

Optimization Algorithms

Gradient descent :

It is Optimization algorithm used in machine learning to tweak the parameters iteratively to minimize the given cost function(differentiable function) to its local minimum.

The idea is to take repeated steps in the opposite direction of the [gradient](#) (or approximate gradient) of the function at the current point, because this is the direction of steepest descent.

Steps :

1. Start with some random values parameters.
2. Keep changing the parameters until we end up at a minimum value of cost function.

Algorithm :

Repeat until converges {

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

}

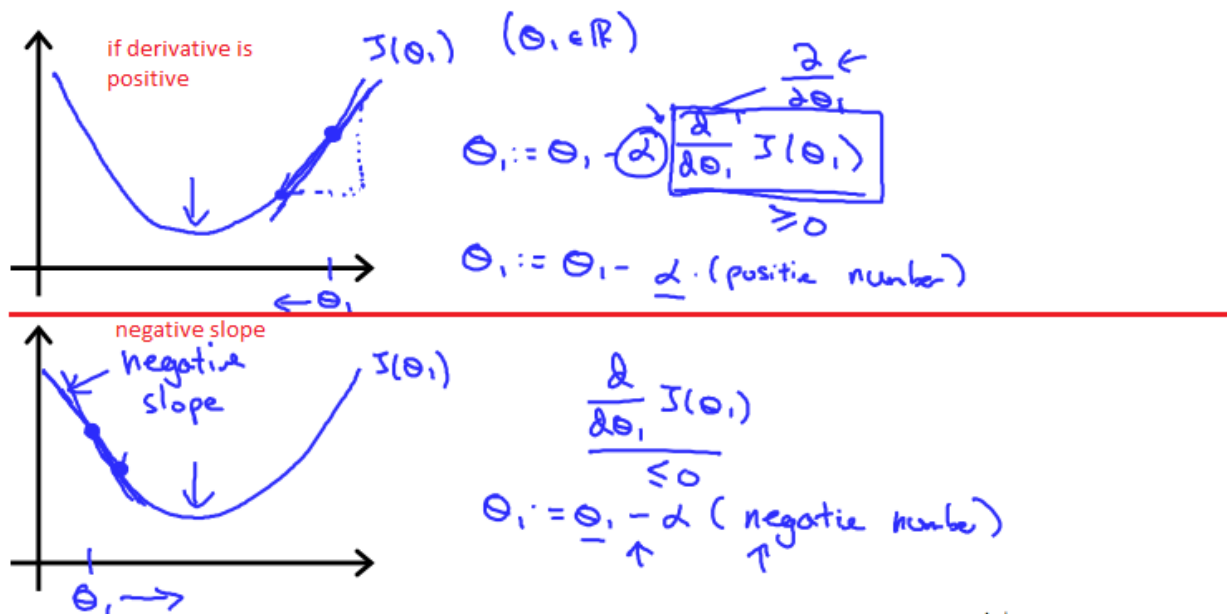
- Here “ α ” is the learning rate, it decides how big the step taken by algorithm to reach local minimum.
- $\frac{\partial J}{\partial \theta} J(\theta_0, \theta_1, \dots, \theta_n)$ is the derivative term.
- We must simultaneously update the parameter values. Suppose we have 2 parameters θ_0, θ_1 , then we must update the parameters as follows :

<u>Correct: Simultaneous update</u>	<u>Incorrect:</u>
$\rightarrow \text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$	$\rightarrow \text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$
$\rightarrow \text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$	$\rightarrow \theta_0 := \text{temp0}$
$\rightarrow \theta_0 := \text{temp0}$	$\rightarrow \text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$
$\rightarrow \theta_1 := \text{temp1}$	$\rightarrow \theta_1 := \text{temp1}$

- In the incorrect one the updated θ_0 will be used in updating the θ_1 which is an error, so we follow simultaneous update in gradient descent algorithm.

<https://coursera.org/share/4e98bc5d188ae89053e81c8101dca65a>

To make it simpler consider only θ_1 , now plot a graph between θ_1 and cost function J .



- If derivative term is positive then the next θ_1 value will be move towards left side to reach local minimum.
- If derivative term is negative then the next θ_1 value will be move towards right side to reach local minimum.
- If derivative term is zero then it reached the local minimum.

On the other side we also need to adjust the value of our learning Parameter to ensure that the gradient descent algorithm converges at a resonable time.

Failure to converge or too much time to obtain minimum value imply that our step size is wrong.

- If our learning parameter value is too small, then it takes more time to reach minimum.
- On the other hand if its value is large it will fail to converge.

Why there is no need to update the learning parameter along with the parameter ?

Because as we approach to the local minimum the differentiation value will be decreased and the total product value will also decreases significantly, that is why we can approach local minimum even with fixed value of α .

Note: For the specific choice of cost function $J(\theta_0, \theta_1)$ used in linear regression, there are no local optima (other than the global optimum).

Normal Equation :

Instead of running iterative algorithm, we can solve theta values analytically, i.e. we get weight parameters all in one go.

The equation is $\theta = (X^T X)^{-1} (X^T y)$

In the "Normal Equation" method, we will minimize J by explicitly taking its derivatives with respect to the θ_j 's and setting them to zero.

Derivation:

Suppose we have "m" training examples (x_i, y_i) where $i = 1..m$, and we have 'n' features for every training example i.e. $x^i = [x_i^1, x_i^2, \dots, x_i^n]$, we add $x_i^0 = 1$ (n X 1) vector, for easy matrix manipulations, like it acts as coefficients for our bias θ_0 , then our feature vector is of size $(n+1, 1)$.

Now combining all these feature vectors to form a matrix representation, i.e.

$$X = \begin{bmatrix} x_1^0 & \dots & x_1^n \\ \vdots & \ddots & \vdots \\ x_m^0 & \dots & x_m^n \end{bmatrix} \text{ a } (n+1, m) \text{ matrix and our corresponding weight matrix } \theta = \begin{bmatrix} \theta^0 \\ \vdots \\ \theta^n \end{bmatrix} \text{ an}$$

$n+1$ dimensional vector. Each row of X represents one data sample.

Now we need to find for what value of theta does $X \theta = y$,

So we in order to achieve that

$$e = (y - X\theta)^T (y - X\theta)$$

$$e = (y^T - \theta^T X^T) (y - X\theta)$$

$$e = (y^T y - y^T X\theta - \theta^T X^T y + \theta^T X^T X \theta)$$

$$\parallel (A^T B = B^T A) \\ \Rightarrow y^T X \theta = \theta^T X^T y$$

$$e = (y^T y - 2\theta^T X^T y + \theta^T X^T X \theta)$$

$$\frac{\partial e}{\partial \theta} = \frac{\partial}{\partial \theta} (y^T y - 2\theta^T X^T y + \theta^T X^T X \theta) = 0$$

$$\left(\frac{df}{dx} x^T b = b \right)$$

b - scalar
vector

B - Matrix

$$\left(\frac{df}{dx} x^T B x = B x \right)$$

$$= \frac{\partial}{\partial \theta} (y^T y) - 2 \frac{\partial}{\partial \theta} (\theta^T X^T y) + \frac{\partial}{\partial \theta} (\theta^T X^T X \theta)$$

$$\frac{\partial e}{\partial \theta} = 0 - 2X^T y + 2X^T X \theta$$

now

$$\frac{\partial e}{\partial \theta} = 0$$

$$\Rightarrow 2X^T X \theta - 2X^T y = 0$$

$$2X^T X \theta = 2X^T y$$

$$(X^T X)^{-1} (X^T X) \theta = (X^T X)^{-1} X^T y$$

$$\theta = (X^T X)^{-1} X^T y$$

References :

1. http://mlwiki.org/index.php/Normal_Equation
2. <https://www.youtube.com/watch?v=uoeit0FCWWA>
3. <http://www.rmi.ge/~kade/LecturesT.Kadeishvili/MathEconomics/Term3/Week3QuadraticLEC.pdf>
4. http://www.gatsby.ucl.ac.uk/teaching/courses/sntn/sntn-2017/resources/Matrix_derivatives_cribsheet.pdf

