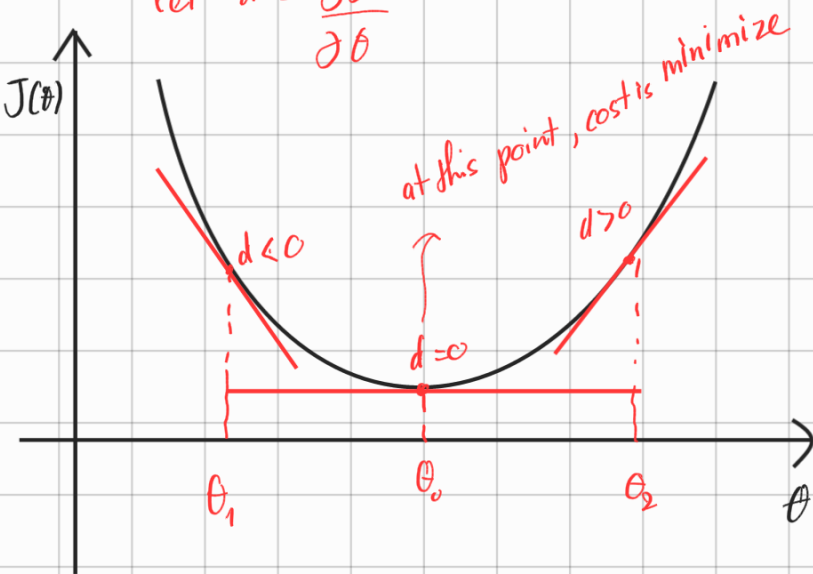


Gradient Descent

- The idea is If you change a little bit model parameter $\theta \rightarrow$ your cost function will change
- And How you change θ to make your cost reduce.

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

let $d = \frac{\partial J}{\partial \theta}$



• $d < 0$: If θ belong to this part then we want it to increase to θ_0 .

• $d > 0$: If θ belong to this part then we want to it to decrease to θ_0 .

- Then we can illustrate this as: $\theta = \theta - \alpha \cdot d$
 - + α : learning rate
 - + d : change and sign of θ

- First, we initialize value for θ
- Second, we choose appropriate α

+ If too large $\rightarrow \theta$ can not converge
+ If too low $\rightarrow \theta$ converge very slow

- Third, Calculate gradient descent (d)

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

$$\rightarrow \frac{\partial J}{\partial \theta_j} = \frac{2}{m} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)}) x_j^{(i)} \quad (j = 0 \dots n)$$

temporary forget this constant

For example with $j = 0$:

$$\frac{\partial J}{\partial \theta_0} = (\theta^T x^{(1)} - y^{(1)}) x_0^{(1)} + (\theta^T x^{(2)} - y^{(2)}) x_0^{(2)} + \dots + (\theta^T x^{(m)} - y^{(m)}) x_0^{(m)}$$

Rewritten in matrix form:

$$\text{gradient}^T = \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \vdots \\ \frac{\partial J}{\partial \theta_n} \end{bmatrix}_{n \times 1} = \begin{bmatrix} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)}) x_0^{(i)} \\ \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)}) x_1^{(i)} \\ \vdots \\ \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)}) x_n^{(i)} \end{bmatrix}_{n \times 1} \quad (I)$$

we have

$$X = \begin{bmatrix} 1 & x_1^1 & \dots & x_n^1 \\ 1 & x_1^2 & \dots & x_n^2 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^m & \dots & x_n^m \end{bmatrix}_{m \times n}$$

$$X\theta - \hat{y} = \begin{bmatrix} \theta_0 x_0^{(1)} + \theta_1 x_1^{(1)} + \dots + \theta_n x_n^{(1)} - y^{(1)} \\ \theta_0 x_0^{(2)} + \theta_1 x_1^{(2)} + \dots + \theta_n x_n^{(2)} - y^{(2)} \\ \vdots \\ \theta_0 x_0^{(m)} + \theta_1 x_1^{(m)} + \dots + \theta_n x_n^{(m)} - y^{(m)} \end{bmatrix}$$

$$\begin{bmatrix} \vdots \\ \theta_0 x_0^{(m)} + \theta_1 x_1^{(m)} + \dots + \theta_n x_n^{(m)} - y^{(m)} \\ \vdots \end{bmatrix}$$

• then (I) can be written as:

$$\text{gradient} = \underset{n \times m}{X^T} (\underset{m \times n}{X} \underset{n \times 1}{\theta} - \underset{m \times 1}{\hat{y}})$$

• Remember put $\frac{2}{m}$ again!

$$\text{gradient} = \frac{2}{m} X^T (X\theta - \hat{y})$$

• Finally, we have the formular:

$$\theta = \theta - \underset{\substack{\uparrow \\ \text{choosing}}}{\alpha} \cdot \text{gradient}$$

