

Sid

Multiple features (variables)

n = # of features

$x^{(i)}$ = input features of i^{th} training example - vector

$x_j^{(i)}$ = value of feature j in i^{th} training example

Hypothesis : $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n$

$x_0 = 1$ ($x_0^{(i)} = 1$)

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

$$h_{\theta}(x) = \theta^T \cdot x$$

$$\begin{bmatrix} \theta_0 & \theta_1 & \dots & \theta_n \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

Hypothesis : $h_{\theta}(x) = \theta^T \cdot x = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$ ↗ $x_0 = 1$

Parameters : $\vec{\theta}$ ($n+1$ -dimensional vector)

$$\begin{aligned} \text{Cost Function : } J(\vec{\theta}) &= \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \\ &= \frac{1}{2m} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)})^2 \quad \text{or} \quad \frac{1}{2m} \sum_{i=1}^m \left[\left(\sum_{j=0}^n \theta_j x_j^{(i)} \right) - y^{(i)} \right]^2 \end{aligned}$$

Gradient Descent : $\theta_j := \theta_j - \alpha \frac{d}{d\theta_j} J(\vec{\theta})$

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

Feature Scaling $-1 \leq x_i \leq 1$, $x_0 = 1$

Mean normalizing: Replace x_i w/ $x_i - \mu_i$

$$x_i \leftarrow \frac{x_i - \mu_i}{s_i}$$

$J(\theta)$ should decrease after every iteration

Convergence if $J(\theta)$ decreased by less than 10^{-3} / iteration

Normal Equation: Solve for $\vec{\theta}$ analytically

$$X = \begin{bmatrix} \end{bmatrix}$$

$m \times (n+1)$

$$\vec{y} = \begin{bmatrix} \end{bmatrix}$$

m -dimensional vector

$$\theta = (X^T X)^{-1} X^T y$$

$$\text{pinv}(X' * X) * X' * y$$

m examples: $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$

n features

$$\vec{x}^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \quad X = \begin{bmatrix} \text{---} (x^{(1)})^T \text{---} \\ \text{---} (x^{(2)})^T \text{---} \\ \vdots \\ \text{---} (x^{(m)})^T \text{---} \end{bmatrix} \quad \vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

Gradient Descent

- Choose α
- many iterations
- ✓ for large n

Normal Equation

- no need to choose α
- no iterations
- need to compute $(X^T X)^{-1}$
 $O(n^3)$
- ✗ for large n