

Parameter Learning

Gradient Descent

- It is used for minimizing cost function.
- It is a more general algorithm and is used not only in linear regression but all over the different places in Machine Learning.
- Ex: - Have some function $J(\theta_0, \theta_1)$
Want $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

Outline \rightarrow

- Start with some θ_0, θ_1 (random).
- Keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$ until we hopefully end up at a minimum (local).

Gradient Descent Algorithm :-

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\text{for } j=0 \& j=1)$$

}

(Next page)

Also update θ_0 and θ_1 simultaneously: —

Simultaneous Update: —

$$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp0}$$

$$\theta_1 := \text{temp1}$$

↑
This is correct method.

$$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp0}$$

$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0 - \theta_1)$$

$$\theta_1 := \text{temp1}$$

↑
This is incorrect method, as you can see it uses updated θ_0 for cost function for getting temp1. This method may also work fine but it should be avoided as this is some other algorithm, not gradient descent. It can also cause some weird behaviour.

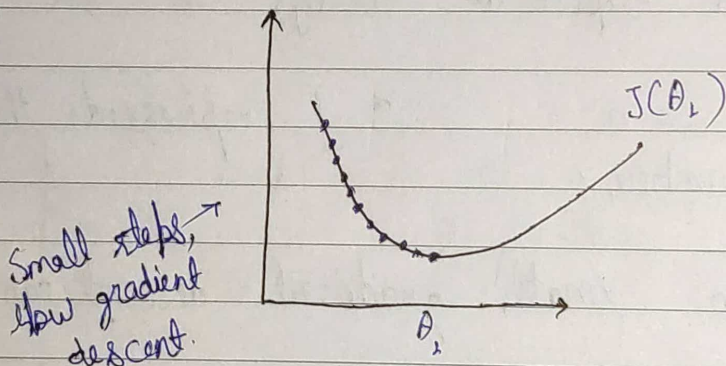
- α is learning rate. It ^{determines} ~~tells~~ how much the θ_0 or θ_1 or θ_n should be changed.
- $:=$ is assignment operator.
- Above algorithm (gradient descent) can be applied on cost function with any (n) no. of parameters by slightly modifying above formulae.
- $\frac{\partial}{\partial \theta_j}$ ~~is~~ is symbol of partial derivative with respect ~~to~~ to θ_j .
- $j = 0, 1, \dots, n$. j represents the feature index number.
- If α is too small, gradient descent can be ~~at~~ slow.
- If α is too large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge.
- If parameter θ_1 is at local minimum already, then gradient descent will leave it unchanged. Because we will have 0 slope, ~~so~~ ~~the~~ derivative will be ~~the~~ 0, hence parameter will be unchanged.

- Gradient descent converge to a local ~~min~~ minimum even with the learning rate α fixed. Because,

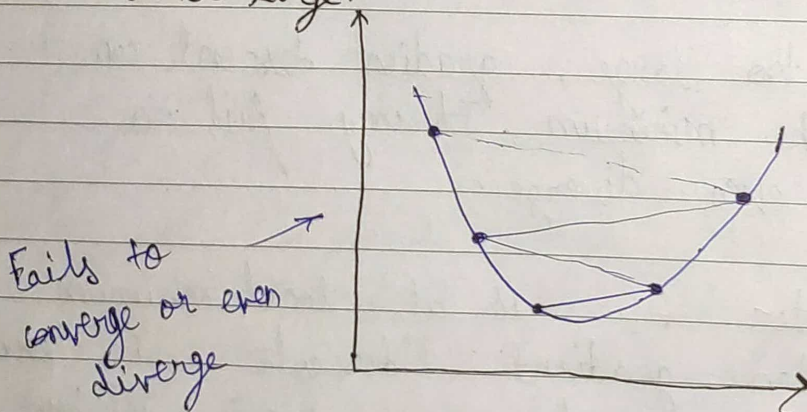
As we approach a local minimum, gradient descent will automatically take smaller steps.
So, no need to decrease α over time.

Examples:-

α too small :-



α too large :-



Applying Gradient Descent to our Cost Function:-

Gradient Descent Algorithm

repeat until convergence {

$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

(for $j=1$ and $j=0$)

Linear Regression Model

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

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

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) &= \frac{\partial}{\partial \theta_j} \cdot \frac{1}{2m} \sum_{i=1}^m (h_\theta(x_i) - y_i)^2 \\ &\Rightarrow \frac{\partial}{\partial \theta_j} \cdot \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x_i - y_i)^2 \end{aligned}$$

Now,

For:-

$$\theta_0 \rightarrow j=0 : \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x_i) - y_i)$$

$$\theta_1 \rightarrow j=1 : \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x_i) - y_i) \cdot x_i$$

So, our case:-

Gradient Descent:-

repeat until convergence $\left\{ \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) \right.$

$$\theta_0 := \theta_0 - \alpha \frac{1}{n} \sum_{i=1}^n (h_0(x_i) - y_i)$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{n} \sum_{i=1}^n (h_0(x_i) - y_i) \cdot x_i$$

}

$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

Also ~~update~~ update θ_0 & θ_1 simultaneously,

Note:- Cost function for linear regression is always going to be a bowl-shaped function.

Technical term for this is convex function.

Convex function means a bowl-shaped function, so this function doesn't have any local optima except for one global optimum.

Note: - Above implementation is also called Batch Gradient Descent.

It basically means that it uses all training examples in each step.

Other implementations of gradient descent which use subsets of data are also possible.