

## Gradient Descent for Multiple Variables

**Hypothesis:**  $h_{\theta}(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$

Our formal hypothesis in multivariable linear regression where we've adopted the convention that  $x_0 = 1$ . The parameters of this model are  $\theta_0$  through  $\theta_n$ , but instead of thinking of this as  $n$  separate parameters, which is valid, we are going to think of the parameters as  $\theta$  where theta is a  $n+1$ -dimensional vector.

**Parameters:**  $\theta_0, \theta_1, \dots, \theta_n = \theta$  ( $n + 1$  dimensional vector)

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

Our cost function is  $J$  of  $\theta_0$  through  $\theta_n$  which is given by this usual sum of square of error term. But again instead of thinking of  $J$  as a function of these  $n + 1$  numbers, we're going to more commonly use  $J$  as just a function of the parameter vector  $\theta$ .

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

**Video Question:** When there are  $n$  features, we define the cost function as

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

For linear regression, which of the following are also equivalent and correct definitions of  $J(\theta)$ ?

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

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m ((\sum_{j=0}^n \theta_j x_j^{(i)}) - y^{(i)})^2 \text{ (Inner sum starts at 0)}$$

- $J(\theta) = \frac{1}{2m} \sum_{i=1}^m ((\sum_{j=1}^n \theta_j x_j^{(i)}) - y^{(i)})^2$  (Inner sum starts at 1)
- $J(\theta) = \frac{1}{2m} \sum_{i=1}^m ((\sum_{j=0}^n \theta_j x_j^{(i)}) - (\sum_{j=0}^n y_j^{(i)}))^2$

The gradient descent equation itself is generally the same form; we just have to repeat it for our ' $n$ ' features:

$$\begin{aligned} &\text{repeat until convergence: } \{ \\ &\quad \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_0^{(i)} \\ &\quad \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_1^{(i)} \\ &\quad \theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_2^{(i)} \\ &\quad \dots \\ &\quad \} \end{aligned}$$

In other words:

$$\begin{aligned} &\text{repeat until convergence: } \{ \\ &\quad \theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} \quad \text{for } j := 0 \dots n \\ &\quad \} \end{aligned}$$

The following image compares gradient descent with one variable to gradient descent with multiple variables:

## Gradient Descent

Previously ( $n=1$ ):

Repeat {

$$\rightarrow \theta_0 := \theta_0 - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})}_{\frac{\partial}{\partial \theta_0} J(\theta)}$$

$$\rightarrow \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

(simultaneously update  $\theta_0, \theta_1$ )

}

New algorithm ( $n \geq 1$ ):

Repeat {

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

(simultaneously update  $\theta_j$  for  $j = 0, \dots, n$ )

}

$$\rightarrow \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\rightarrow \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

$$\rightarrow \theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_2^{(i)}$$

...