

Gradient Descent (GD) Algorithm

Zengchang Qin (Ph.D.)

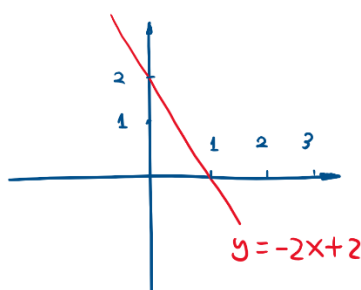
March 26, 2018

Given a database $D = \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(n)}, y^{(n)})\}$
Where $x^{(i)}$ is a vector in n -dimensional space and $y^{(i)}$ is a scalar
i.e.: $x_i \in \mathbb{R}^n$, $y_i \in \mathbb{R}$

We hope to learn $f(x^{(i)}) \rightarrow y^{(i)}$.

This is a typical **machine learning** problem

1) If our hypothesis of relation function is a linear model.



$$y = ax + b$$

$$y = \theta_1 x + \theta_0$$

if we set $x_0 = 1$

$$y = \theta_1 x_1 + \theta_0 x_0 \\ = \theta^T x$$

$$\theta^T = (\theta_0, \theta_1)$$

$$x = \begin{pmatrix} x_0 \\ x_1 \end{pmatrix}$$

In general form

$$y = \theta^T x \quad \begin{cases} \theta \in \mathbb{R}^{n+1} \\ x \in \mathbb{R}^{n+1} \end{cases}$$

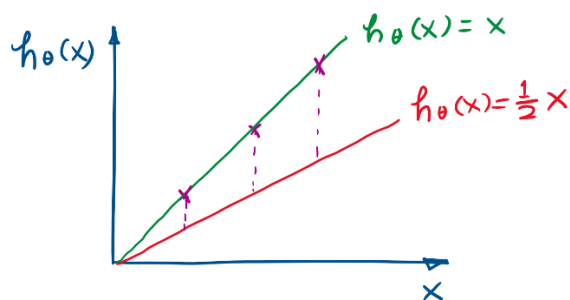
$$y^{(i)} = h_{\theta}(x^{(i)}) = \theta^T x^{(i)} = \sum_{j=1}^n \theta_j x_j^{(i)}$$

Let us consider the simplest case that

$$h_{\theta}(x) = \theta_1 x_1 + \theta_0, \text{ where } \theta_0 = 0$$

If given the training data

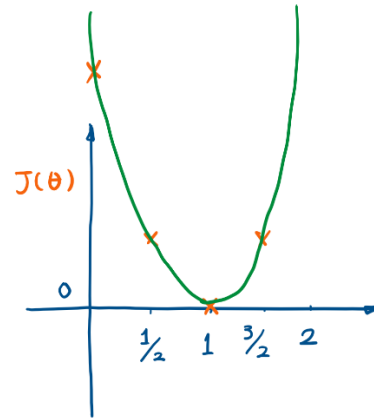
$x^{(1)} = 1$	$y^{(1)} = 1$
$x^{(2)} = 2$	$y^{(2)} = 2$
$x^{(3)} = 3$	$y^{(3)} = 3$



$$J(\theta) = \frac{1}{2N} \sum_{i=1}^N (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$J\left(\frac{1}{2}\right) = \frac{1}{2 \times 3} \left(\left(-\frac{1}{2}\right)^2 + (-1)^2 + \left(-\frac{3}{2}\right)^2 \right) \\ = \frac{1}{6} \left(\frac{1}{4} + \frac{4}{4} + \frac{9}{4} \right) = \frac{14}{24} = \frac{7}{12}$$

$$J(0) = \frac{1}{2 \times 3} \left((-1)^2 + (-2)^2 + (-3)^2 \right) \\ = \frac{1}{6} \times 14 = \frac{7}{3}$$

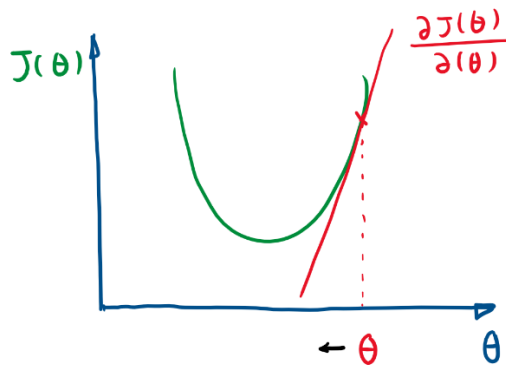


2) How to minimize Cost function $J(\theta)$?

$$J(\theta) = \frac{1}{2N} \sum_{i=1}^N (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$\text{If } h_{\theta}(x) = \theta x$$

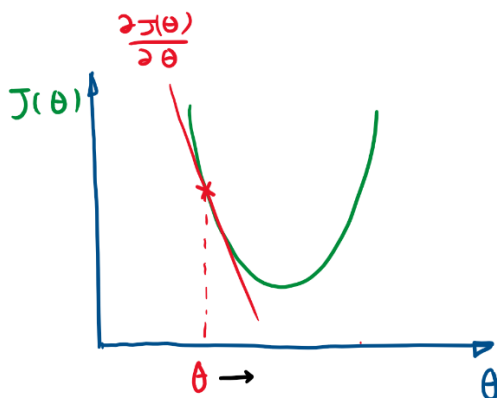
$$\frac{\partial J(\theta)}{\partial \theta} = \frac{\partial}{\partial \theta} \left(\frac{1}{2N} \sum_{i=1}^N (\theta x^{(i)} - y^{(i)})^2 \right) \\ = \frac{1}{N} \sum_{i=1}^N (\theta x^{(i)} - y^{(i)}) x^{(i)}$$



$$\theta \leftarrow \theta - \alpha \frac{\partial J}{\partial \theta}$$

if $\frac{\partial J}{\partial \theta} > 0$ (positive)

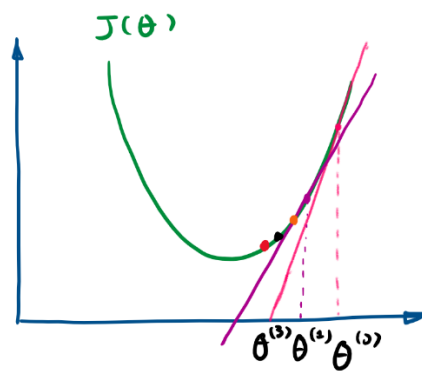
$$\theta \leftarrow \theta - \alpha (\text{positive number})$$



if $\frac{\partial J(\theta)}{\partial \theta} < 0$ (negative)

$$\theta \leftarrow \theta - \alpha (\text{negative number})$$

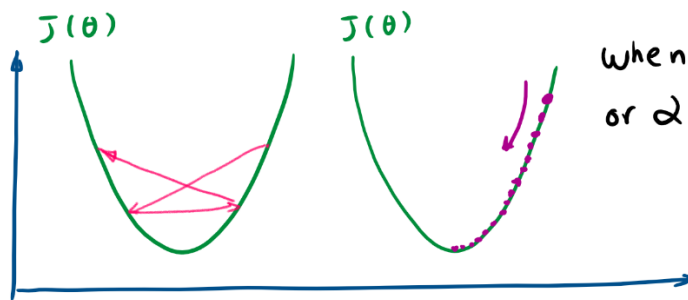
How Gradient change?



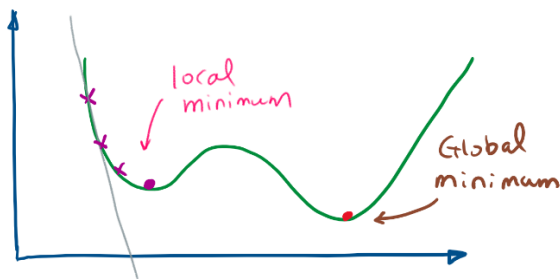
$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \alpha \frac{dJ(\theta)}{d\theta}$$

Gradient is diminishing in the process of iterative updating.

α can be fixed!



When α is too big.
or α is too small.



Gradient descent does not have the guarantee to converge to the global minimum. Be aware of the local minimum.

3) Gradient descent :

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

E.g. for the above case of $h_\theta(x) = \theta x$

$$\theta \leftarrow \theta - \alpha \frac{dJ(\theta)}{d\theta}$$

If $h_\theta(x) = \theta_1 x + \theta_0$, GD is updated *Simultaneously* by:

$$\begin{aligned} \theta_0 &\leftarrow \theta_0 - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0} \\ \theta_1 &\leftarrow \theta_1 - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1} \end{aligned} \quad (\alpha > 0) \text{ is the learning rate.}$$

$$\begin{aligned} \text{where } \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0} &= \frac{\partial}{\partial \theta_0} \left(\frac{1}{2N} \sum_{i=1}^N (h_\theta(x^{(i)}) - y^{(i)})^2 \right) \\ &= \frac{1}{N} \sum_{i=1}^N (h_\theta(x^{(i)}) - y^{(i)}) \end{aligned}$$

$$\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1} = \frac{1}{N} \sum_{i=1}^N (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}$$

4) For a more general case of multivariate linear regression where $h_\theta(x) = \theta^T x$

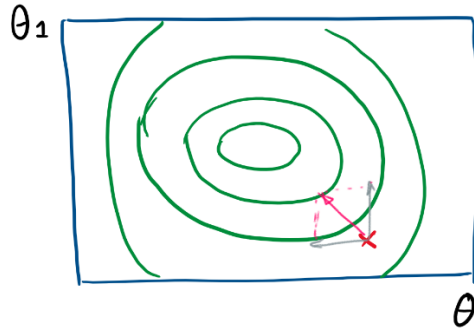
The gradient descent algorithm becomes :

$$\begin{aligned} \theta_j &\leftarrow \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j} \\ &= \theta_j - \alpha \frac{1}{N} \sum_{i=1}^N (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \end{aligned}$$

$$\begin{aligned} \theta_0 &\leftarrow \theta_0 - \alpha \frac{1}{N} \sum_{i=1}^N (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)} \quad x_0 = 1 \\ &= \theta_0 - \alpha \frac{1}{N} \sum_{i=1}^N (h_\theta(x^{(i)}) - y^{(i)}) \end{aligned}$$

A batch of data $(x_i : i=1 \dots N)$ is used for updating $\theta_j (j=1 \dots n)$

- 5) Contour with two variables (gradient is the direction, α is actually the step forward)

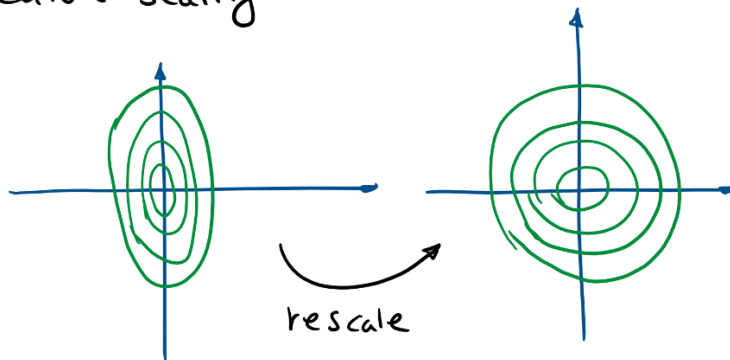


$$\alpha \frac{\partial J}{\partial \theta_1}$$

$$\alpha \frac{\partial J}{\partial \theta_0}$$

It is possible to be trapped into a local optima.
We can try different initialization

- 6) Feature Scaling



$$(1) x_i \leftarrow \frac{x_i}{\max(x_i) - \min(x_i)}$$

$$(2) x_i \leftarrow \frac{x_i - \bar{x}}{\max(x_i) - \min(x_i)}, \text{ where } \bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

- 7) A linear combination of features tells the importance of each feature towards the target.

$$y = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 \dots \theta_n x_n$$

↑
weight

$\theta_j > 0$, positive influence

$\theta_j < 0$, negative influence.