

In statistics we call it ordinary least square.....

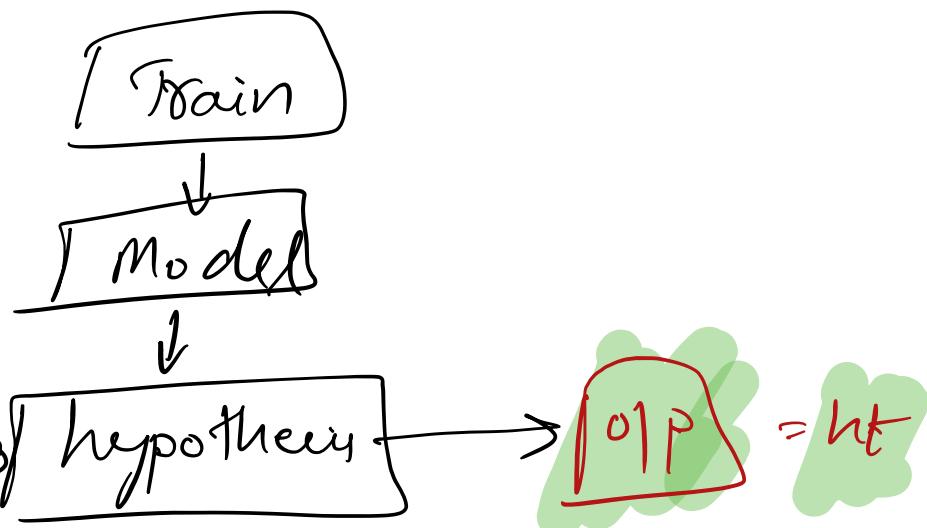
→ Used in Regression (Supervised Learning)

### Mathematical Intuition:-

Ex:-

wt	ht (cm)
73	150
63	155
78	170
82	165

Train data

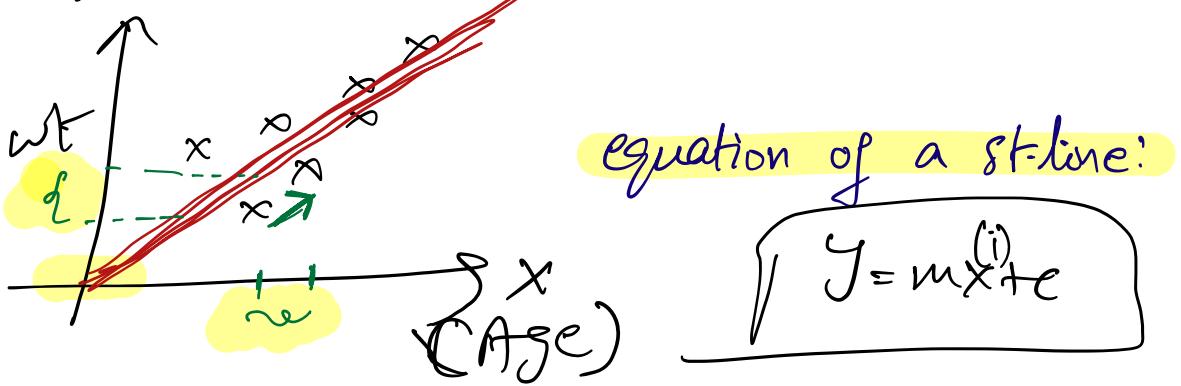


Note:

ht

is a continuous variable

Krish Naik



it can also be written in the form of

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

Aim:

we should find out the best fit line, in such way that the distance b/w data points and the predicted values is min.

[we take the sum of all these distances, then it should be minimal]

$\Rightarrow$  Now, there might be many such lines which can be drawn across the points...

$\rightarrow$  so, for this confusion, we start at point and we go in a path such that we reach (or) get the best-fit line at the end of all the data points.

$$\text{intercept} = \theta_0$$

where is it touching the y-axis

$$\theta_1 = \text{slope of } \text{coefficient}$$

→ So to guide us, what path (or) what way we should follow we have smtg called as **Cost function**.

### Cost function:-

the dist formula b/w the predicted (and) the real data

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

↳  $\theta_0 + \theta_1 x$

*hypothesis function*

$m$  = all the points present

we need to keep on updating  $\theta_0$  and  $\theta_1$  unless we get the best fit line

so, (i) will run from  $i=1$  till  $m$

Note: we can just divide by 'm' but we are taking 1 so that it will help us in derivation.

The above formula  
is also called as  
**Squared Error Function**

What we need to solve:

minimize  $\theta_0, \theta_1$

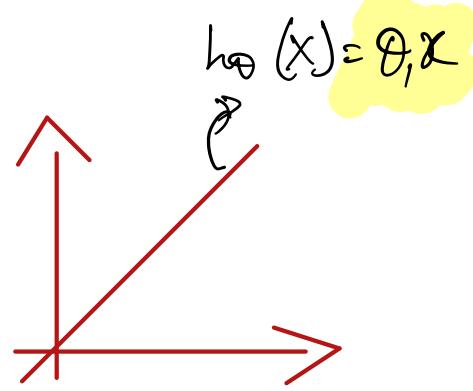
$$\frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

minimize  
 $\theta_0 \theta_1$

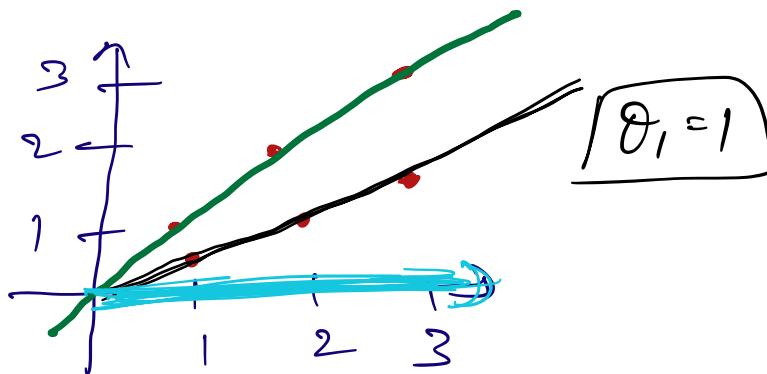
$$J(\theta_0, \theta_1)$$

$$\textcircled{2} \quad h_{\theta}(x) = \theta_0 + \theta_1 x \quad \text{if } \theta_0 = 0$$

intercept is passing through origin.



ex!  $(1,1) (2,2) (3,3)$  — data points



Scenario = 1

best fit line is given by green line

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

$$= \frac{1}{2m} \left[ \theta_0 + \theta_1 x_i + (1 \cdot 1) - 1 \right]^2 + (2-2)^2 + (3-3)^2$$

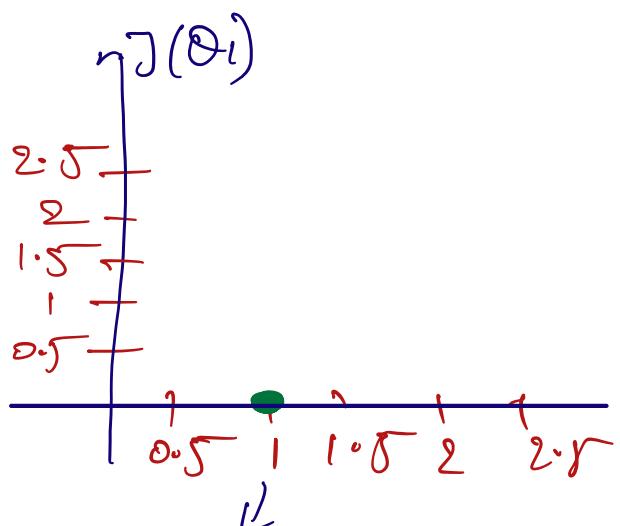
$$\boxed{J(\theta_1) = 0}$$

Cost func  
graph

Scenario = 2

$$\boxed{\theta_1 = 0.5}$$

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



$$\begin{aligned}
 &= 0 + 0.5(1) \\
 &= 0.5
 \end{aligned}
 \quad \left. \begin{array}{l} \text{cost funct}^n \\ \Rightarrow \frac{1}{2m} [(0.5-1)^2 + (1-2)^2 + (1.5-3)^2] \end{array} \right\}$$

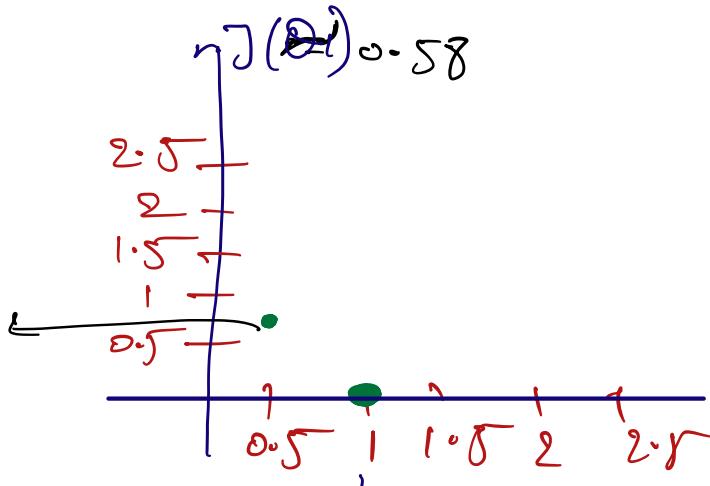
My  $x=2 \Rightarrow 1$   
 $x=3 \Rightarrow 1.5$

predicted  
actual

Plotting the best fit line in the above graph with black color.

$$= \frac{1}{2 \times 3} [0.25 + 1 + 2.25]$$

(our second point)



(our first point)

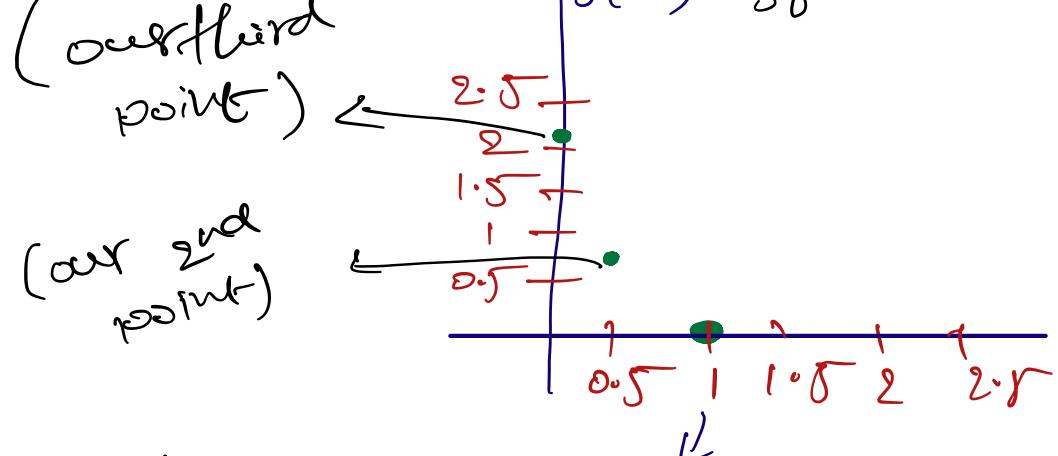
Scenario = 3

$\boxed{\theta_0, \theta_1}$

Then in the plot our best fit line is given by blue color, calculating cost function.

$$\begin{aligned}
 J(\theta_0, \theta_1) &= \frac{1}{2 \times 3} [(0.7)^2 + (0-2)^2 + (0-3)^2] \\
 &= \frac{1}{6} (1 + 4 + 9) \approx 2.3
 \end{aligned}$$

$J(\theta_0, \theta_1) \approx 0.58$



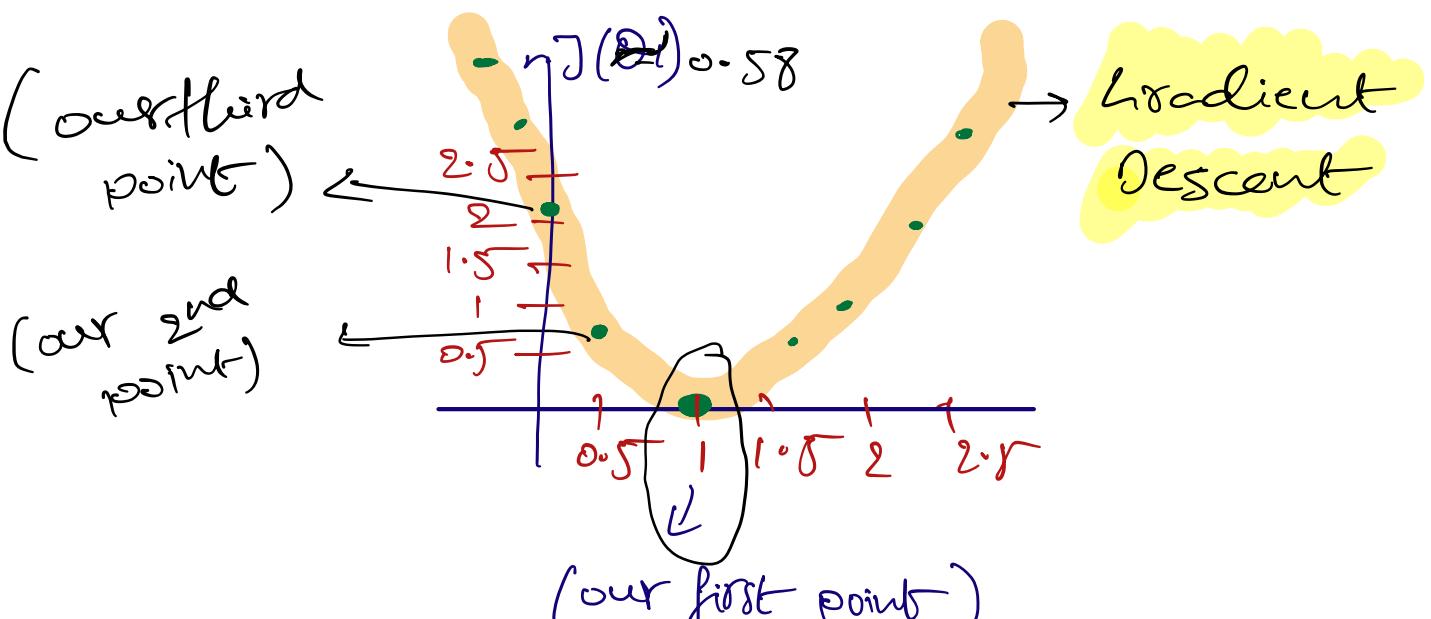
note:

⇒ So, for diff values of  $\theta_1$  we get diff points on the graph of the cost function.

⇒ Joining all those points give us a kind of curve, which we call as

## Gradient Descent

→ This GDF play a main role in getting the right  $\theta_1$  value



The point on X-axis ( $x=1$ ) is called as Global minima

→ Our total concentration will be on getting this global minima, where we get the best fit line. If we see, our  $\boxed{\partial J = 0}$  gave us the best fit line.

Note: we have actually assumed  $\boxed{\partial J = 0}$  for the first time and we got the best-fit directly....

→ Else, we pick some value for  $\boxed{\partial J}$  and check do we get best fit line or not  
→ then we update our Value of  $\boxed{\partial J}$  accordingly. In a way, that our slope in the Gradient Descent will give us the Global minima

→ to assume diff values of  $\boxed{\partial J}$  we can randomly pick values, for this reason we take the help of learning rate and try to move towards our global minima

## Conversion Algorithm:

→ after getting a particular  $\theta$  value we just need to get close to our global minima. So, in this process we are actually leaving out few values of  $\theta$ , which may not guide us to global minima

[Life got simple :)]

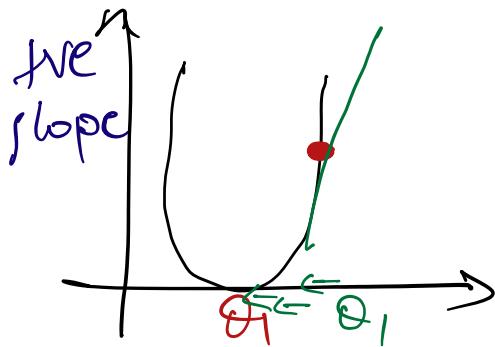
repeat until convergence  $\nexists$  Suppose a while loop

$$\theta_j := \theta_j - \alpha \left[ \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \right]$$

}

learning rate

- we use this formula to find global  
minima.
- derivative to find slope



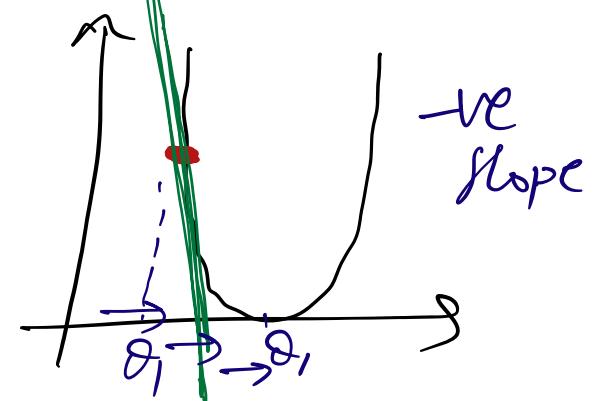
right hand side of slope is pointing top.

$$\theta_1 : \theta_1 - \alpha (\text{tve})$$

$\downarrow$  *one*

$\downarrow$  *tve number*

So, after few iterations, we come to our global minimum



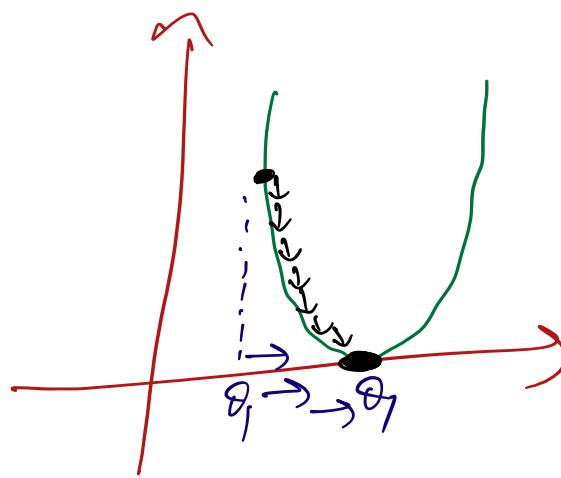
$$\theta_1 : \theta_1 - \alpha (-\text{ve})$$

$$\Rightarrow \theta_1 : \theta_1 + \alpha (\text{ve})$$

So, after few iterations, we come to our global minimum

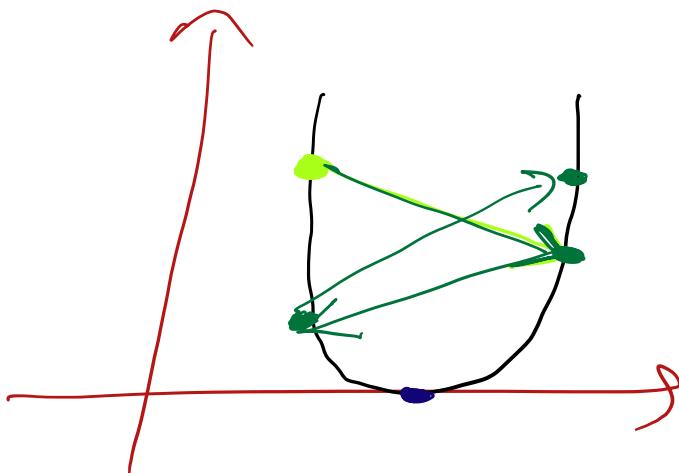
learning rate: By what speed we should come towards our global minimum

ex/ for  $\alpha = 0.01$



it actually takes small steps to reach our global minimum which is so time consuming

for  $\alpha = \text{large}$



This creates a mess and it keeps on jumping and we can't reach global minimum....

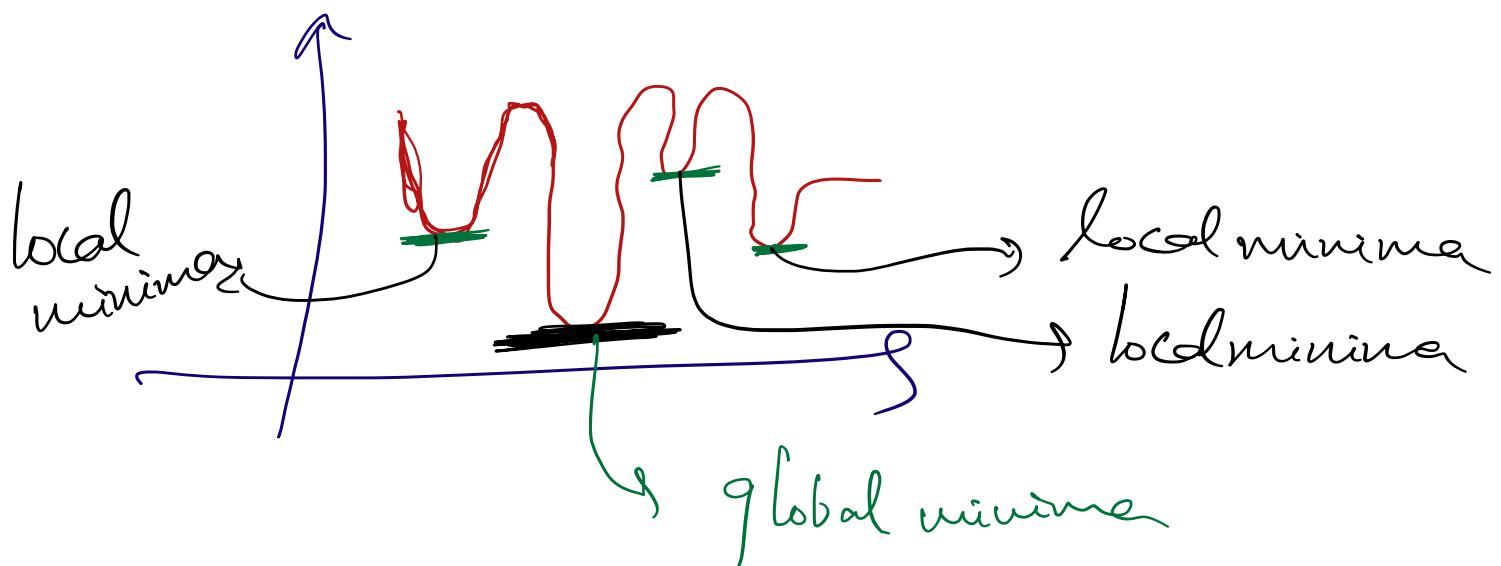
Note! So, acc to the dataset, acc to our data values we should decide what value of  $\alpha$  is best for our model.

global minima is a point where the slope=0

Local minima Vs Global minima:

→ In a Cost function graph when take a big dataset and we have to find more values of  $\theta_1$

- Then we might get stuck at a local minima, but we choose the first  $\theta_i$  value randomly
- There might be a chance for some other  $\theta_i$  value there is a global minima, which we don't know.



Note:

These ML algorithms can't deal with this issue, we need neural networks to solve this, which we will learn in AI.

Deep learning. [GD Algo's  
RMS prop  
Adam optimizer  
etc...]

Gradient Descent Algo:-

} repeat until convergence ~~#~~ suppose a while loop

$$\theta_j := \theta_j - \alpha \left[ \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \right]$$

}

learning rate

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

$$\begin{aligned} \text{if } J=0 &\Rightarrow \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \\ \text{if } J=1 &\Rightarrow \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \end{aligned}$$

↳ called as **Convergence Algorithm**

Repeat until convergence

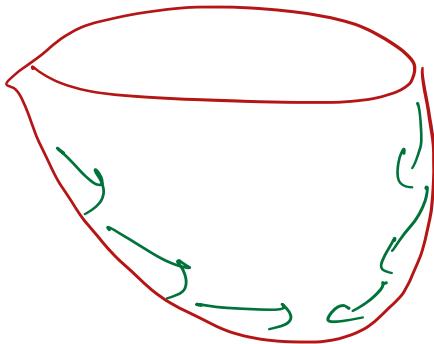
{

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (\text{h}_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$\beta$

Convex:



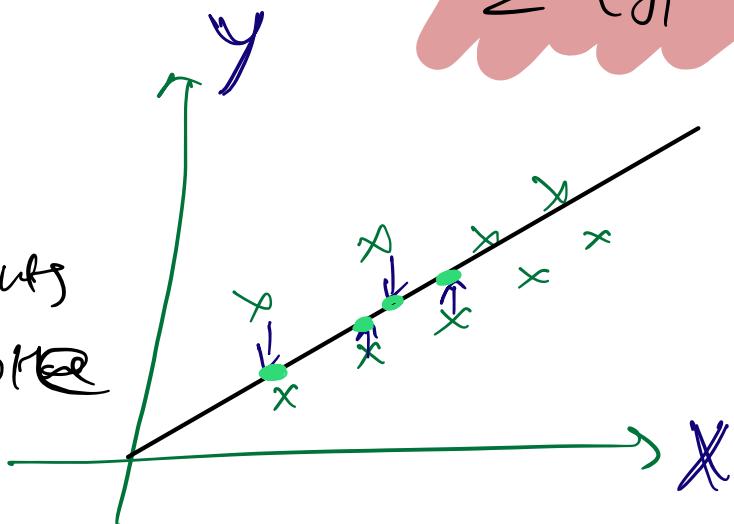
for a 3-d  
GD

performance matrix  $x' \begin{bmatrix} R^2 \\ \text{Adjusted } R^2 \end{bmatrix}$

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

$\hat{y}_i = h_\theta(x)$

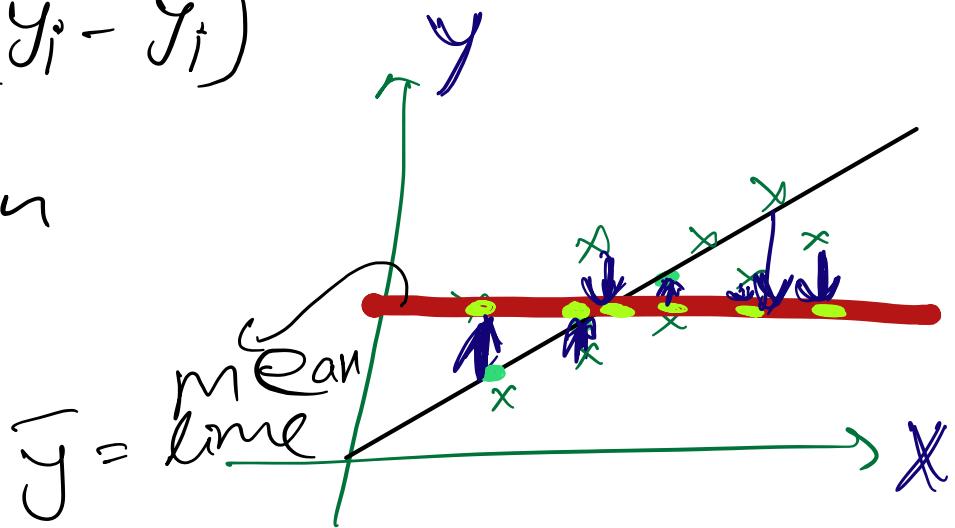
- = predicted points which we denote



$SS_{\text{Res}} = \left[ \text{diff b/w } \text{point } x \text{ and } \text{predicted}(x) \hat{y}_i \right] \uparrow \text{summation of these}$

$$SS_{\text{Total}} = (y_i - \bar{y}_i)^2$$

$\bar{y}_i$  = mean



Note: The diff of dist is high in  $SS_{\text{Total}}$

But, the mean value diff will be higher.

$\Rightarrow 1 - \frac{\text{low}}{\text{high}}$   $\Rightarrow$  small num  
 $\Rightarrow 1 - \text{small num}$   
 $\Rightarrow$  big number.

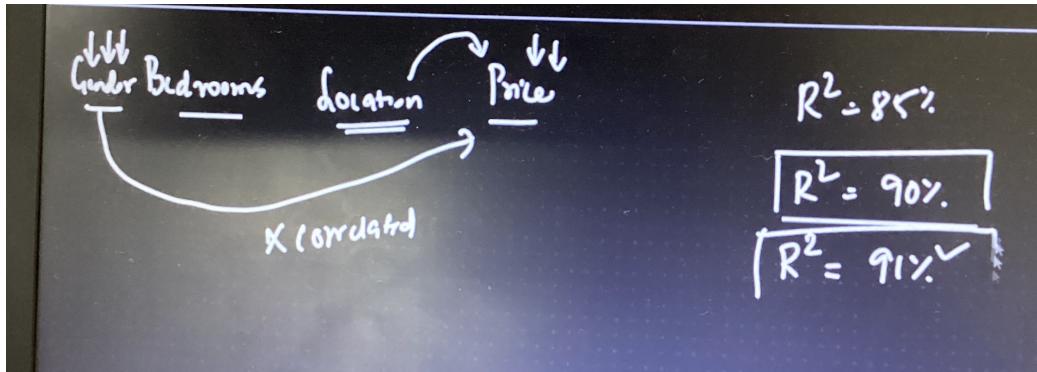
$\Rightarrow R^2$  Value is good (if fitted well)

Note:

If we add values in our dataset, even though they are not at all correlated,

$R^2$  Value will keep on increasing.

exr



we use adjusted  $\overbrace{R^2}$  for this issue.

$\swarrow$

$$\boxed{\frac{1 - (1-R^2)(N-1)}{N-p-1}}$$

$N$  = total no. of samples / data points

$p$  = no. of features (or) predictors

exr  $p=2 \Rightarrow R^2 = 90\% \Rightarrow R_{adj}^2 = 86\%$ .

$p=3 \Rightarrow R^2 = 91\% \Rightarrow R_{adj}^2 = 82\%$ .

$\Rightarrow$  This is because of the formula we took, which balances the performance acc to our  $N$  and  $p$  values