



NAME: DUONG DUY CHIEN

## I. Theory:

**Hypothesis:**

$$h_{\theta}(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n \quad \text{with } x_0 = 1 \quad J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

**Cost function:**

**Gradient descent algorithm:**

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad \text{Update theta}$$

with  $m$  = number of samples  
 $j = 1, \dots, n$  : number of features

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

## II. Implement part supports multiple variables in vectorized version

**Calculate hypothesis:** with  $m$  samples  $x$  (each sample is multiple variables –  $n$  features).

Input:

$$X = \begin{bmatrix} -x^{(1)T} & - \\ -x^{(2)T} & - \\ \vdots & \vdots \\ -x^{(m)T} & - \end{bmatrix} \quad \text{with } x^{(i)T} = [x_0, x_1, \dots, x_n] \text{ } n \text{ features} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

$X$  (vector column) store  $m$  samples  $x$

$$\text{Hypothesis: } H = \sum_{i=1}^m h^{(i)} = \sum_{i=1}^m x^{(i)T} \theta = \begin{bmatrix} -x^{(1)T} \theta - \\ -x^{(2)T} \theta - \\ \vdots \\ -x^{(m)T} \theta - \end{bmatrix} = \begin{bmatrix} -\theta^T x^{(1)} - \\ -\theta^T x^{(2)} - \\ \vdots \\ -\theta^T x^{(m)} - \end{bmatrix}$$

$h_{\theta}(x^i) = \theta^T x = x^T \theta$   
 Note:  $\theta^T x = x^T \theta$   
 or  $\theta^T x^{(i)} = x^{(i)T} \theta$   
 when  $x^{(i)}$  and  $\theta$  are vectors

**Cost function:**  $J(\theta) = \frac{1}{2m} (X\theta - \vec{y})^T (X\theta - \vec{y})$

**Gradient descent:**

$$\begin{aligned} \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \frac{\partial J}{\partial \theta_2} \\ \vdots \\ \frac{\partial J}{\partial \theta_n} \end{bmatrix} &= \frac{1}{m} \begin{bmatrix} \sum_{i=1}^m ((h_{\theta}(x^{(i)}) - y^{(i)})x_0^{(i)}) \\ \sum_{i=1}^m ((h_{\theta}(x^{(i)}) - y^{(i)})x_1^{(i)}) \\ \sum_{i=1}^m ((h_{\theta}(x^{(i)}) - y^{(i)})x_2^{(i)}) \\ \vdots \\ \sum_{i=1}^m ((h_{\theta}(x^{(i)}) - y^{(i)})x_n^{(i)}) \end{bmatrix} \quad \text{where} \end{aligned}$$

$$h_{\theta}(x) - y = \begin{bmatrix} h_{\theta}(x^{(1)}) - y^{(1)} \\ h_{\theta}(x^{(2)}) - y^{(2)} \\ \vdots \\ h_{\theta}(x^{(m)}) - y^{(m)} \end{bmatrix}.$$

$$\sum_i \beta_i x^{(i)} = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & & | \\ | & | & & | \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_m \end{bmatrix} = X^T \beta,$$

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

$$= \frac{1}{m} X^T (h_{\theta}(x) - y).$$

where the values  $\beta_i = (h_{\theta}(x^{(i)}) - y^{(i)})$ .

Gradient descent =  $X^T \beta$

### III. Results:

**Part 1:** One variable

- When Gradient descent is working correctly, cost function should be decrease with each step.

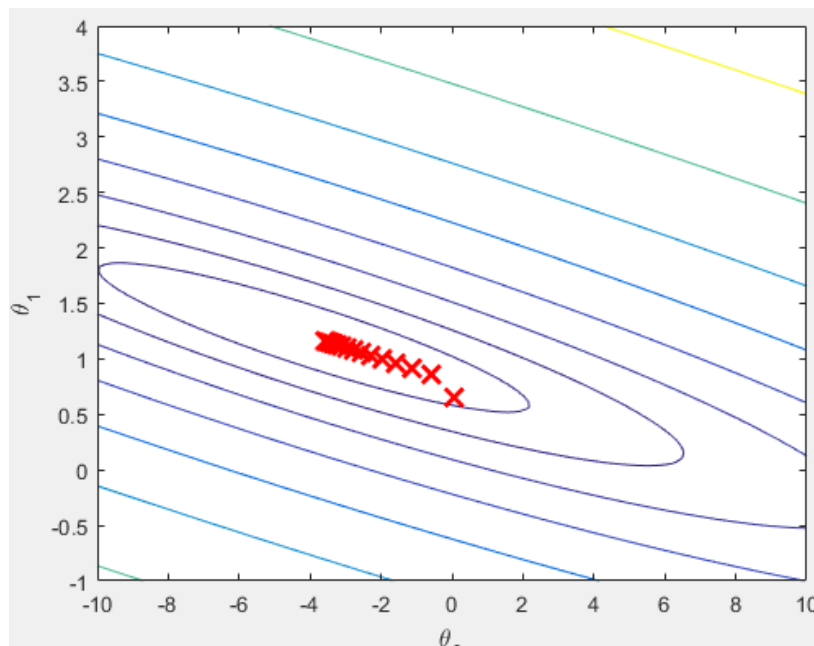


Figure 1-Gradient descent work correctly so the cost function reduce and converge at the global minimum .

- When we find the good parameter theta. The linear regression fit with our data

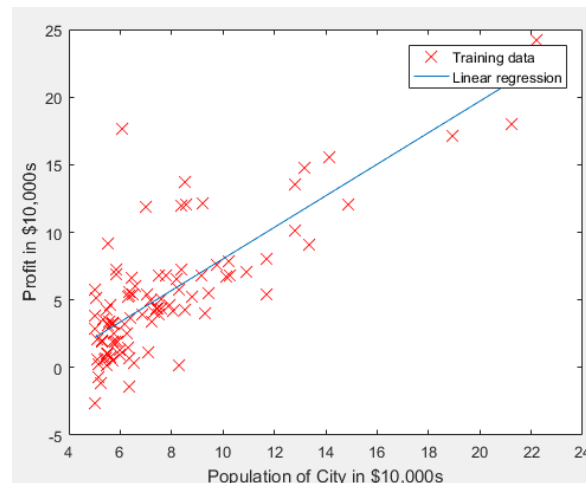


Figure 2- Training data with linear regression fit

## Part2:

**Learning rate** are very important in training. If we choose small learning rate, gradient descent takes a very long time to converge to the optimal value. In the other hand, with a large learning rate, gradient descent might not converge or might even diverge (as shown in Figure.3 and Figure 4)

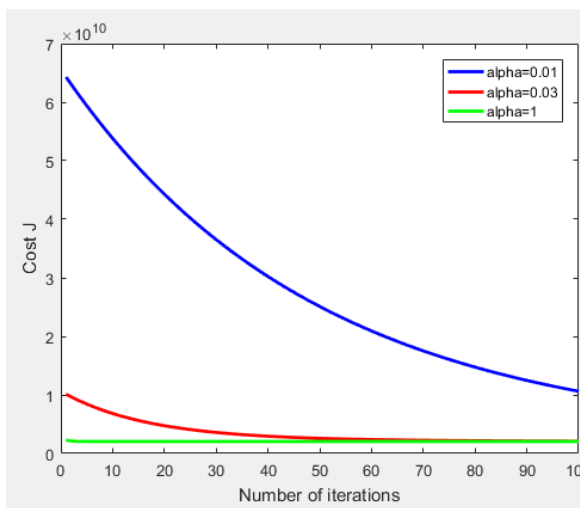


Figure 3 Cost function with different leaning rate  $\alpha$

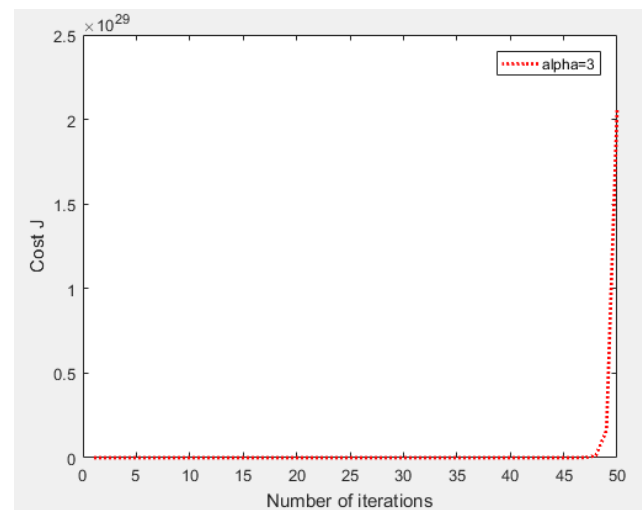


Figure 4 Gradient descent **diverge** when **leaning rate** is **too large**.

Testing: Predicted price of a 1650 sq-ft, 3 bedroom house (using gradient descent):

\$293081.464335