

# Gradient Descent Algorithm

Dr. Muhammad Wasim

# Intuition behind Gradient Descent Algorithm

$$\min_{\theta_1} J(\theta_1)$$

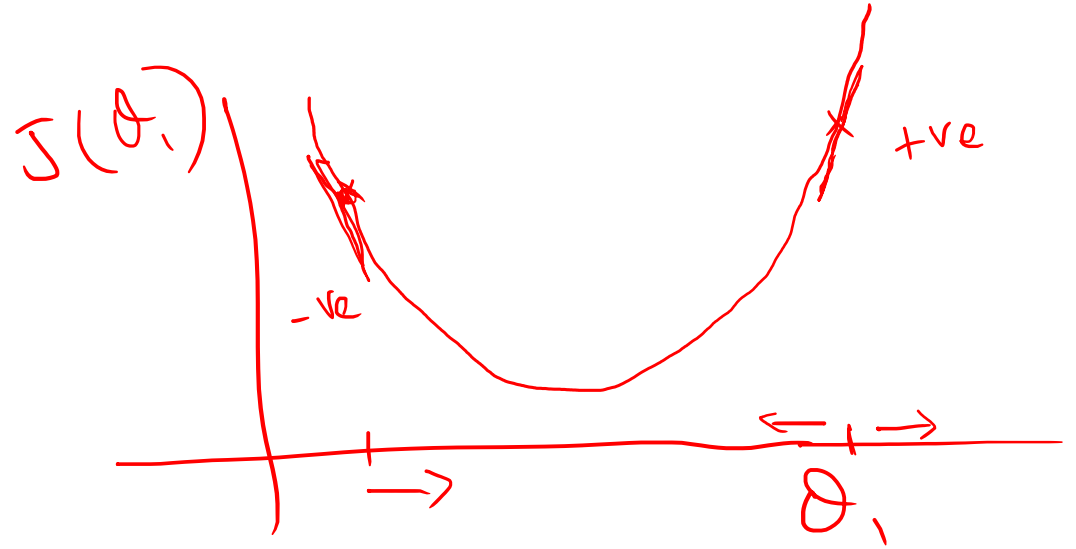
$$\theta_1 \in \mathbb{R}$$

$$\theta_1 = \theta_1 - \alpha \frac{d}{dx} J(\theta_1)$$

$-\alpha \times -ve = +ve$

$\alpha$

$+ve$



# The role of Alpha – The Learning Rate

$$\theta_1 = \theta_1 - \alpha \frac{d}{dx} J(\theta_1)$$

If  $\alpha$  is too small, gradient descent can be slow

If  $\alpha$  is too large, gradient descent can overshoot the minimum. It may fail to converge



# Computing the Derivative

$$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

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

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

# Gradient Descent for Linear Regression with one Feature

Repeat until convergence {

$$\underline{\theta_0} = \underline{\theta_0} - \underline{\alpha} \left[ \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \right]$$

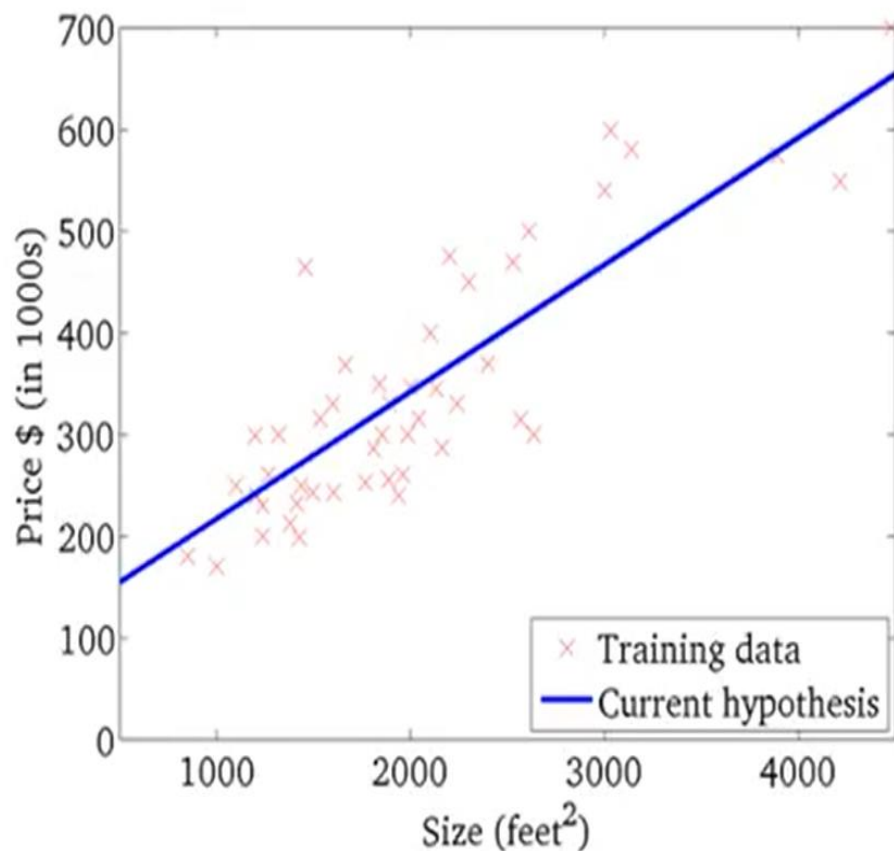
$$\underline{\theta_1} = \underline{\theta_1} - \underline{\alpha} \left[ \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x^{(i)} \right]$$

}

# How gradient descent for regression works

$$h_{\theta}(x)$$

(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



$$J(\theta_0, \theta_1)$$

(function of the parameters  $\theta_0, \theta_1$ )

