

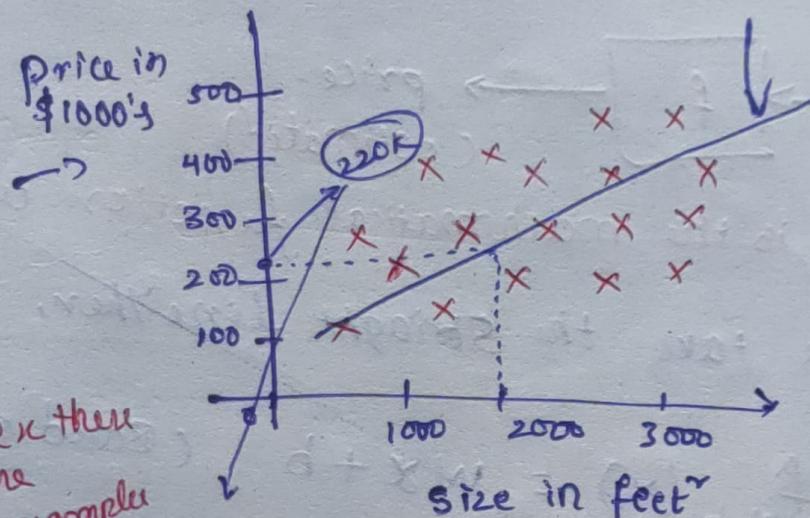
# Linear Reg<sup>n</sup> Model part - I

Best fit line

Data table:

	(x)	(y)
size in feet	price in \$1000's	
1 → 2104	440	
2 → 1416	232	
3 → 1534	315	
852	178	
⋮	⋮	
47 → 3210	870	

For ex there  
are  
47 examples



we will find corresponding priu for  
size of house using  
Regression.

This Dataset Used  
to train the model called "Training set"

$x \rightarrow$  "input" variable (feature)

$y \rightarrow$  "output" variable (target Variable)

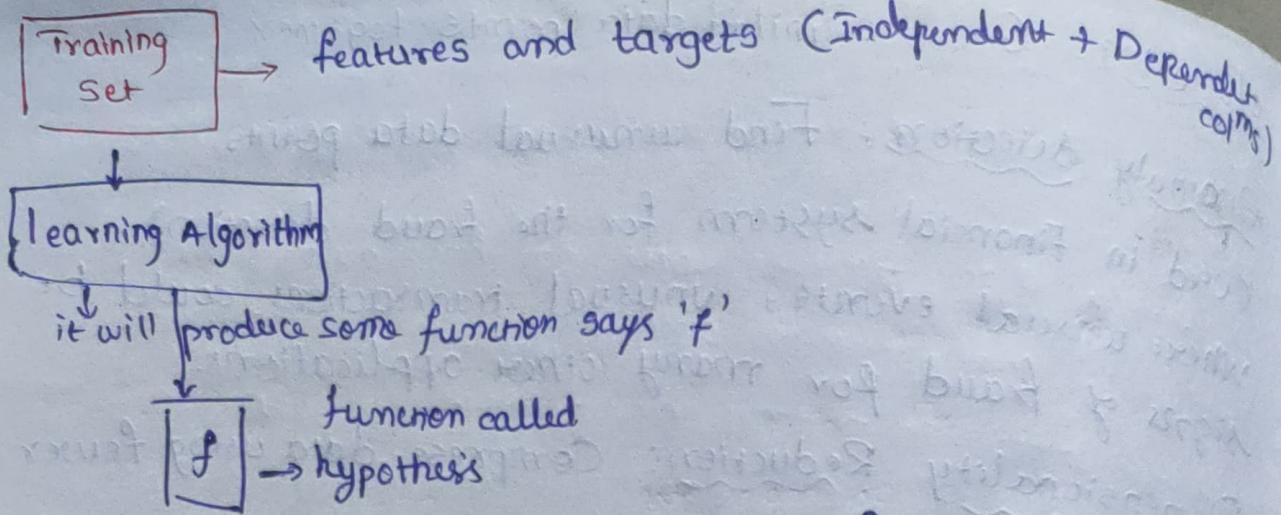
$$x = 2104 \quad y = 440$$

$(x, y) \rightarrow$  single training example like  $(2104, 440)$

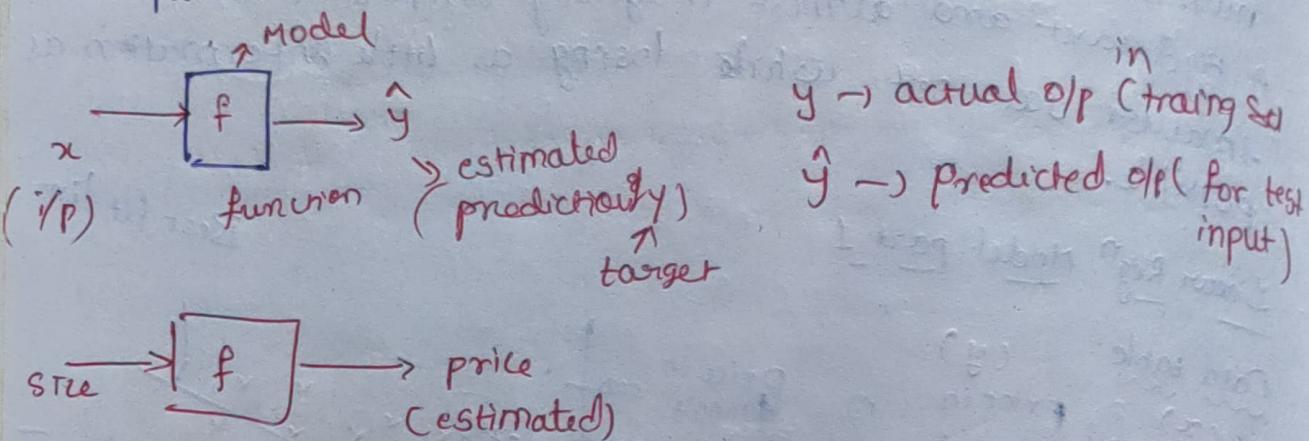
$(x^{(i)}, y^{(i)}) \rightarrow$   $i^{\text{th}}$  training example like 2nd training example

$$(x^{(2)}, y^{(2)}) \rightarrow (1416, 232)$$

$$i = 1, 2, 3, \dots$$



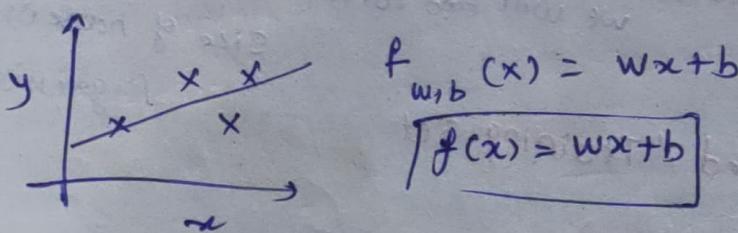
That function will generate new output for the test data



What is the mathematical function we used in the 'f' to compute the output?

Let's take the straight line then,

$$f_{w,b}(x) = w \cdot x + b \quad (\text{Based on } x \text{ we predict } \hat{y})$$

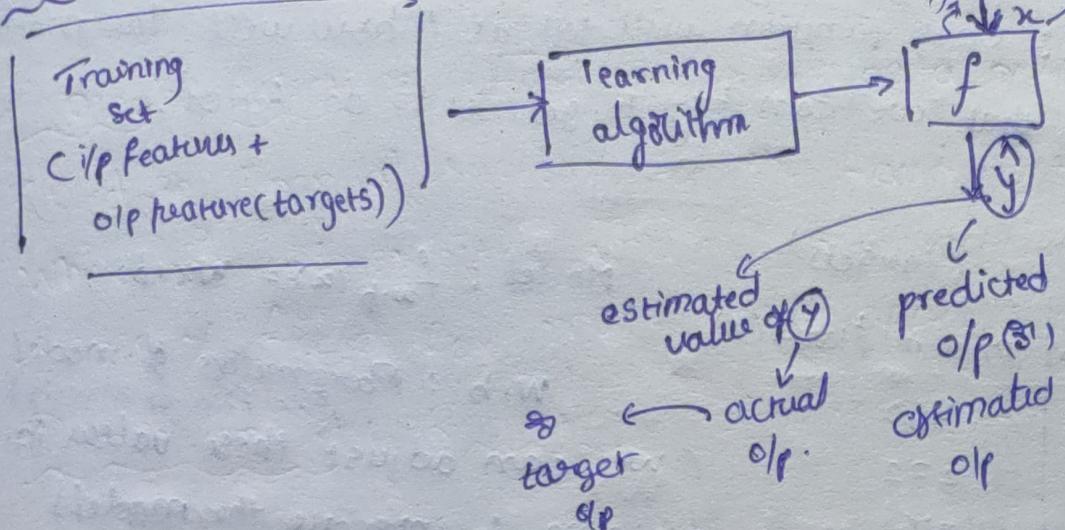


Linear Regn with one variable (single feature x)

(8)  
Univariate Linear Regn

(Later we will how we can find price of house, depending on multiple feature like (size, area, ..., ))

## Linear Reg<sup>n</sup> Part 2:

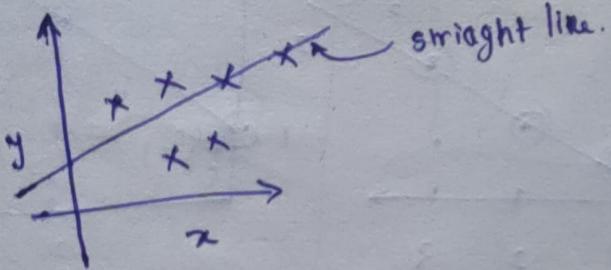


How to Represent  $f$ ? (Model) i/p feature i/p feature

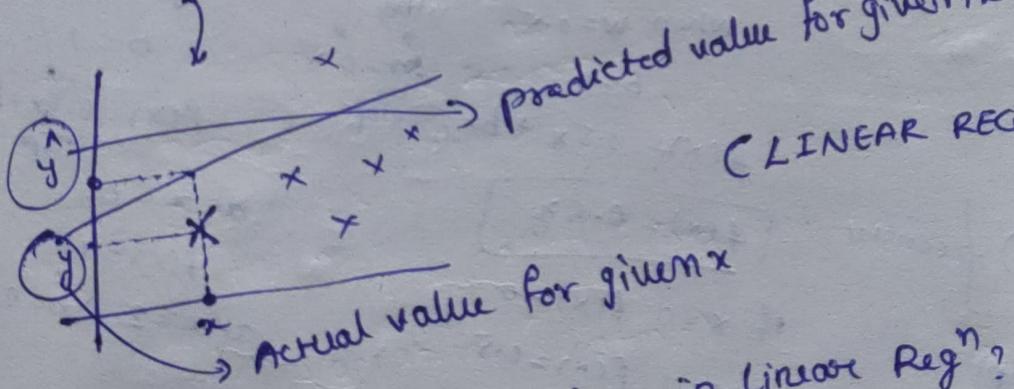
(Mathematical eq<sup>n</sup>)  $\rightarrow f_{w,b}(x) = w \cdot x + b$

(cor)  $\hat{y}$  Based on  $w, b$  value the  $\hat{y}$  values will be determined

$$f(x) = w \cdot x + b$$



we predict value of  $y$  based on straight line



Single feature  $x$   
ie size =  
(LINEAR REG<sup>N</sup> WITH ONE VARIABLE)

Q1 Why you choose straight line in Linear Reg<sup>n</sup>?

It's not like that we need to use some complex  $f^n$  also like curve

(Learn this later)



C polynomial Reg<sup>n</sup> i think

# LINEAR REGN WITH ONE VARIABLE

Cost Function: Cost function will tell us "how well the model is doing so that we can try to get it to do better".

Training Set:

(features)

see in feet<sup>2</sup>) price in \$1000's (y)

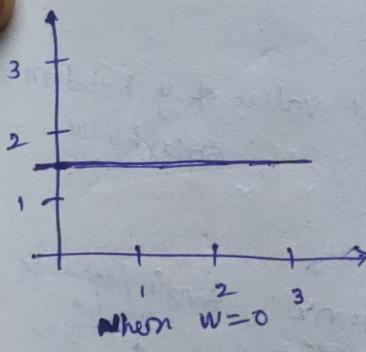
x	y
8104	460
1416	232
1534	315
852	178
⋮	⋮

$$\text{Model: } f_{w,b}(x) = wx + b$$

w, b : Parameters of model  
coefficient & weight

(we can adjust these values in training to improve the model)

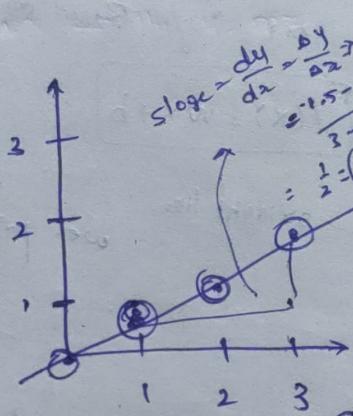
what w, b do? (Based on w, b we have diff f(x) then it generate diff line on the graph)



$$f(x) = wx + b$$

$$\hat{y} = b = 1.5 \quad \text{constant}$$

↓  
y intercept

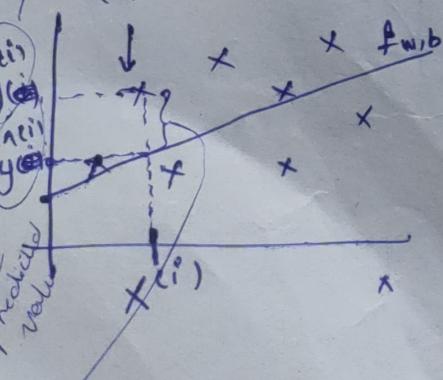


$$f(x) = 0.5x + b < 0 \quad \text{slope}$$

$$\boxed{f(x) = 0.5x}$$

if  $x=0$   $f(x)=0$   
 $x=1$   $f(x)=0.5$   
 $x=2$   $f(x)=1$   
 $x=3$   $f(x)=1.5$

at  $x=0$   
 $f(x)=1$   
at  $x=2$   $f(x)=2$



$$\hat{y}^{(i)} = f_{w,b}(x^{(i)})$$

$$\boxed{f_{w,b}(x^{(i)}) = wx^{(i)} + b}$$

To know how well a line fits the data?

→ This is known by Error b/w  $\hat{y}$  &  $y$  (pred vs actual) for all pts whatever the line gives less error then that will be Best fit line.

cost f^n:  
 $m \rightarrow$  no. of training examples

$$\sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2 \rightarrow \text{Sum of squared errors for all training examples.}$$

Suppose if we have more training examples then  $m$  i.e. so that the sum of square of errors are also increased. So for that for large training set examples we take average like cost function as

$$\left| \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2 \right|$$

some ML Engr use  $\frac{1}{2m}$  instead of  $\frac{1}{m}$  to make neat calculation.

We can use  $\frac{1}{2m}$  or  $\frac{1}{m}$  both are okay.

So, Finally

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

$m = \text{no. of training examples.}$

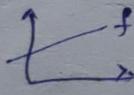
→ Squared error cost function.

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w, b}(x^{(i)}) - y^{(i)})^2$$

Cost Function Intuition:

If we are looking to fit a straight line to training data, then we have below:

model:  $f_{w, b}(x) = w x + b$



parameters:  $w, b$

cost  $f^n$ :  $J(w, b) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$  (To know how best our line fits the data)

goal: minimize  $J(w, b)$  (Whatever line gives less  $J(w, b)$  value that will be the best fit line)

If you want to better visualize the cost function  $J$ , this work of a simplified version of Linear Reg "model". We are going to use  $f_w(x) = wx$

$$f_w(x) = wx$$

$\leftarrow$  simplified model.

Now we have only one parameter 'w' & cost  $J^n$

$$J(w) = \frac{1}{2m} \sum_{i=1}^m (f_w(x^{(i)}) - y^{(i)})^2$$

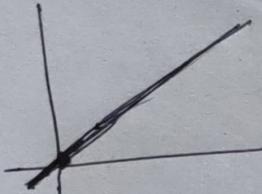
with this simplified model our goal is to find minimize value of  $J(w)$

$$\underset{w}{\text{minimize}} J(w)$$

{ we need to find value of  $w$  that minimize the value of  $J(w)$  }

if  $b=0$ , then  $f_w(x) = wx \rightarrow$   $y = mx$

$\downarrow$   
straight line passing from origin



Let's see how cost function changes for different values of 'w'

~~$f_w(x)$~~

$f_w(x) = wx$

input

for fixed  $w$ , function of  $x$

Like the  $f_w(x)$  &  $\hat{y}$  is dependent on 'x'  
ends on 'x'  
input

(function of input  $x$ )

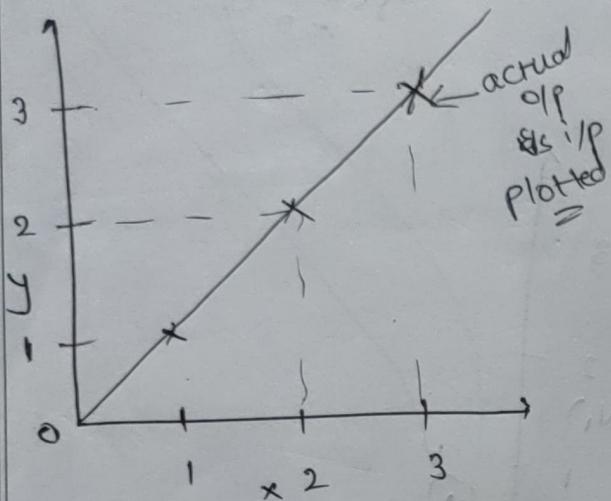
$$J(w)$$

$$J(w) = \frac{1}{2m} \sum_{i=1}^m (y^{(i)} - \underbrace{wx^{(i)}}_{f_w(x)})^2$$

(function of  $w$ )

where 'w' controls the slope of line defined by  $f_w(x)$ . so cost defined on  $J$  depends on parameter 'w'.

(function of parameter  $w$ )

$f_w(x)$ (for fixed  $w$ , function of  $x$ )  
↓  
inputLet's say  $w=1$ 

$$f_w(x) = w x^w = x$$

$$\underbrace{f_w(1)}_g = 1 \quad \underbrace{f_w(2)}_y = 2$$

Let's calculate  $J$  when  $w=1$ 

$$J(1) = \frac{1}{2(3)} \sum_{i=1}^3 (\hat{y}^{(i)} - y^{(i)})^2$$

$\downarrow^{(i)}$   
 $f_w(x)$   
 $\downarrow^{w x^{(i)}}$

$$= \frac{1}{2m} \sum_{i=1}^m (wx^{(i)} - \hat{y}^{(i)})^2$$

here the error is

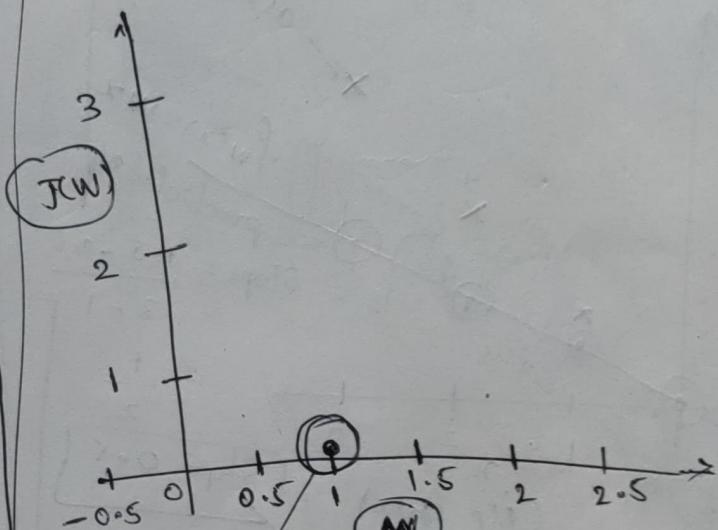
$$= \frac{1}{2m} (0 + 0 + 0) = 0$$

$\downarrow^{(1-1)(1-1)(1-1)}$

$$J(1) = 0$$

 $J(w)$ (function of  $w$ )  
↑ parametercost  $f^n$ 

J

As seen here at  $w=1$ , the  $J=0$ 

$$\text{at } w=1 \quad J(w) = 0$$

$$\boxed{J(1) = 0}$$

$f_w(x)$

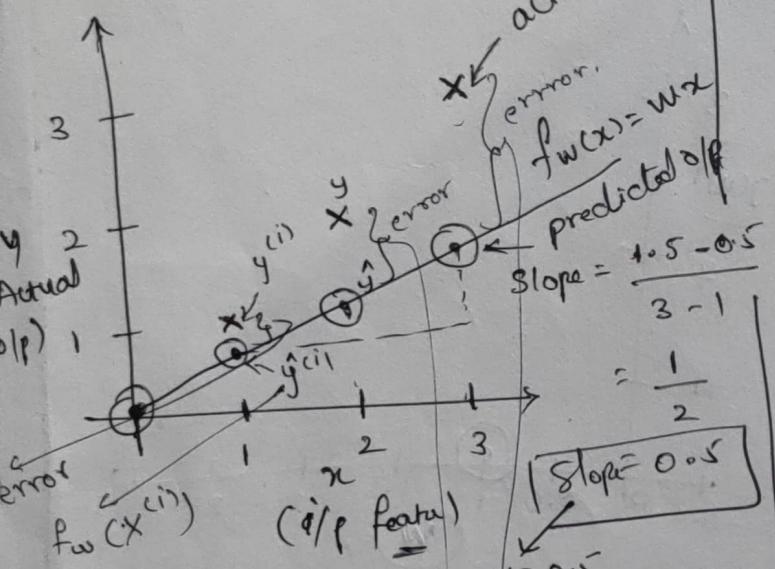
Now take  $w=0.5$

$$f_w(x) = 0.5x$$

$$\downarrow \text{slope} = 0.5$$

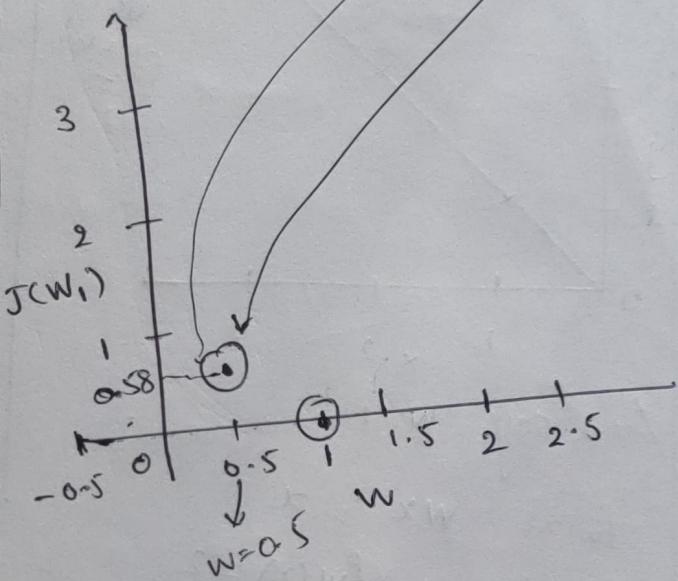
$$\begin{array}{ll} \text{at } x=1 & f_w(x) = 0.5 \\ x=0 & f_w(x) = 0 \\ x=2 & f_w(x) = 1 \end{array}$$

$$x=3 \quad f_w(x) = 1.5 \\ \text{P/L}$$



$J(w)$

Now plot  $J$  value,  
 $J(0.5) \text{ when } w=0.5 = 0.58$



$J(0.5)$

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

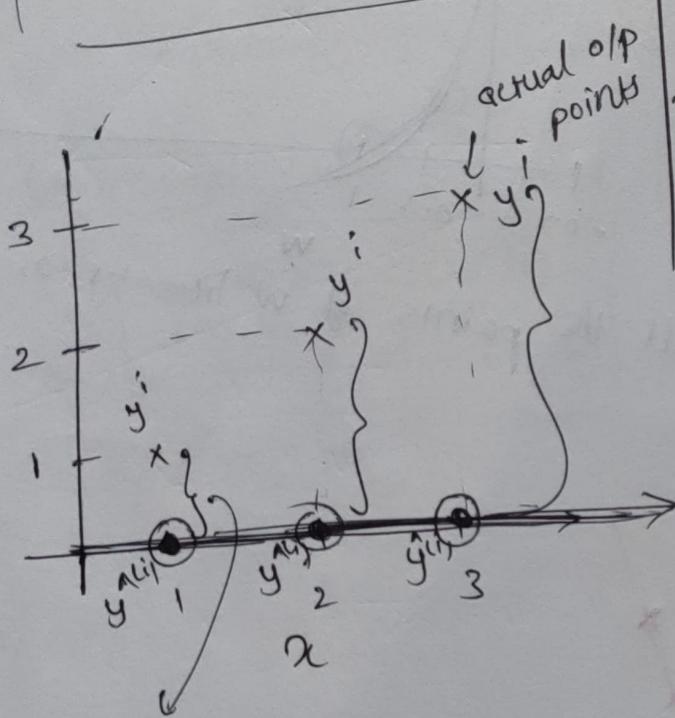
$$= \frac{1}{2 \times 3} ((0.5 - 1)^2 + (1 - 2)^2 + (1.5 - 3)^2)$$

$$= \frac{1}{2 \times 3} [3.5] = \frac{3.5}{6} = 0.58$$

$f_w(x)$ Take  $w=0$ 

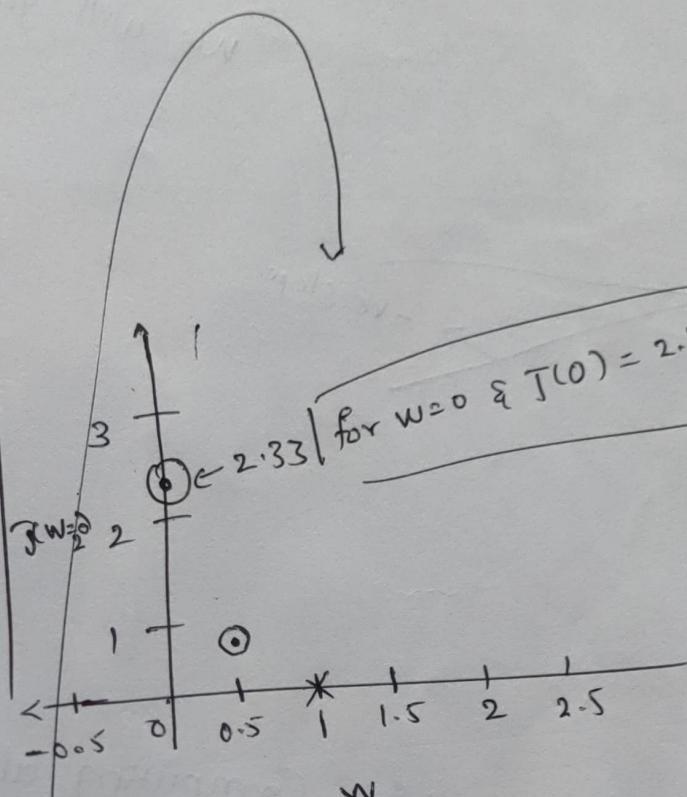
$$f_w(x) = w \cdot x \\ = 0$$

for any value of  $x$ , the  $f_w(x)$  is zero



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

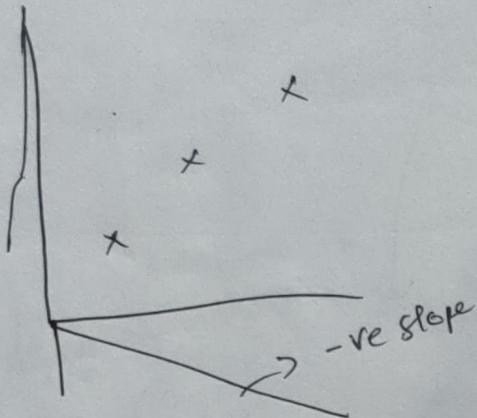
$$= \frac{1}{6} (14) = 2.33 \dots$$

 $J(w)$ 

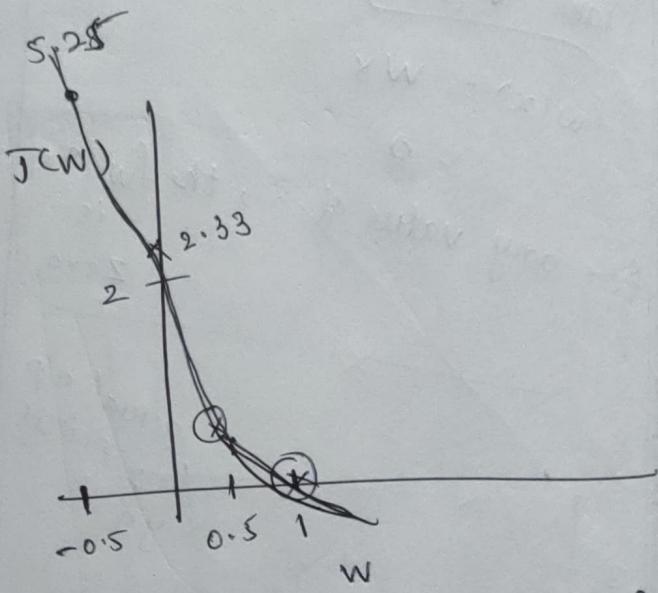
$e^{-2.33}$  for  $w=0 \Rightarrow J(0) = 2.33$

If we repeat the process for multiple points like

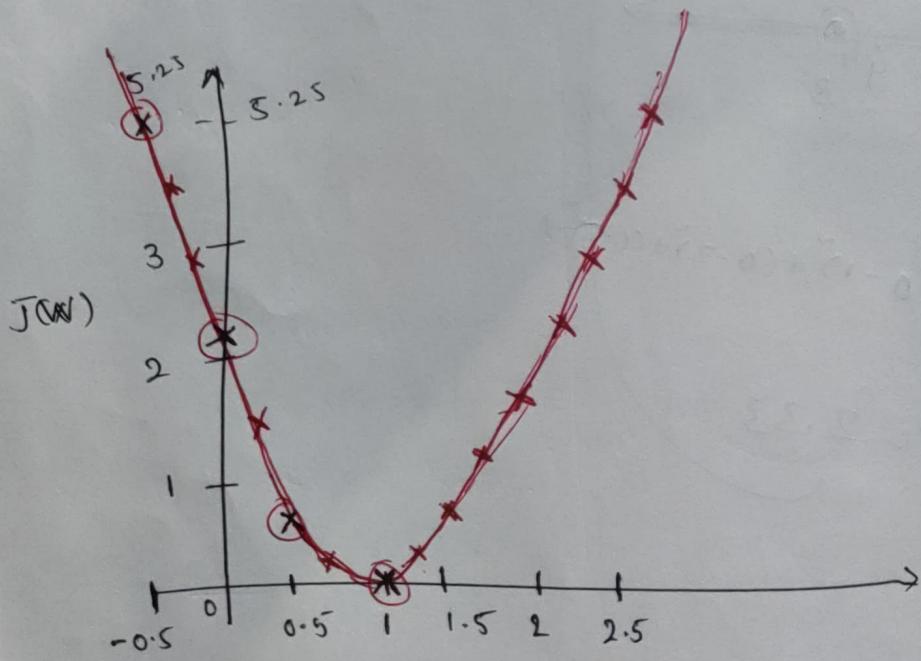
$$w = -0.5, \text{ then } f_w(w) = -0.5x$$



We will get  $J(w) = 5.25$  something  
then



And Finally After Computing all the points of  $w$  like  $1.5, 2, 2.5 \dots$   
then we will get  $J(w)$  like below



(\*) → we plotted prev.

So Here the question is how can you choose value of  $w$  to Minimize  $J(w)$ ?

We need value of  $w$ , where we have less  $J(w)$ .

In this example, From graph at  $w=1$ , we have  $J(w) = 0$

so,  $w=1$  is the Best value for the Best fit Line.

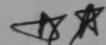


Goal of Linear Regression

minimize  $J(w)$

$w$

General Case:



minimize  $J(w, b)$

$w, b$

## Visualizing the Cost Function:

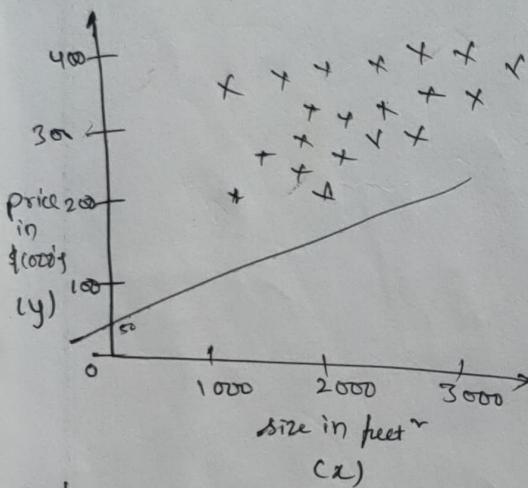
Till now, we seen  $f_w(x) = wx$  i.e at  $b=0$

Now we will see when  $b \neq 0$

$$\text{i.e } f_{w,b}(x) = wx + b$$

$$f_{w,b}$$

(function of  $x$ )



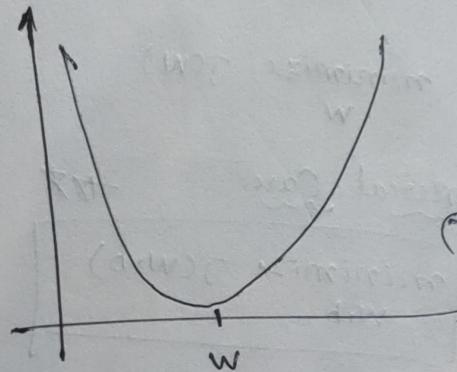
Let's take like

$$f_{w,b}(x) = 0.06x + 50$$

$w$   
Slope       $b$   
                intercept

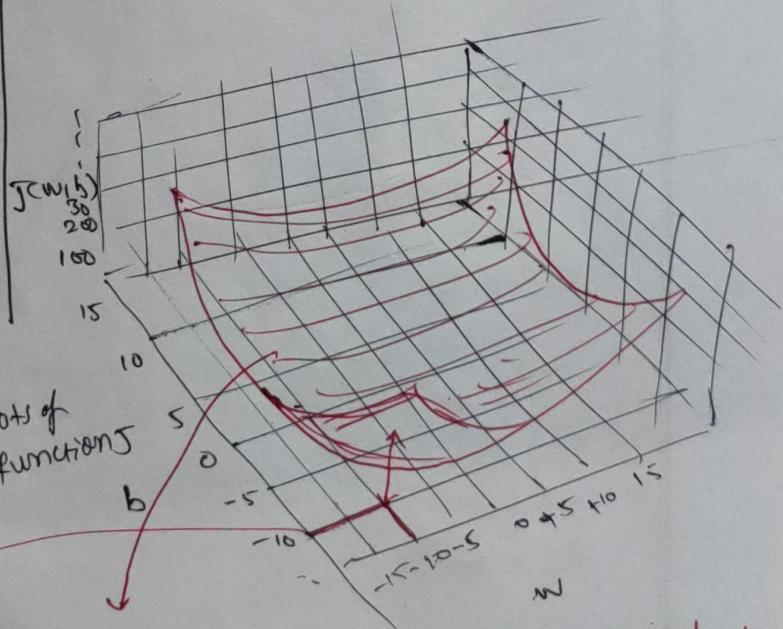
We know this Model is worst because it Underestimate the price of house  
let's consider as ex.

(function of  $w, b$ )  
Both parameters



When  $b=0$ ,  $J$  is function of  $w$  only then we got a function like above

But now we have two parameters ( $w, b$ )



Contour

Contour plots are best to visualize 3D plots of cost functions

at  $b = -10$  &  $w = 10$

the  $J$  is height of arrow

Now also it is like U-shaped but in 3D shape  
As we vary  $w$  &  $b$  we get diff values for  $J(w,b)$

\* What we really need is a efficient algorithm that you can write in code for automatically finding the value of parameters  $w$  &  $b$  they give you the Best fit Line that minimizes cost function. There is an algorithm for doing this called "GRADIENT DESCENT".

TRAIN THE MODEL WITH GRADIENT DESCENT:

GRADIENT DESCENT

Gradient descent not only used in Linear Reg<sup>n</sup> But also used in some deep learning models.

Have some function  $J(w, b)$

Want to minimize  $J(w, b)$

Gradient descent is used to minimize any  $f^n$  (not only Linear Reg<sup>n</sup>)  
 " " is applied to more general functions, including other cost functions that work with models that have more than two parameters.

for ex for instance we have cost  $f^n$   $J(w_1, w_2, w_3, \dots, w_n, b)$

Our objective is to minimize

$\underset{w_1, w_2, w_3, \dots, w_n, b}{\text{minimize}} J(w_1, w_2, w_3, \dots, w_n, b)$

Outline :

$\boxed{\underset{w, b}{\text{Minimize}} J(w, b)}$

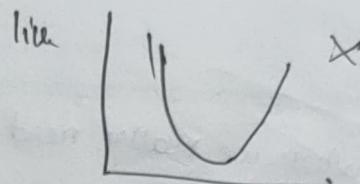
1) Start with some  $w, b$  (like start with  $w=0, b=0$ ) as initial guess

2) Keep changing  $w, b$  to reduce  $J(w, b)$

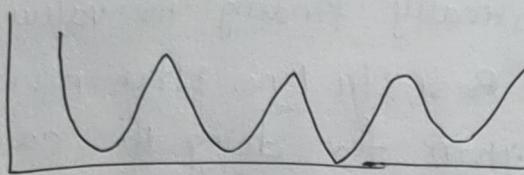
Until we settle at or near a minimum

One thing we need to note, the J is not always U-shaped

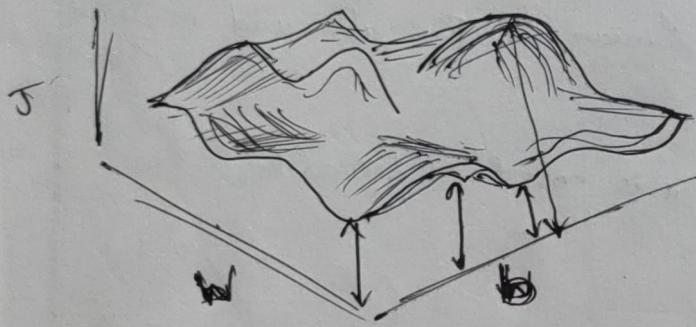
Sometime we may have  $>1$  minimum values



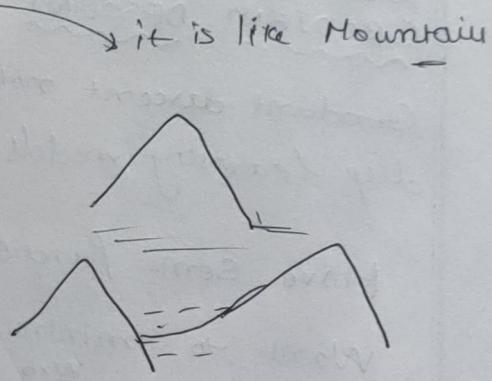
for ex:



Multiple low value

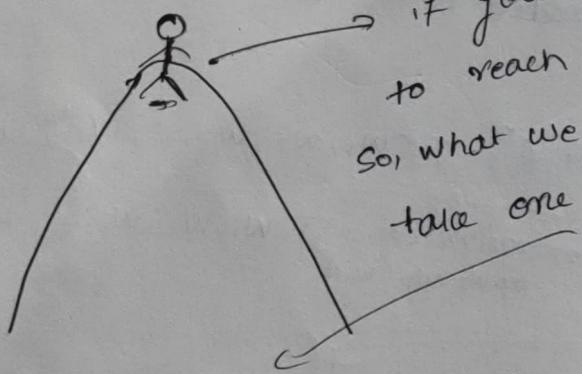


3D plot

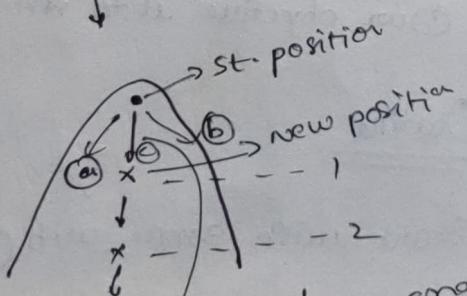


How Gradient descent Algo works?

if you are on the Mountain & you want to reach bottom of mountain efficiently as possible  
So, what we will do, we will look around & take one tiny step towards bottom.



we have 3<sup>rd</sup> possible to take tiny Step But we choose Middle one because Middle one is efficient we can reach faster if we select middle



reached end (c)  
This will take downhill faster than a tiny little baby step you could have taken in other direction (a, b)

After taking one step, we need to repeat process

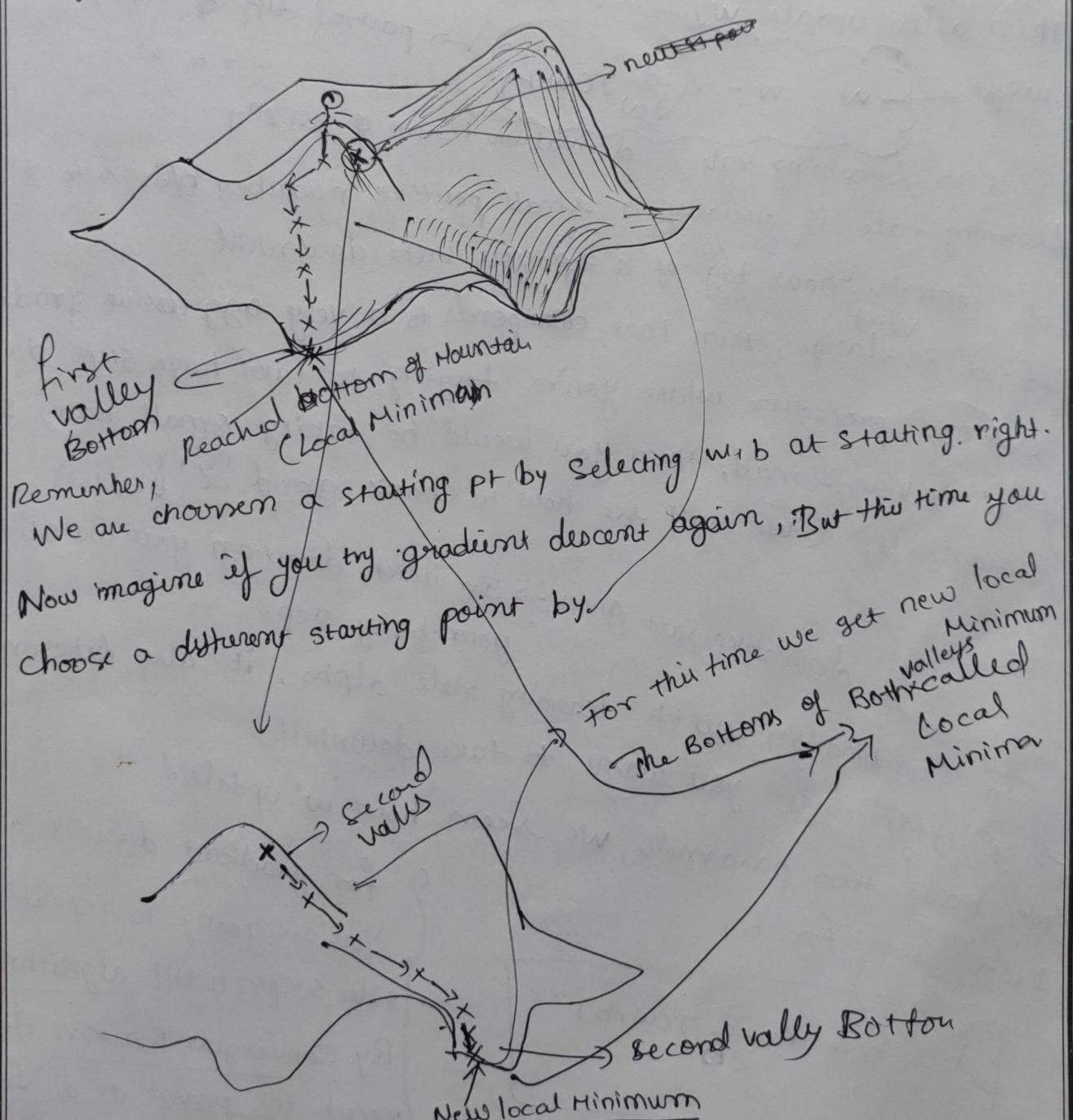
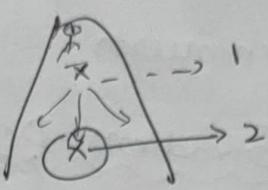
→ Need to see  $360^\circ$

→ choose the Best path

(What next little baby step in order  
to move downhill)

If you do that and take another step.... until end of Mountain

$\curvearrowleft 360^\circ$



Because if you start going down the first valley, Gradient descent won't lead you to second valley, & the vice versa is also true. & same is true if you started going down the second valley, you stay in that second minimum and not find your way into first local minimum.

### Implementing Gradient Descent Algorithm:

1) First take ~~one value~~ of we have taking the value for w right. Then after update w.

~~Weight~~  $w = w - \alpha \frac{\partial J(w, b)}{\partial w}$  partial diff of cost  $J^n$  w.r.t. w.

$\alpha \rightarrow$  Learning rate derivative term of cost  $J^n$

Learning rate is usually a small positive number b/w 0 to 1 (say 0.01)

$\alpha'$  controls how big of a step you take downhill

If  $\alpha'$  is very large, then that corresponds to a very aggressive gradient descent procedure where you're trying to take huge steps downhill.

If  $\alpha'$  is very small, then you would be taking small baby steps downhill. (We will see how to choose good ' $\alpha'$ ' later)

In combination with learning rate alpha, it also determines the size of the steps you want to take downhill.

We have two parameters, we seen how 'w' updated let's see for b,

$$b = b - \alpha \cdot \frac{\partial J(w, b)}{\partial b}$$

$$w = w - \alpha \frac{\partial J(w, b)}{\partial w}$$

} for gradient descent algorithm, you are going to repeat the side steps until algorithm converges. By Converges, I mean that you reach the point at a local minimum where the parameters w & b no longer change much with each additional step that you take.

Now, there's one more subtle detail about how to correctly implement semantic gradient descent. You are going to update two parameters  $w$  and  $b$ . This update takes place for both parameters,  $w$  and  $b$ . One important detail is that for gradient descent, you need to simultaneously update  $w$  and  $b$  (need to update both parameters at the same time) (old  $w$  to New  $w$  & old  $b$  to New  $b$ )

code for update

$$\text{tmp\_}w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

$$\text{tmp\_}b = b - \alpha \frac{\partial}{\partial b} J(w, b)$$

$$w = \text{tmp\_}w$$

$$b = \text{tmp\_}b$$

Correct:

simultaneous update ✓

(✓ do this one) ✓

Incorrect

$$\text{tmp\_}w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

$$w = \text{tmp\_}w$$

$$\text{tmp\_}b = b - \alpha \frac{\partial}{\partial b} J(w, b)$$

$$b = \text{tmp\_}b$$

here  $w$  changing

so, this is incorrect notation.

here  $\rightarrow$   $\text{tmp\_}b$  is diff & here  $\text{tmp\_}b$  value is also diff

NOTE:

Gradient descent is an algo for finding values of parameters  $w$  and  $b$  that minimize the cost function  $J$ . What does this update statement do? (Assume  $\alpha$  is small)

$$w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

Ans) Update parameter  $w$  by small amount.

GRADIENT DESCENT INTUITION: (Let's see how w & derivative & Multiplication of both is ~~actually~~ actually worse)

Repeat until convergence {

→ Learning rate (It control the how big of a step you take) when updating w & b

$$w = w - \alpha \frac{\partial J(w, b)}{\partial w}$$

Derivative part → It will

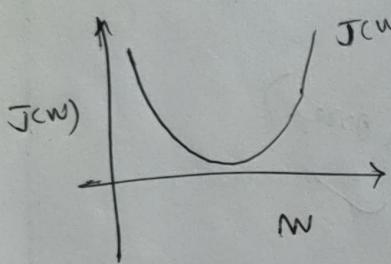
$$b = b - \alpha \frac{\partial J(w, b)}{\partial b}$$

partial

Let's see small example how updation is done & all

Let's take  $J(w)$  instead of  $J(w, b)$  →  $w = w - \alpha \frac{\partial J(w)}{\partial w}$

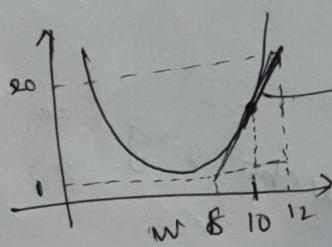
Our aim is to minimize  $J(w)$



→ it is like our prev example where we had temporarily set  $b=0$

Now  $\frac{\partial}{\partial w} J(w)$  means slope of tangent at ~~one particular~~ value of w.

Let's take assume  $w = 10$



$$\text{slope} = \frac{20-1}{12-8} = \frac{19}{4} = 4.75 \rightarrow$$

↑  
+ve  
> 0

slope =  $\frac{\partial}{\partial w} J(w)$

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

$$\text{So, } w = w - \alpha (+ve \text{ number})$$

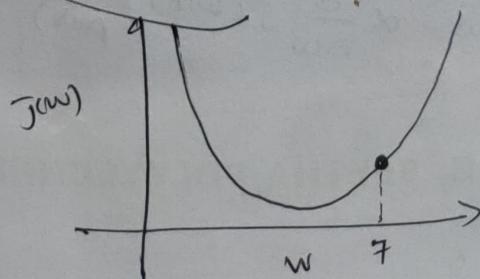
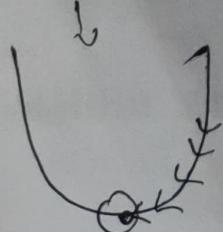
$\alpha$  is always +ve so,

$$\text{finally, } w = w - (+ve \text{ number})$$

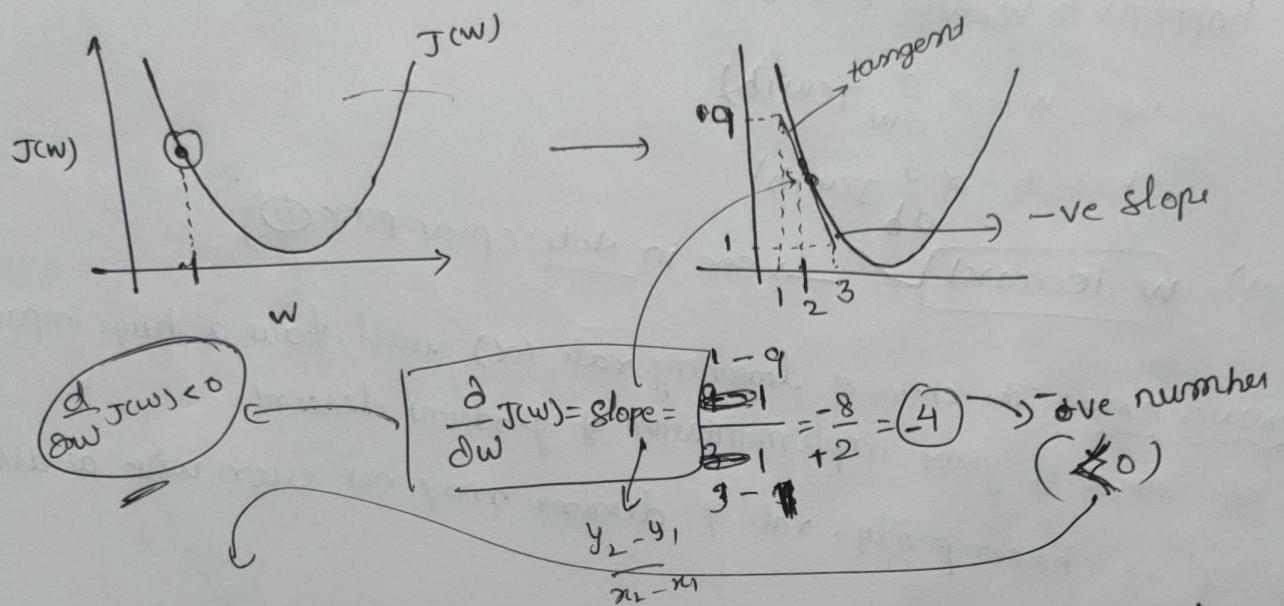
final w value is lower.

i.e. in graph we are moving to left (say  $w=7$ )

So after moving to left again & again it will reach to minimum point.



Let's take another point of  $w$  (this time at right)

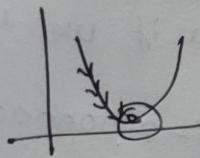


$$w = w - \alpha \text{ (-ve number)}$$

Then  $w$  value increases then it move towards right side  
 (so to know the Minimum point)

$$w = w + \overset{\text{+ve}}{\curvearrowleft}$$

w increases



NOTE :

If  $\frac{d}{dw} J(w) > 0 \rightarrow$  i.e slope is (+ve) then  $w$  value will decreases

$\frac{d}{dw} J(w) < 0 \rightarrow$  i.e slope is (-ve) then  $w$  value increases

Ques

- i) Assume ' $\alpha$ ' is very small +ve number. When  $\frac{d}{dw} J(w, b)$  is a +ve number ( $> 0$ ) → as in example in the upper part of slide down above, what happens to  $w$  after one update step?

$$w = w - \alpha \frac{d}{dw} J(w, b)$$

$$b = b - \alpha \frac{d}{db} J(w, b)$$

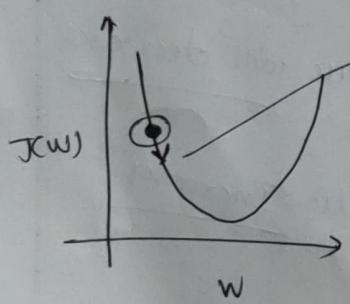
Ans) w decreases → see in Note: (prev page 😊)

LEARNING RATE: The choice of learning rate ( $\alpha$ ) will have a huge impact on the efficiency of your implementation of gradient descent  
If ' $\alpha$ ' is chosen poorly, rate of descent may not even work at all  
How to choose ' $\alpha$ '?

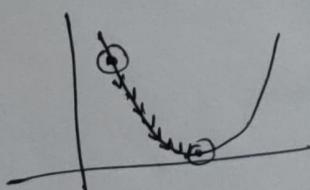
Let's take,  $w = w - \alpha \frac{d}{dw} J(w)$

Let's see what happens if we take either too small (or) too high.

If  $\alpha$  is too small: ( $\alpha=0.0000001$ )



if we take  $\alpha$  as small then the  $w$  is increased by very very small value that means we take a little baby step.  
So from that step we will repeat process, we will take so many tiny steps to reach minimum point.



Here we ~~are~~ The process is done like  $J$  is decreasing slowing upto its minimum value by taking tiny steps)

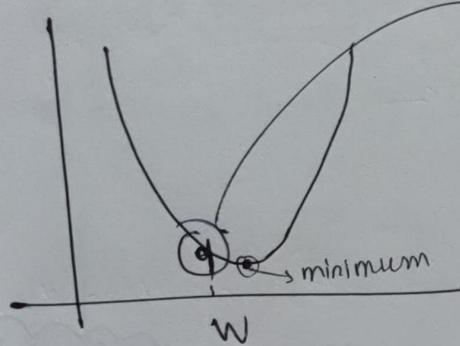
But we ~~need~~ notice that you're going to need a lot of steps to get minimum

To Summarize,

If ' $\alpha$ ' is too small, Gradient descent will be boric, but it will be slow

Because it take so many steps  
(it takes more time)

If  $\alpha$  is too large:



Let's say we assume 'w' here.  
(It is pretty close to minimum)

But if  $\alpha$  is high then,

$$w = w - \alpha (-\text{ve})$$

↓  
high

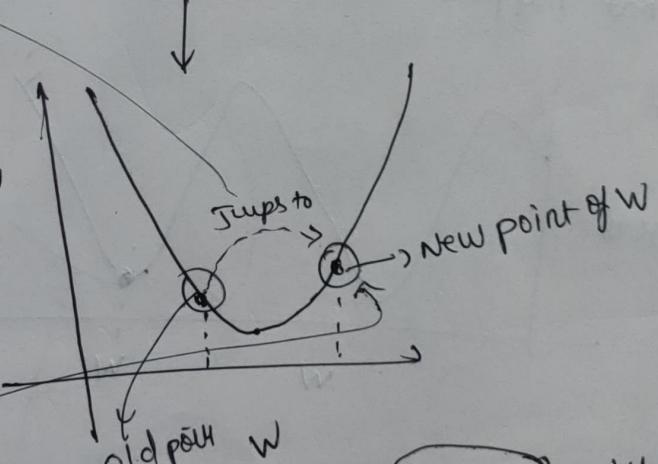
-ve slope no  
for left side  
pt

$$w = w + \alpha$$

↑  
high

w will increase (a giant step will taken)

It jumps to giant step because  
of  $\alpha$  is too high. So,  
we missed the minimum point &  
J value keep on increase.  $J(w)$   
cost will worse



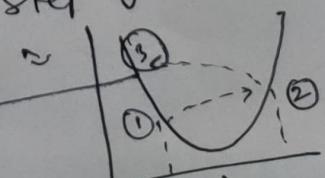
Then after reaching the derivative of  $w = w - \alpha \frac{d}{dw} J(w)$  will  
(slope)

be (+ve) then again  $w = w - \alpha (+\text{ve})$

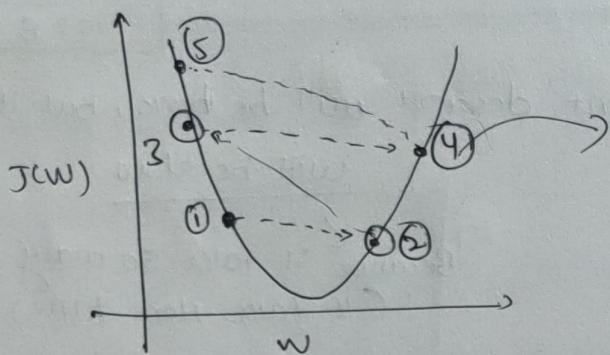
↑  
high

$w =$  decrease by giant step again

Then new point will be here



If  $\alpha$  is too high (too large) the steps will like



here we get +ve slope  
&  $\alpha$  is also large  
then,  
 $w = w - \alpha (\text{+ve})$

~~w is very less~~  
w reduces very high  
(like w takes  
giant step  
to ⑤)

So, If  $\alpha$  is too large,

Gradient descent may

→ Overshoot, never reach minimum

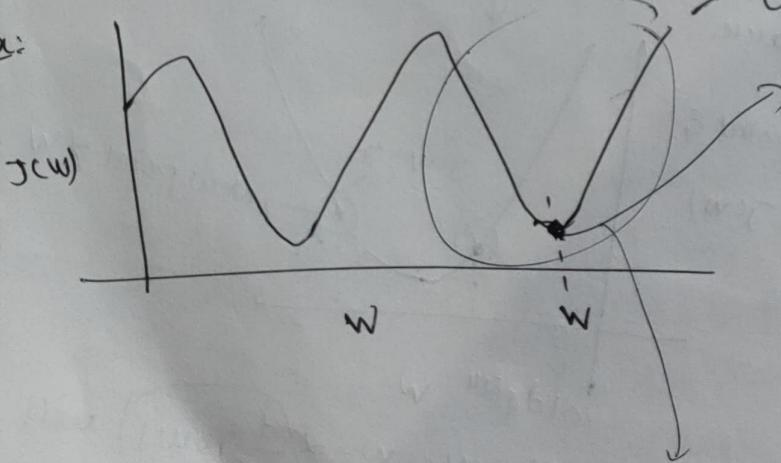
→ Fail to converge, diverge

IMPORTANT

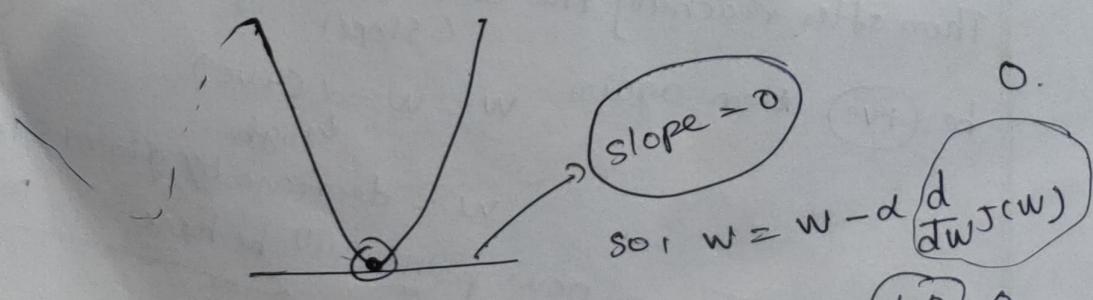
(Used for interview)

What happens Once we reach minimum point.  
squared error cost function

Ex:



let's say after doing all steps  
we reached to Minimum  
so, after that what happens?



slope = 0

$$\text{so, } w = w - \alpha \frac{d}{dw} J(w)$$

$$w = w - \alpha \cdot 0$$

$$w = w$$

So, w won't change

If we are at local minima, then we get slope = 0 then

$$w = w - \alpha \frac{d}{dw} J(w) \rightarrow 0$$

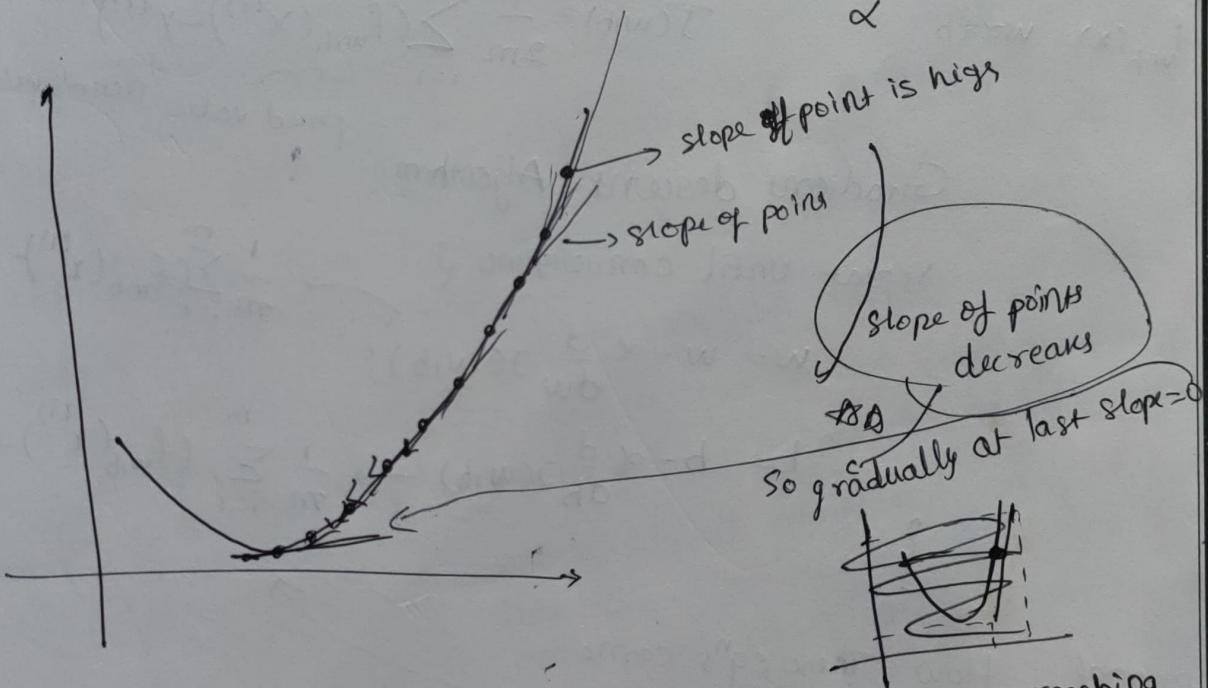
$$\Rightarrow w - \alpha \cdot 0 = w$$

The Gradient descent leaves 'w' unchanged when we are at minimum point. because at that point  $\boxed{\text{slope} = 0}$  (horizontal line)

$$w = w$$

The w won't update.

Q) Can reach local minimum with fixed learning rate?



$$w = w - \alpha \frac{d}{dw} J(w)$$

$\frac{d}{dw} J(w)$  getting Reduces when we reaching towards local minimum

Near local minimum,

→ Derivative becomes smaller.

→ Update steps become smaller.

Update steps small because ( $\alpha$  is fixed)

$\alpha$ -fixed so,  $w$  depends on slope only

If slope  $\downarrow$  then  $w$  also  $\downarrow$  with small step.

So we can reach local minimum with fixed  $\alpha$ . ✓

## GRADIENT DESCENT FOR LINEAR REGRESSION:-

Here we use  $J$  as exact linear Reg<sup>n</sup> cost function (mean squared error cost function) putting together Gradient descent with cost function will give you first learning algorithm, the Linear reg<sup>n</sup> algorithm.

Linear Reg<sup>n</sup> model

cost function

$$f_{w,b}(x) = wx + b$$

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

pred value      act<sup>n</sup> val

Gradient descent Algorithm

repeat until convergence {

$$w = w - \alpha \frac{d}{dw} J(w, b)$$

$$\frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$$

$$b = b - \alpha \frac{d}{db} J(w, b) \rightarrow \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$$

Optional : How These eq<sup>n</sup>'s came

$$\frac{\partial}{\partial w} J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 \rightarrow \text{partial derivative w.r.t derivative}$$

$$\begin{aligned} &= \frac{\partial}{\partial w} \frac{1}{2m} \sum_{i=1}^m (wx^{(i)} + b - y^{(i)})^2 \rightarrow \frac{\partial}{\partial w} (wx^{(i)} + b - y^{(i)}) \\ &= \frac{1}{2m} \sum_{i=1}^m \cancel{2} \cdot (wx^{(i)} + b - y^{(i)}) \cancel{\frac{\partial}{\partial w} (wx^{(i)} + b - y^{(i)})} \rightarrow 2(x+2+3) \cdot \frac{\partial}{\partial x} (x+2+3) \\ &= \frac{1}{m} \sum_{i=1}^m (wx^{(i)} + b - y^{(i)}) (x^{(i)} + 0 - 0) \rightarrow 2(x+2+3)(1+0+0) \\ &\quad + \frac{1}{m} \sum_{i=1}^m f_{w,b}(x^{(i)}) - y^{(i)}) \end{aligned}$$

Hy,  $\frac{d}{db} \text{scw}(b)$  also  $\frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$

↳ here we done partial diff w.r.t.  $b'$   
except  $b$  all will consider as constant.

Gradient descent algo for Linear Reg<sup>n</sup> is

repeat until convergence {

$$w = w - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

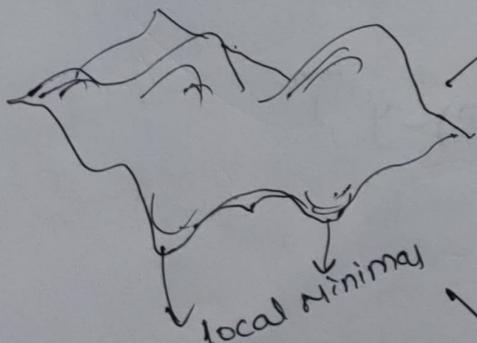
$$b = b - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$$

}

$$f_{w,b}(x^{(i)}) = w x^{(i)} + b *$$

Need to update  
 $w, b$  simulta  
neously

Let's familiar with how gradient descent works



local minima  
local minima  
(2 local minima)

One ~~the~~ stroke we saw with gradient descent is that it can lead to local minimum instead of global minimum where global minimum means the point that has the lowest possible value for cost function of all possible pts

depending on where we initialize  $w, b$ , we can end up with different local minima.  
(We seen this earlier also)