

Gradient descent:

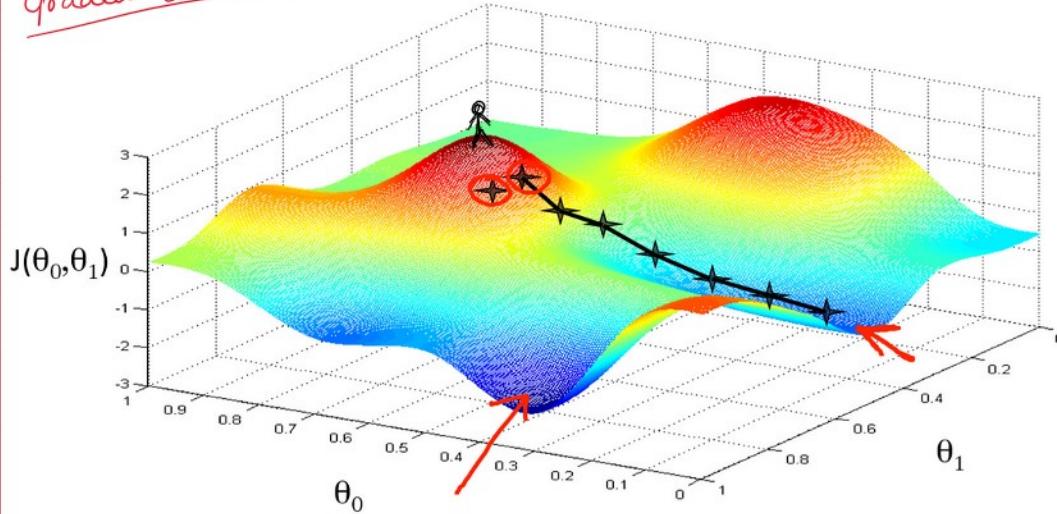
Monday, August 28, 2023 11:44 AM

gradient descent :-

WE USE GRADIENT DESCENT TO KEEP CHANGING OR ADJUSTING THE PARAMETER $J(w,b)$ TO MINIMIZE THE "ERROR" OR "COST FUNCTION".

- This helps the model get better making at predictions, classifications or other tasks based on the given data
- we keep changing the parameters (w, b) to reduce $J(w,b)$ until we settle at or near a minimum.

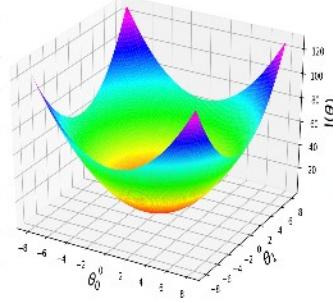
Gradient descent



→ this type of Cost Function you only get when you training a neural network.

→ Above function is not a squared error cost

For linear regression with Squared Error Cost Function, you will always end up with graph of hammock shape like this.



* Gradient Descent Algorithm *

$$W = W - \alpha \frac{\partial J(w,b)}{\partial w} \quad \begin{matrix} \rightarrow \text{derivative} \\ \rightarrow ① \end{matrix}$$

$\alpha \Rightarrow$ learning rate, small +ve number which 0 and 1.

→ What α basically does?

→ It basically controls how big of a step you take downhill.

→ If α is big then it will be aggressive steps downhill.

→ " " " small " " " " small baby steps downhill.

$$b = b - \alpha \frac{\partial J(w,b)}{\partial b} \quad ②$$

* We will be constantly repeat both the function until convergence.

→ We will be keep updating the parameters w and b .

* Correct simultaneous update *

* Correct Simultaneous update *

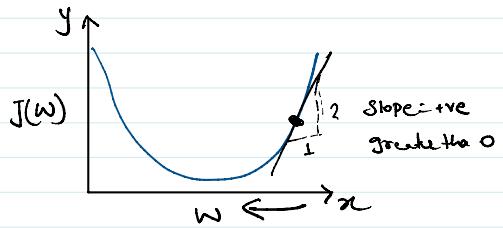
$$\text{tmp_W} = W - \alpha \frac{\partial}{\partial w} J(w, b)$$

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

$$W = \text{tmp_W}$$

$$b = \text{tmp_b}$$

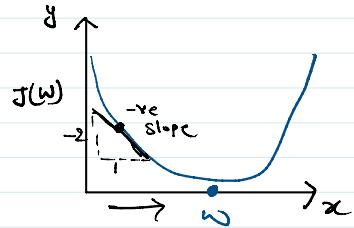
* Gradient descent Intuition *



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

$$w = w - \alpha \cdot (+ve\ number)$$

⇒ you will move to the left while performing this. and you decreasing Cost of $J(w)$ and you are getting closer to minima.



$$\frac{dJ(w)}{dw}$$

$$w = w - \alpha \cdot (-ve\ number)$$

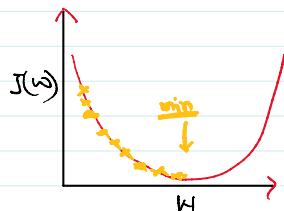
⇒ you will end up increasing w which mean you moving to the right to the graph and your Cost J decreasing which gets you closer to minimum.

* Learning rate (α) *

* Two Cases of learning Rate :-

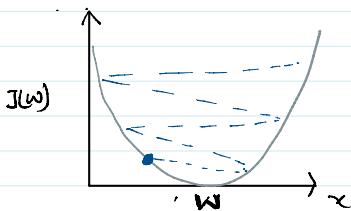
1. If α is too small.
2. If α is too big.

* If learning rate is small (α)



⇒ If the learning rate is small then the G.T will work but it will be slow. because here tiny steps takes long time to get to the minima.

* If α is too big (learning rate).

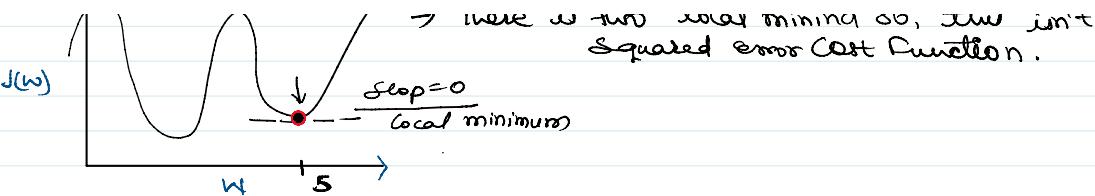


⇒ If α is large then G.T may overshoot and may never reach the minimum.

⇒ Fail to Converge and may even diverge



⇒ There are two local minima so, this isn't squared error Cost Function.



What IF Parameter w is already at local minima.

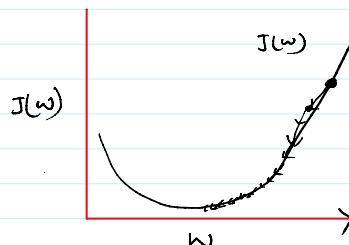
$$w_t = w - \alpha \cdot \frac{d J(w)}{d w} \rightarrow 0$$

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

$$\Rightarrow w_t = w \Rightarrow \text{unchanged}$$

Can reach local minima with fixed learning rate:

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



- Derivative becomes smaller
- update steps become smaller

Gradient descent for linear regression

* linear regression model

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

* Cost Function

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

* Gradient descent algorithm.

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

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

If we use these function instead of derivative term it will work as well

* Derivation of $\frac{\partial}{\partial w} J(w, b)$

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

$$\Rightarrow \frac{\partial}{\partial w} \frac{1}{2m} \sum_{i=1}^m (w x^{(i)} + b - y^{(i)})^2 / 2$$

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

* Derivative of $\frac{\partial}{\partial b} J(w, b)$

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

$$\Rightarrow \frac{\partial}{\partial b} \frac{1}{2m} \sum_{i=1}^m (w x^{(i)} + b - y^{(i)})^2 / 2$$

$$\Rightarrow \frac{\partial}{\partial b} \frac{1}{2m} \sum_{i=1}^m (w^T x^{(i)} + b - y^{(i)})^2$$

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

Batch \Rightarrow Search step of gradient descent uses all the training examples.