

## 5. Gradient Descent.

Have some function  $J(w, b)$   
Want  $\min_{w, b} J(w, b)$ .

We use Gradient Descent to minimize any cost function.

Our objective is to minimize.

$$\min_{w_1, w_2, \dots, w_n, b} J(w_1, w_2, \dots, w_n, b).$$

Outline:

Start with some  $w, b$ . (set  $w=0, b=0$ )  
keep changing  $w, b$  to reduce  $J(w, b)$ .  
until we settle at or near a minimum.

$J$  may have more than 1 possible minimum.

An example can be given below, we see an example in 3D.

We can see that it has more than 1 minimum values.

We want to take a decision of steepest descent.  
and we take this step everytime. Slowly, we take these steps, to find ourselves at the bottom of  $J$  at a local minimum.

Depending on the initial position, we might end up on a different local minima.

## II. Implementing Gradient Descent.

$$w = w - \alpha \frac{d}{dw} J(w, b) \quad \left. \right\} \text{assign value of } w.$$

$\alpha$  is called a Learning Rate  
 $0 < \alpha < 1 \rightarrow$  how big of a step we take.

$\frac{d}{dw} J(w, b) \rightarrow$  Derivative term of cost function  $J$ .

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

We repeat these steps until convergence.

We need to simultaneously update both  $w$  and  $b$ .

Correct way:

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

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

$$w = \text{temp\_} w ; b = \text{temp\_} b.$$

Incorrect way:

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

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

$$\text{temp\_} b = b - \alpha \frac{d}{db} J(w, b) \quad b = \text{temp\_} b.$$

We are updating  $w$  first and then using it for calculating  $b$ . This is wrong.

## II. Gradient Descent Intuition.

Algorithm: Repeat until convergence:

$$w = w - \alpha \frac{d}{dw} J(w, b) \quad \alpha \rightarrow \text{Learning rate.}$$

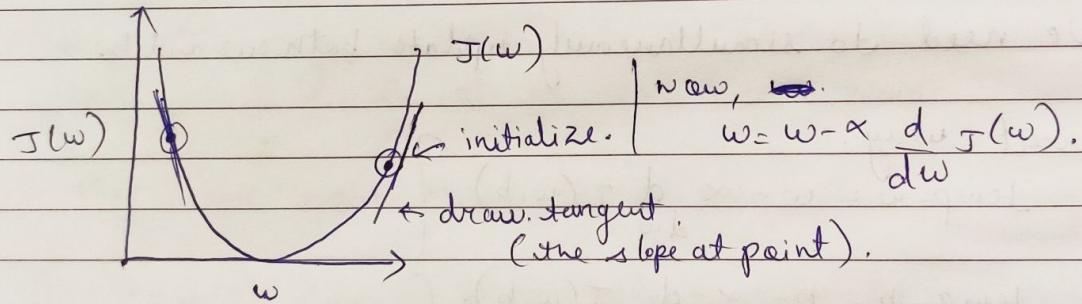
$$b = b - \alpha \frac{d}{db} J(w, b)$$

An example:

$J(w)$ ; if  $J$  is a number.

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

This our example, where we set  $b=0$ .



if  $\frac{d}{dw} J(w) > 0$ ,  $w = w - \alpha \times \text{positive number}$ .

so we are getting closer to min. value.

If have starting point in the left. The slope of this line will have negative slope.

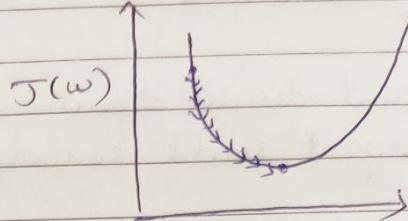
So,  $w = w - \alpha \times (\text{negative number})$ .

So, subtracting a -ve no., we adding a +ve number to  $w$ . So, we are moving in the right direction.

### III. Learning Rate.

Now, we have  $w = w - \alpha \frac{d}{dw} J(w)$ .

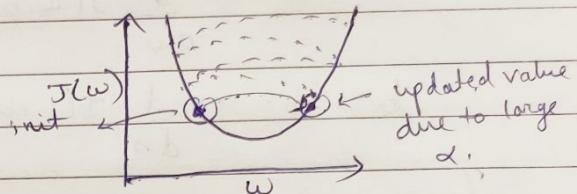
If  $\alpha$  is too small, we take a very small step. We reach bottom but very slowly.



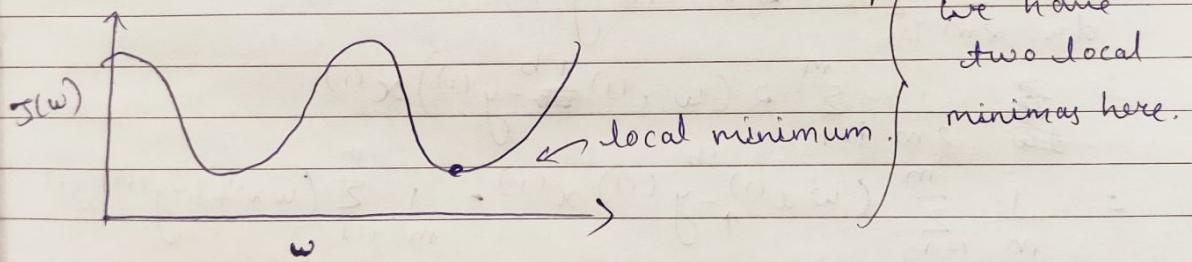
If  $\alpha$  is too large, we might keep on increasing the cost function.

Gradient descent may overshoot and never reach the minimum.

It may fail to converge, or even diverge.



- What if we already reach a local minimum?



So, we have current value of  $w$ .

Since slope is 0,  $w = w - \alpha \times 0 \rightarrow$  set  $w$  to  $w$ .

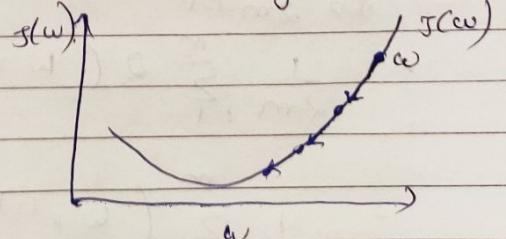
$w$  value is unchanged.

Can reach local minimum with fixed learning rate ( $\alpha$ ).

~~Let  $\alpha$  be large.~~

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

as we approach the minimum, the derivative approaches 0. Update steps becomes smaller.



## IV. Gradient Descent for Linear Regression.

Linear regression model.  
 $f_{w,b}(x) = wx + b$

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

Gradient descent algorithm.  
repeat until convergence {

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

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

$$\begin{aligned} \text{Now, } \frac{d}{dw} J(w, b) &= \frac{d}{dw} \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 \\ &= \frac{d}{dw} \frac{1}{2m} \sum_{i=1}^m 2(w x^{(i)} + b - y^{(i)})^2 \\ &= \cancel{\frac{1}{2m}} \sum_{i=1}^m 2(w x^{(i)} + b - y^{(i)}) \cancel{x^{(i)}} \\ &= \frac{1}{m} \sum_{i=1}^m (w x^{(i)} + b - y^{(i)}) x^{(i)} = \frac{1}{m} \sum_{i=1}^m (wx^{(i)} + b - y^{(i)}) x^{(i)} \end{aligned}$$

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

$$\begin{aligned} &= \frac{d}{db} \frac{1}{2m} \sum_{i=1}^m (w x^{(i)} + b - y^{(i)})^2 \\ &= \frac{1}{2m} \sum_{i=1}^m 2(w x^{(i)} + b - y^{(i)}) x_1 \end{aligned}$$

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

So, the gradient descent algorithm for linear regression.

repeat until convergence {

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

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

we need to update  $w$  and  $b$  simultaneously.

As we already saw,  $J(w, b)$  may have more than 1 local minima.

But when using squared error cost function, our  $J$  function is bowl shaped. So it will have only one local minima.  
It is a convex function.

Convex function is a bowl shaped function. It cannot ~~slope~~ have any local minima other than a single global minima.

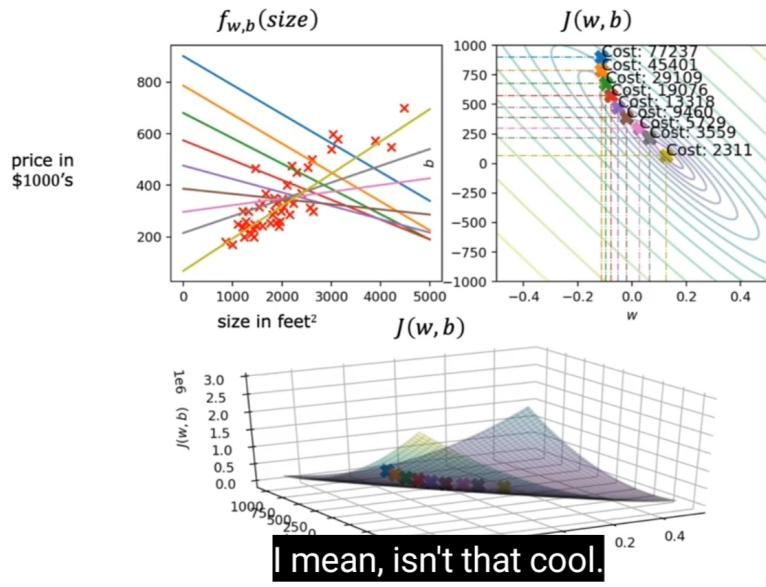
## I. Running Gradient Descent.

We can have a look at our 3D plot for squared error cost function.

As we initialize it at some value, we can have a look at the figure.

This is actually called "Batch" Gradient Descent.

"Batch": Each step of gradient descent uses all the training examples.



for updating value of  $w$  and  $b$ , we are looking at the entire batch of training data.

$$\sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2. \quad \{ m = 47 \}$$

other gradient descent algorithms also uses subsets rather than the whole batch for training.

## VI. Quiz.

1. Gradient descent is an algorithm for finding values of parameters  $w$  and  $b$  that minimize the cost function  $J$ .

$$\text{repeat until } \begin{cases} w = w - \alpha \frac{d}{dw} J(w, b) \\ b = b - \alpha \frac{d}{db} J(w, b). \end{cases}$$

when  $d/dw J(w, b)$  is a negative number, what happens to  $w$  after one update step?

- a)  $w \uparrow$
- b) Not possible to tell
- c)  $w \downarrow$

d) stays same

→ Value of  $w$  increases.

2. For linear regression, what is the update step for parameter  $b$ ?

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

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

→ a.