetes  $\cap$  wrong  
 model  $\rightarrow$  No diabetes  $\cap$  scenario  
 (blunder)  
 because person health get

$\rightarrow$  So, we use "Recall". [To calculate accuracy]

$$\frac{TP}{TP+FN}$$

problem: Tomorrow the stock market  
 we invest and  
 losses

$\Rightarrow$  we have to calculate "Recall"  
 as "FN" is important

	1	0	adult
1	TP	FP	
0	FN	TN	
1			Predict

because if our model predict  
 market "crash" (1) but  
 our model predict market  
 will "not crash" (0)  $\rightarrow$  blunders

## ⑤ F-Beta Score

$$F\text{-Beta Score} = \frac{(1+\beta^2) \cdot \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}}{\beta}$$

Scenario

- ① If FP and FN are equally important

$$\boxed{\beta = 1}$$

$$F_1\text{-Score} = \frac{(1+1) \cdot \frac{P * R}{P+R}}{2} = \frac{2 P * R}{P+R}$$

↓

"Harmonic Mean"

- ② If FP is more important than FN

$$\boxed{\beta = 0.5}$$

$$F_{0.5}\text{-Score} = \frac{(1+0.5) \cdot \frac{P * R}{P+R}}{1+0.5}$$

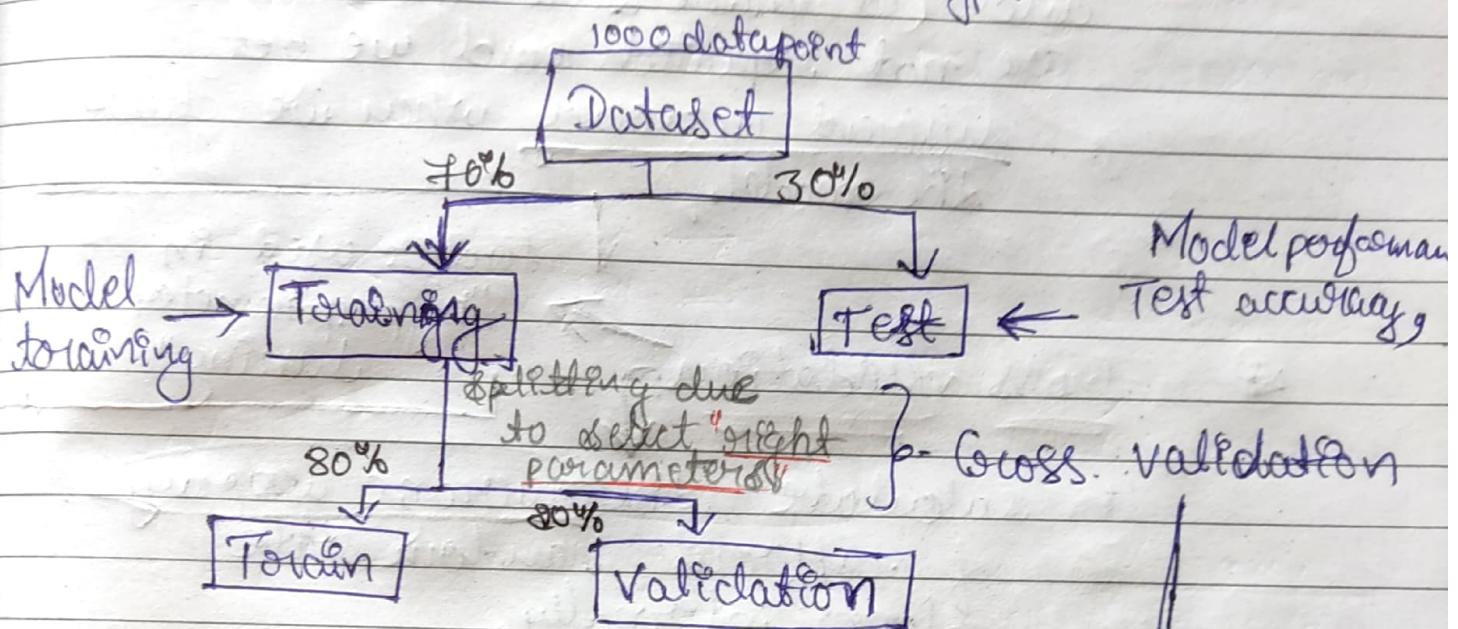
- ③ If FN is more important than FP

$$\boxed{\beta = 2}$$

$$F_2\text{-Score} = \frac{(1+4) \cdot \frac{P * R}{P+R}}{1+4}$$

## \* Cross Validation

→ Used to do hyper parameter tuning.  
2 types



\* If we do split using train-test-split,

we have a parameter  
"random\_state"

done to do  
"hyperparameter  
tuning"

↳ if we change value of this  
we get different values.

↳ if  $= 100$ , acc  $\Rightarrow 80\%$       changes  
 $= 42$ , acc  $\Rightarrow 75\%$       based on  
 $= -$ , acc  $\Rightarrow 90\%$       random state  
value

So, instead of changing values again  
and again we use cross validation

instead of change random state  
we use

In

\* ~~Cross validation~~ it will internally split divide training data into train and validation data and make sure the record are not repeated and we get different accuracy when we check the model.

\*

Cross validation

↓ further used in "hyperparameter tuning"

Calculate the average accuracy of different random state done internally.

\* Training  $\Rightarrow$   $Cv = 5$  experiment

Average of accuracy

Experiment 1

Train

Validation

Acc 1

Experiment 2

Train

Validation

Acc 2

Experiment 5

Train

Validation

Acc 5

## Types of Cross Validation

① leave one out cross validation [Loo cv]

The diagram illustrates the留一出外交叉验证 (Leave-one-out cross-validation) process:

- Training Data** (500 data points) is shown at the top.
- A bracket labeled **Expt** groups the first **1** data point and **499 data points**.
- An arrow labeled **Validation data** points from the validation set to a box below.
- An arrow labeled **As pt** points from the validation set to the same box.
- The box contains the text **use 1 data point for validation**.
- To the right, an arrow labeled **Accuracy of model** points to the text **85%**.

Enc 500  
499 datapoints | 500 |  $\Rightarrow$  Acc = 82%  
validation  
datapoint

Now, we calculate average of all the accuracy.

## \* Diss advantage:

① Time Complexity: Is huge.. for training large dataset.

② Model tends to overfit  $\Rightarrow$  Training accuracy ↑↑  
New data  $\rightarrow$  accuracy ↓↓

### ② Leave p out CV %

P value can be

anything

based on problem

hyperparameters

$$p = 10, 10\% \approx$$

Validation points will be anything. 'p' datapoints

### ③ K-fold Cross Validation:

$$K=5, n=500$$

↓

No of experiments  
to be done

$$\text{Validation data size} = \frac{n}{K} = \frac{500}{5} = 100$$

Exp 1	Valid	Train	Accuracy
	100	400	= 83%

Exp 2	Train	Valid	Train	Accuracy
	400	100	300	= 85%

Exp 5	Train	Valid	Accuracy
	400	100	= 80%

average  
of accuracy

~~more used~~

Contains equal proportion  
of validation data points.

④ ~~Stratified~~ ~~K fold CV~~

Used for "Imbalanced dataset"

Experiment

$$K = 5$$

$$\begin{aligned} n &= 500 \\ \{ &1 \rightarrow 350 \\ &0 \rightarrow 150 \end{aligned}$$

Exp1

Valid	Train
100	400

- = Acc1

Validation data

Contain Equal [0's = 1's]

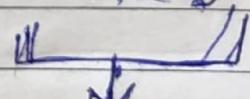
no of data points from  
class

$$\begin{array}{l} \{ 50 \rightarrow 1's \\ 50 \rightarrow 0's \\ \hline 100 \end{array}$$

Exp5

Train	Valid
-------	-------

- = Acc5



Average  
of accuracy

\* Most of the time we perform K fold validation.

\* We use ~~cross validation~~ inside internally in the ~~hyperparameter tuning~~ to select the best parameters

## \* Hyperparameter tuning with cross Validation

↳ Finding the "best parameter" while training the models.

### Types

- ① GridSearch CV.
- ② RandomizedSearch CV.

→ Few parameters of Logistic Regression.

→ penalty { 'l1', 'l2', 'elasticnet', None }

> solver { 'lbfgs', 'liblinear', 'newton-cg', 'newton-cholesky',  
'sag', 'saga' }

Q) Now, Task is to know which parameter do we use and which combination



So, we use hyperparameter tuning

① Grid Search CV [ GridSearch + Cross Validation ]

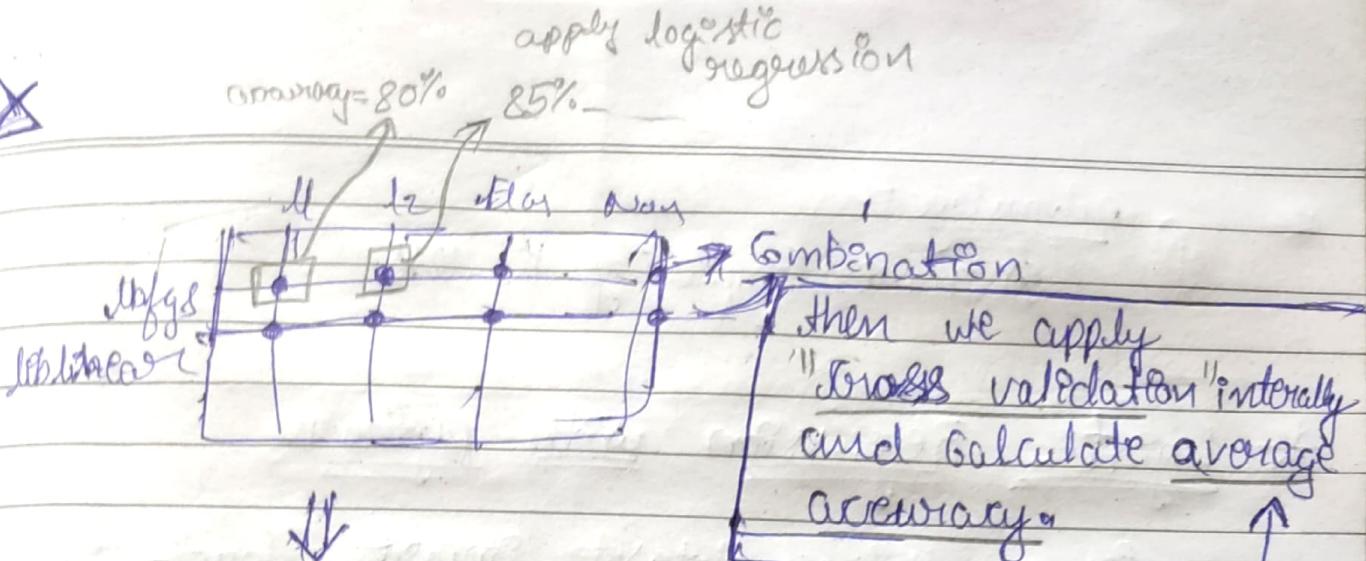
using ↓

we try to

find out!

best combination  
of parameters

all by taking  
all possible combinations ↘



\* Here, using `grid_search_cv` we are going to find out all possible Combination with their accuracy and select the best combination.

### ② Randomized Search CV

+ Disadvantages:

- ① Time Complexity increases with large dataset.  
for. Training the model.

### ③ Randomized Search CV

$$n\_iter = 100$$

stratified k-fold cv  
+  $k \leq 5$

→ Here, we add. iteration parameter which tell randomly select perform

10 different Combination

← 10 iteration

+

Cross Validation

etc

Here, 10 different Combination +  $cv=5 \Rightarrow 50$

↓

select the best parameter

## Advantages

① Time complexity is less

## \* Logistic Regression with Python

> from sklearn.datasets import load\_iris

> dataset = ~~load\_iris()~~

↓  
dictionary contains information about dataset

> df = pd.DataFrame(dataset['data'], columns=dataset['feature\_names'])

① Independent and dependent feature  
(df.iloc[:, :-1])

② Train Test Split

~~Classification~~ [model selection]

③ Standard Scaler

④ Train model

> from sklearn.linear\_model import LogisticRegression  
> classifier.fit(x\_train, y\_train)

> y\_pred = classifier.predict(x\_test)

⑤ Prediction

⑥ Accuracy

> from sklearn.metrics import confusion\_matrix, accuracy\_score,

> accuracy-score (y-test, y-pred)

> classification report (y-test, y-pred) ↴

for F1 score,  
recall, precision

⑥ > from sklearn.model\_selection import  
KFold

> CV = K-Fold (n\_splits=5)

Here, >

import cross\_val\_score

estimator

> cv\_scores = cross\_val\_score (classification, x\_train, y\_train)

& scoring = 'accuracy')

CV = ~~K~~ V



output 5  
different accuracy

logistic  
regression  
mode

taking cross  
validation for  
to split

data within  
given data

> final\_score = np.mean (cv\_scores)

## \* Hyper parameter tuning with Python

> from sklearn.model\_selection import  
GridSearchCV

> classifier = LogisticRegression()

> self = GridSearchCV(classifier, param\_grid = parameters,  
cv = 5)

↑  
internally we use  
KFold cross  
validation

> self.fit(x\_train, y\_train)

split data into  
train & validation  
then training

> self.best\_params\_

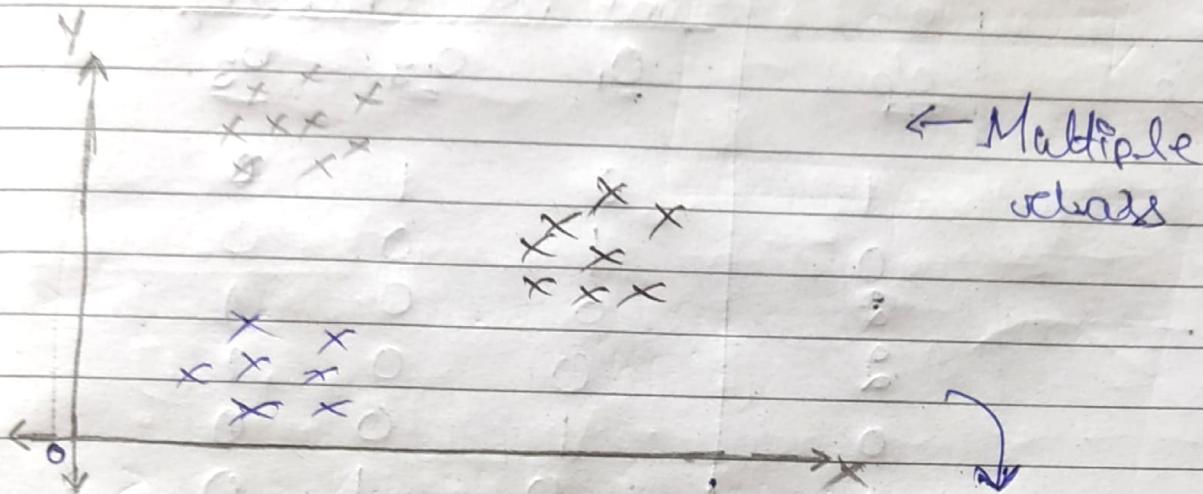
calculate accuracy & ~~precision~~  
model-selection

> import RandomizedSearchCV.

estimation

> clf = RandomizedSearchCV(LogisticRegression(),  
param\_distributions = parameters,  
(cv=5, n\_iter=20))

# \* Logistic Regression for multiclass Multiclass classification



- 2 Technique to classify
- ① OVR  $\rightarrow$  One versus rest
  - ② Multinomial  $\rightarrow$
- we cannot classify using a single best fit line

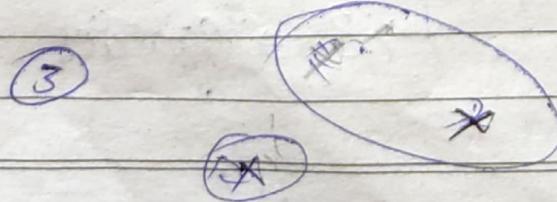
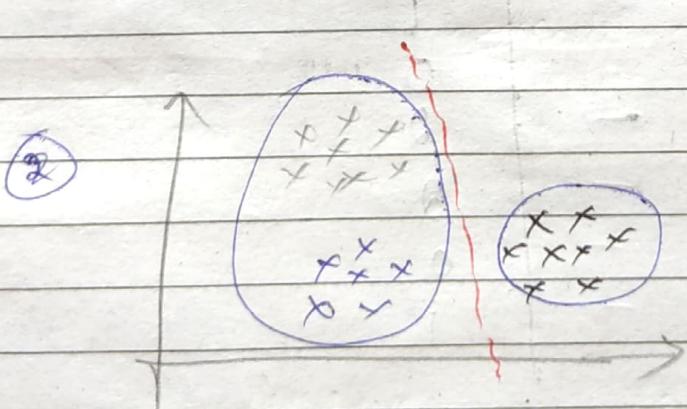
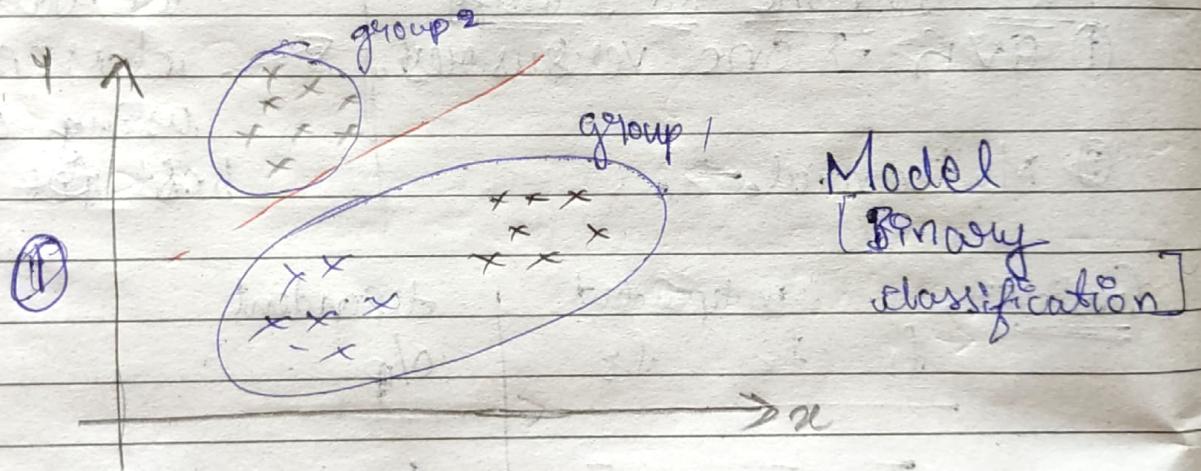
Consider,

	Independent			Dependent O/P
	$f_1$	$f_2$	$f_3$	
-	-	-	-	1
-	-	-	-	0
-	-	-	-	2
-	-	-	-	3
-	-	-	-	0
-	-	-	-	1

$\rightarrow$  3 categories

① OVR [One versus rest]:

	$\text{o}^*$	$\text{o/p}$	<u>One hot encoding</u>		
			$\text{o}_1$	$\text{o}_2$	$\text{o}_3$
1	1	1	0	1	0
2	0	0	1	0	0
3	0	0	0	0	1
4	0	0	0	1	0
5	1	1	0	1	0



M1 model  $\rightarrow$  O/P =  $f_1 + f_2 + f_3$ ;  $O/P = \Theta_1$

XES

$\downarrow$   
Binary  
classification

M2 model  $\rightarrow$  O/P =  $f_1 - f_2 - f_3$ ;  $O/P = \Theta_2$

M3 model  $\rightarrow$  O/P =  $f_1 \cdot f_2 \cdot f_3$ ;  $O/P = \Theta_3$

\* Works similarly if we have many classes.

\* New Test data  $\rightarrow [M_1 \ M_2 \ M_3]$

$$O/P = [0.25, 0.25, 0.5] \quad (\sum = 1)$$

$\downarrow$   
Highest probability

$$O/P =$$

$$0.5 \rightarrow \underline{\underline{C_3}}$$

> In ~~the~~ Logistic regression we have parameters

$\downarrow$

default

> multiclass : { "auto", "over", "multinomial" }

$\Rightarrow$   
we get  
output

$\downarrow$   
we get  
probability

~~M1 model~~  $\rightarrow \text{O/P} = f_1 + f_2 + f_3 ; \boxed{\text{O/P} = \Theta_1}$

~~Binary classification~~

$M_2$  model  $\rightarrow \text{O/P} = f_1 + f_2 + f_3 ; \boxed{\text{O/P} = \Theta_2}$

$M_3$  model  $\rightarrow \text{O/P} = f_1 + f_2 + f_3 ; \boxed{\text{O/P} = \Theta_3}$

\* Works similarly if we have many classes

\* New Test data  $\rightarrow [M_1 \ M_2 \ M_3]$

$O/P : [0.25, 0.25, 0.5] \quad \sum = 1$

$\downarrow$   
highest probability

$O/P : 0.5 \rightarrow [2]$

> In ~~logistic regression~~ we have parameters

$\downarrow$   
default

> multiclass : { "auto", "over", "multinomial" }  $\rightarrow$

we get specific output.

$\downarrow$   
we get probability

\* Usually, in logistic regression we get less accuracy as points will be scattered.

\* Python it's same

①  $x, y$  [segregate dependent & independent]

② train test split

Hyper  
parameters  
Tuning →

③ Model training

> classifier = Logistic Regression(

multi\_class = 'ovr',

solver = 'lbfgs')

optimization

④ prediction

→ model.predict\_proba(\*test)

⑤

accuracy, confusion matrix

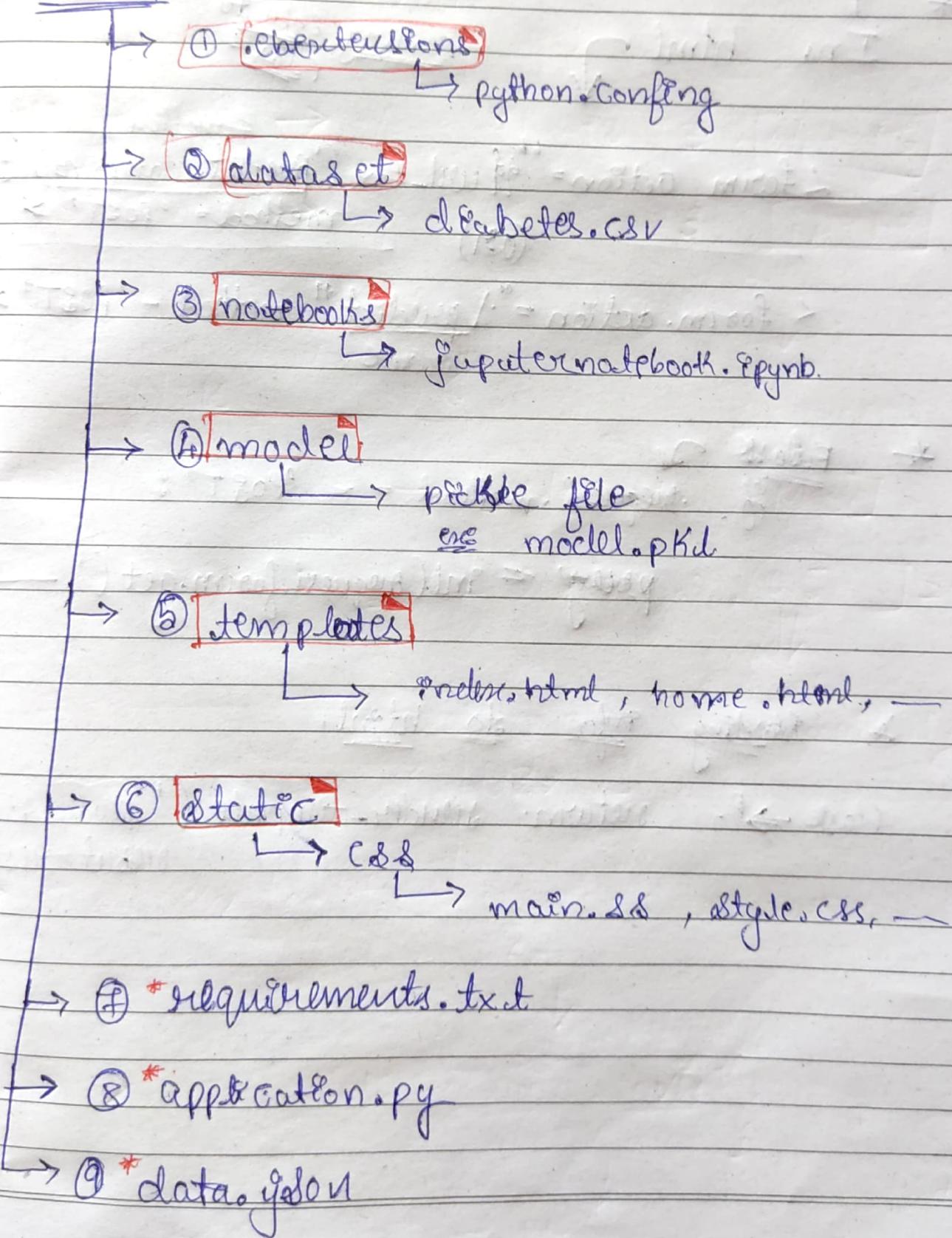
> from sklearn.datasets import make\_classification

>  $x, y = \text{make\_classification}$  (n\_samples = 1000,  
n\_features = 10, n\_informative = 5,  
n\_redundant = 5, n\_class = 3,  
random\_state = 1)

\* WSGI → Web Server Gateway Interface

\* End-to-End project:

Folders



→ "dataset/dataset.csv" ← path  
go back one step

\* In a html ↓ sending info to flask

< form action = "if \_\_name\_\_ == 'predict\_datapoint':  
 method = 'POST'>  
(0.91) function

< form action = "predict" method = "POST">  
route

\* flask if request.method == 'POST':

predy = int(request.form.get(''))

→ sending info to html.

flask → return render\_template('predict.html',  
 result=result)