

ALGORITHMS: \rightarrow One input feature & One output feature

SIMPLE LINEAR REGRESSION: (It is a Supervised ML)

As we know the Supervised ML classified into two types of problems

- 1) Regression \rightarrow o/p feature is continuous in nature
- 2) Classification \rightarrow o/p feature is Binary / Multi-class ^{categories} classification

Let's assume a simple dataset with two features (one dependent & one independent)

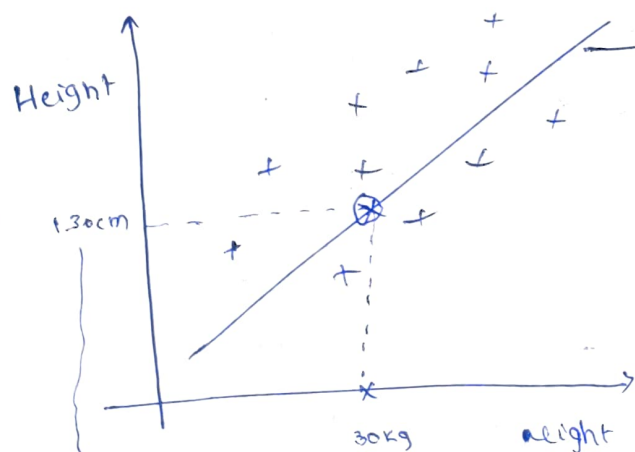
Independent Weight	Dependent height
74	170
80	180
90	190
75	175.5
...	...

\rightarrow If our client wants to predict height based on weight.

\rightarrow Our model needs to find height based on weight so, Weight is Independent feature (input feature), height is dependent feature (o/p or target feature)

- 1) Initially we train our model with the data & we try to predict the height using weights using test data. (Test data only have weights data, based on weight's our model will predict height).

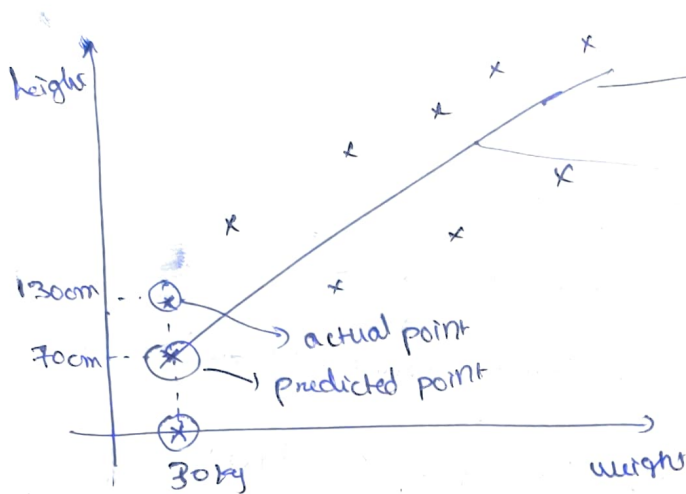
GEOMETRICAL INTUITION BEHIND SIMPLE LINEAR REGRESSION: (It will tell you how our model will find best fit line based on training data)



Linear regⁿ find best fit line
Our Model will ~~create~~ best fit line based on training data

Then based on best fit line we will find the height values for new weight values like

Our Model predicts height = 1.30cm for input weight of 30kg

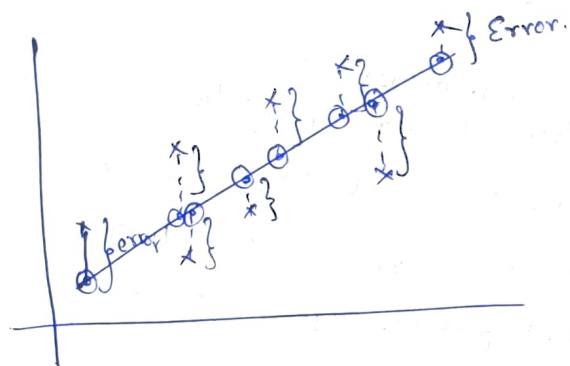


after training the ML model, it creates a Best fit line

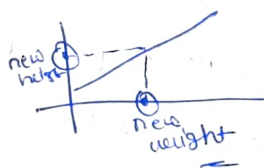
points on best fit line are classified as predicted points

distance (sum of error) ~~sum of~~ b/w ~~the~~ actual data

Best fit line is created in such a way that the ~~error~~ distance b/w actual data point and predicted data point is minimum.

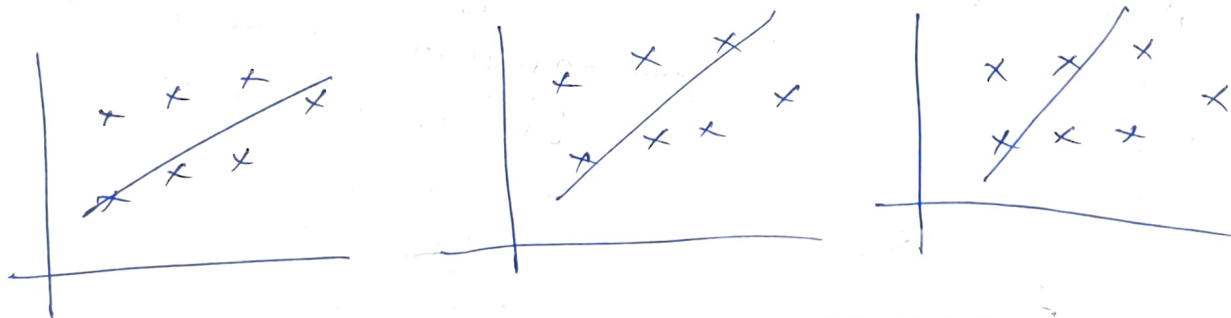


purpose of best fit line → Used to predict height values based on new weight values



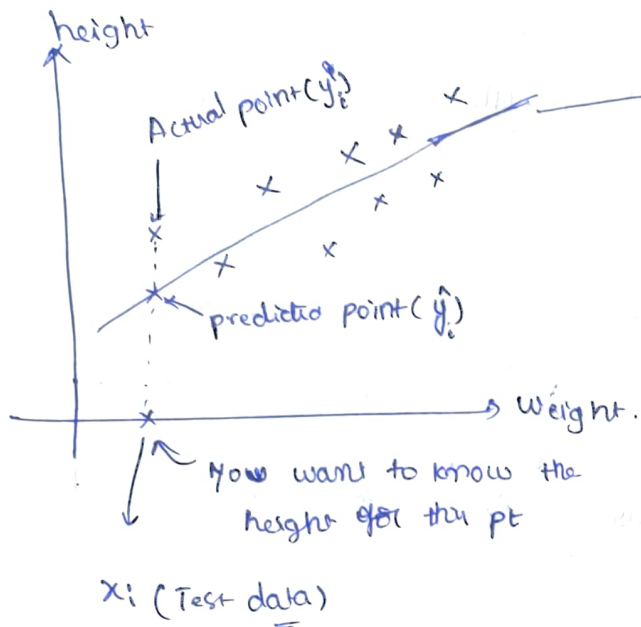
How exactly ~~can~~ the best fit line will be created → will see later.

To be more precise



3 Best fit line's for same data points,

whichever best fit line has less sum of error b/w actual & predicted values it will be considered.



The equation to represent straight-line is

$$y = mx + c \text{ or } y = wx + b$$

Slope \rightarrow data points \rightarrow intercept

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

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

All are same

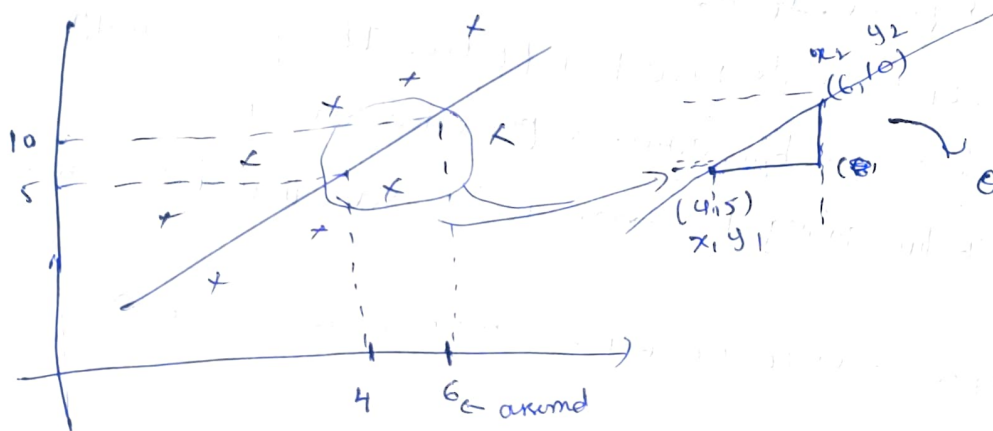
Researchers are using this equation.

$\theta_0 \Rightarrow$ intercept

$\theta_1 \Rightarrow$ slope & co-efficient

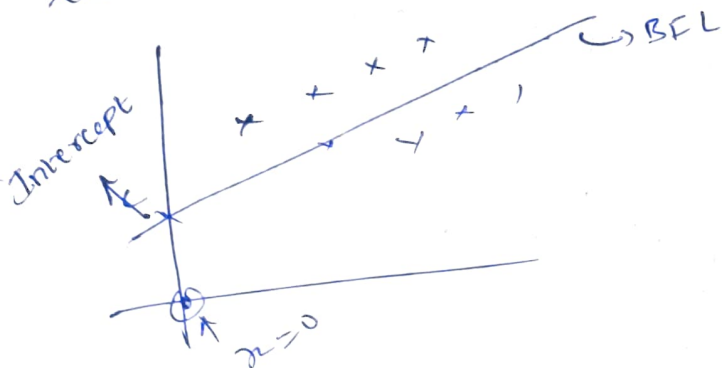
Error $\Rightarrow (y_i^* - \hat{y}_i)$
 \nearrow actual
 \searrow predicted

Slope & co-efficient: with the movement on x-axis how ~~value changing~~ much movement we have in y-axis?



$$\begin{aligned} \theta_1 &= \frac{y_2 - y_1}{x_2 - x_1} \\ &= \frac{10 - 5}{6 - 4} \\ &= \frac{5}{2} = (2.5) \end{aligned}$$

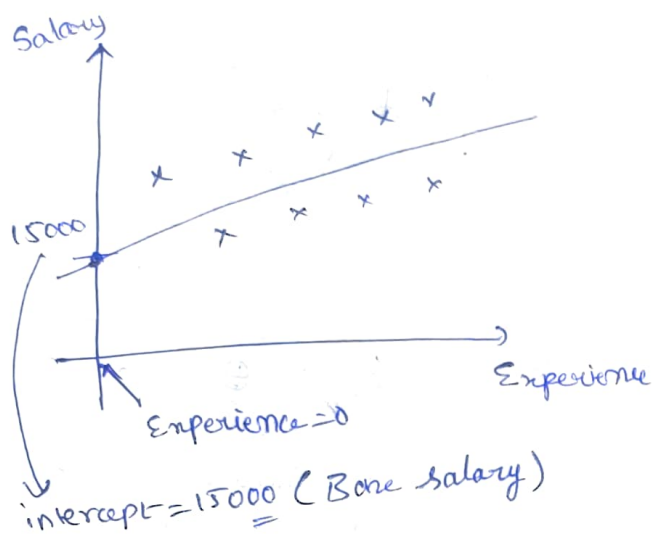
Intercept: if our x-axis value = 0, then the value of y is represented as intercept



$$\begin{aligned} y &= mx + c \\ y &= m(0) + c \end{aligned}$$

$$y = c$$

Small ex: When we plot the graph experience vs salary. even though we joined as a fresher we have some basic salary right



* ~~$h = \theta_0 +$~~ $\boxed{h_0(x) = \theta_0 + \theta_1 x}$ ← data points

So, Based on θ_0 & θ_1 , we will fit & create the best fit line.

Initially, we initiate random values for θ_0 & θ_1 , & we will create a best fit line & based on best fit line we will calculate the ~~sum of~~ error b/w actual & predicted point. If the error is high then again we change the value of θ_0 & θ_1 to create new best fit line. This process will repeat ~~until~~ we get less error b/w actual & predicted data point. Whatever best fit line giving less error, we consider that line as our best fit line.

(THE WHOLE PROCESS WILL BE EXPLAINED IN DETAIL LATER)

↓
(This process is called optimization)

So, In order to derive the optimization we need to derive the

Cost function:

$$J(\theta_0, \theta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - h_0(x)) ^2 \Rightarrow \text{Mean squared error}$$

\downarrow no. of data points. $\underbrace{\text{actual} - \text{predicted}}_{= h_0(x) - y_i}$

Eqn tells you that, Our Aim is to find the best fit line in such a way that our mean squared error should be minimized.

Sum of Squares of Actual - predicted values.

n = no. of data points

y_i = Actual values

$h_0(x)$ = predicted value.

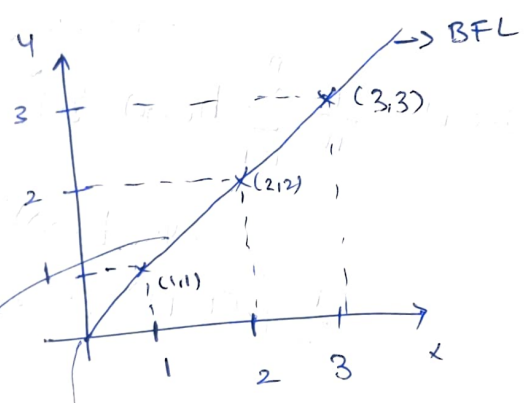
Final Aim : \rightarrow To get Best fit line
Minimize $J(\theta_0, \theta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - h_0(x))^2$

Minimize $J(w, b) = \frac{1}{n} \sum_{i=1}^m (y_i - \hat{y}_i)^2$

Optimization : \rightarrow How do we minimize the cost function?

Dataset :

$x \rightarrow$ independent	$y \rightarrow$ dependent
1	1
2	2
3	3



As we can observe there is no error b/w actual & \hat{y}_i because the BFL is passed exactly through all points

$h_0(x) = \theta_0 + \theta_1 x$ or $y = mx + c$

Here our intercept is zero. \rightarrow i.e. $\theta_0 = 0$ or $c = 0$

$h_0(x) = \theta_1 x$

If we assume $\theta_1 = 1$

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

$h_\theta(x) = 1$

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

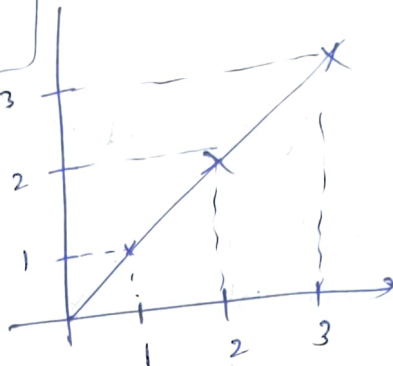
$h_\theta(x) = 2$

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

$= 3$

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

\Rightarrow If we plot



After creating BFL, we need to find cost function

cost f^n : ~~$\theta_0 = 0$~~ $\xrightarrow{\text{Actual.}}$ $\boxed{\theta_0 = 0} \Rightarrow \boxed{h_\theta(x) = \theta_1 x_1}$

error b/w actual & predicted

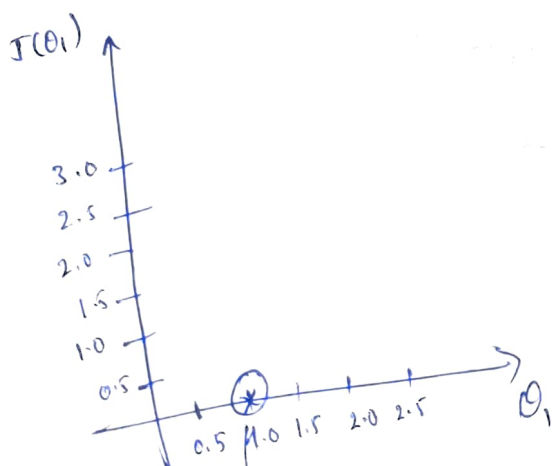
$$J(\theta_1) = \frac{1}{n} \sum_{i=1}^n (h_\theta(x_i) - y_i)^2$$

$$(y_i - h_\theta(x_i))^2$$

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

$= 0$

Creating another graph that plotting $J(\theta_1)$ vs θ_1



for $\theta_1 = 1.0$, $J(\theta_1) = 0$ ✓

* Now assume $\theta_1 = 0.5$

at $x=1$

$$h_0(x) = 0 + 0.5(1)$$

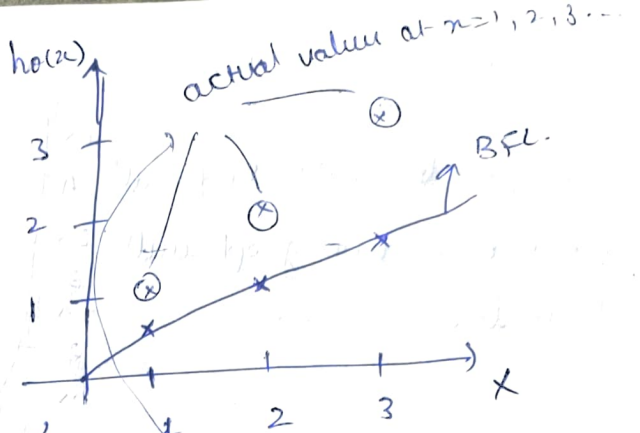
$$= 0.5$$

at $x=2$ $h_0(x) = 0 + 0.5(2)$

$$= 1$$

at $x=3$ $h_0(x) = 0 + 0.5(3)$

$$= 1.5$$



Now the BFL is changed

Now cost f^n :

$$\Rightarrow \frac{1}{n} \sum_{i=1}^n (y_i - h_0(x))^2$$

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

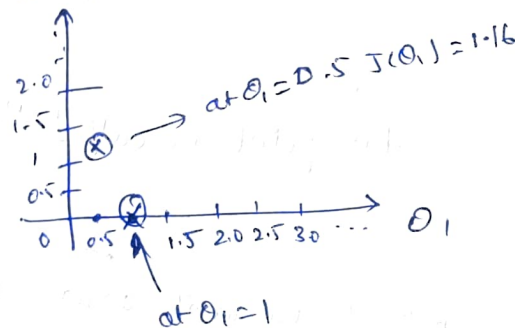
$$= \frac{1}{3} (0.25 + 1 + 2.25)$$

$J(\theta_1) = 1.16$ when $\theta_1 = 0.5$

x	y	$h_0(x)$
1	1	0.5
2	2	1
3	3	1.5

\downarrow actual \downarrow predicted

$J(\theta_1)$



cost f^n :

Now assume $\theta_1 = 0$

at $x=1$ $h_0(x) = 0 + 0(1)$

$$= 0$$

at $x=2$ $h_0(x) = 0 + 0(2)$

$$= 0$$

at $x=3$ $h_0(x) = 0$

x	y	$h_0(x)$
1	1	0
2	2	0
3	3	0

cost f^n :

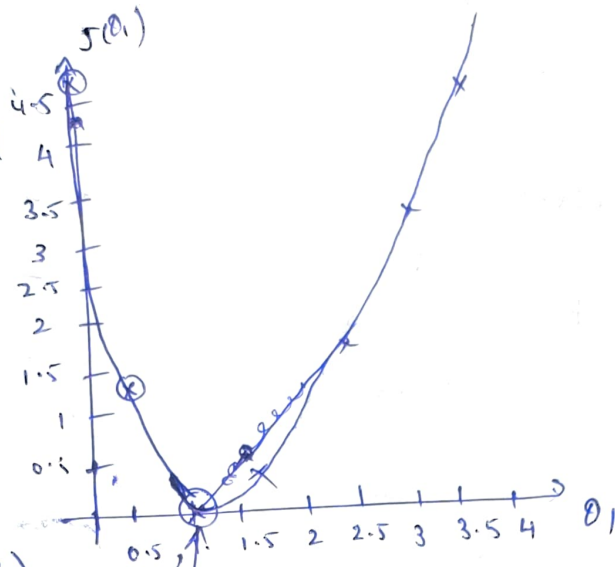
$$= \frac{1}{3} ((1-0)^2 + (2-0)^2 + (3-0)^2)$$

$$= \frac{1}{3} (1 + 4 + 9)$$

$$= 14/3 \Rightarrow 4.66$$

at $\theta_1 = 0$, $J(\theta_1) = 4.66$..

if we plot ~~the~~ $J(\theta_1)$ vs θ_1 for all values of θ_1 then graph will look like. \rightarrow



We will get a parabola like this. (U-shape)
GRADIENT DESCENT CURVE

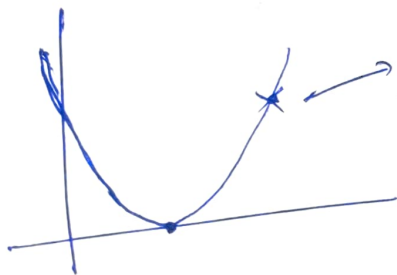
At this point only our cost f^n is zero
i.e error is zero

So, This is point where we get our best fit line for the training data
at error is less b/w actual & predict points.

This point is called "Global minima" (cost f^n is less)
error is less

So,
~~Initially~~ we won't assign the multiple θ_1 values but we will change
either increment or decrement the θ_1 values based on our requirement

i.e



for ex; if you are here then to find the global minima then we need to decrement the θ_1 value right so like that we need to develop our algorithm.

it is CONVERGENCE ALGORITHM

Convergence algorithm Moving towards global minima

In the convergence algorithm, we optimize the change of θ_1

Algorithm:

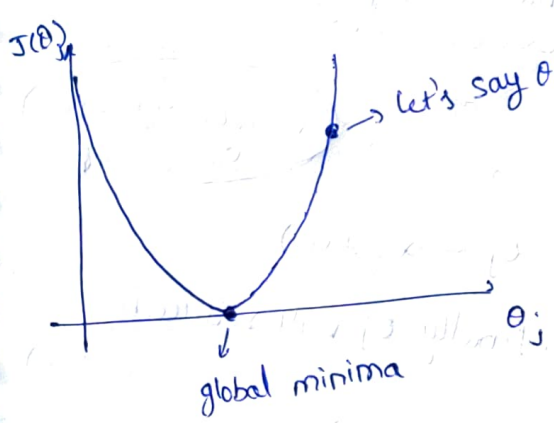
Repeat until convergence

{ $j = 0, 1$

$$\theta_j = \theta_j - \underset{\substack{\uparrow \\ \text{Learning rate}}}{\alpha} \frac{\partial}{\partial \theta_j} (J(\theta_j))$$

}

Let's consider

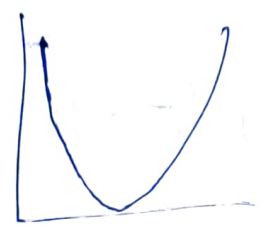
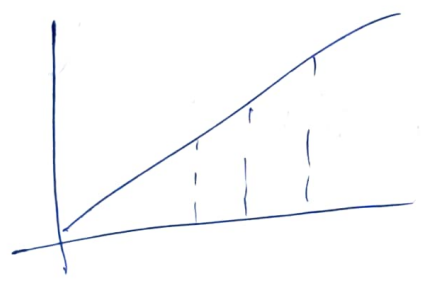


→ So now we need to move towards global minima.
So, How we can move θ_j value towards global minima.

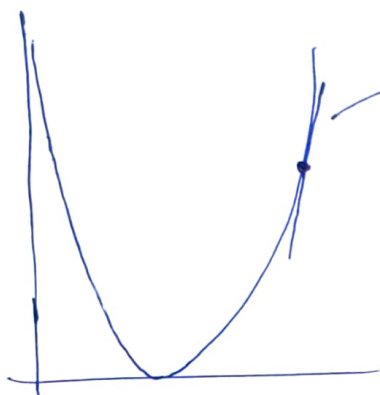
(By seeing graph we can say that we need to reduce the value of θ_j to reach global minima but for algorithm we need to find ~~derivative~~ ^{slope} at that particular point)

As the curve is parabola, we need to find ~~derivative~~ ^{slope} to decide whether we need to \uparrow or \downarrow the θ_j value

Ex: for straight line slope is same at every point but for parabola \uparrow the slope will change at each point.



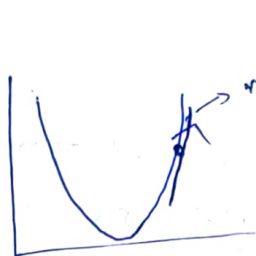
To find the slope value of parabola curve, we need to find derivative at that particular point



if you want to find derivative at particular point then we need to draw a tangent then we need to find whether it is +ve or -ve slope

$$\text{derivate} = \frac{\partial}{\partial \theta_j} J(\theta_j)$$

as a notation, if our tangent is ~~at that~~ showing up, then it has +ve slope



right side of tangent showing upwards, i.e.

it has +ve slope



so, new old

$$\theta_j = \theta_j - \alpha (+ve)$$

α always +ve
Say $\alpha = 0.01$

So, finally θ_j will reduce by $\alpha (+ve)$

$$\theta_{j, \text{new}} < \theta_{j, \text{old}} \checkmark$$



if the tangent right side is downwards then we get -ve slope

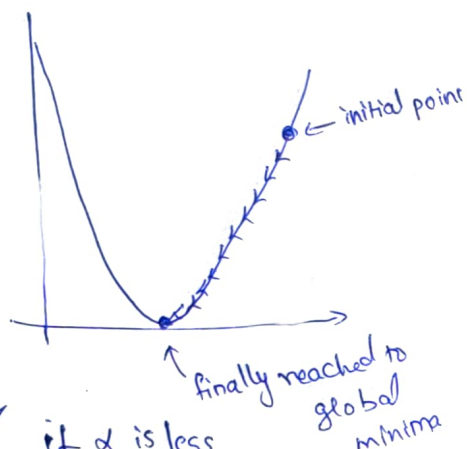
so, $\theta_j = \theta_j - \alpha (-ve)$ final θ_j value will increase.

Finally,

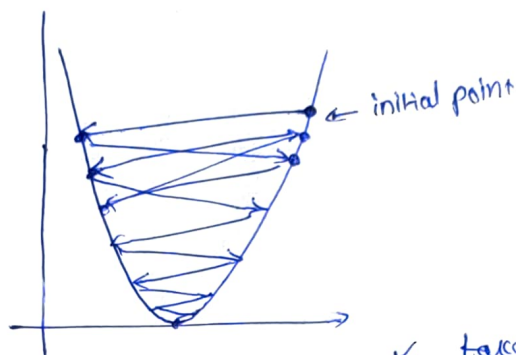
After reaching global minima, Slope = 0 because of horizontal line

Then $\theta_j = \theta_j - \alpha(0) \Rightarrow \theta_j$ (only θ_j will be there)

$\alpha \rightarrow$ It will be very small, It is used to converge the θ_j value quickly,
if it is large value then it will take so much of time to converge.



★ ★ if α is less
★



if α is large then it ~~will~~ [✓] take
so much of time to converge

★ Sometimes it won't reach ^{very} global
minima if α is large

NOTE: ★ ★

if α is less, converge ~~not~~ time is less ★

if α is very large, convergence takes very large amount of time.

Once we reach the global minima, then we consider the θ_1, θ_0 values at ^{Global} _{minim}

& we create our best fit line. ★ ★
★