

Machine Learning - Week 2

1. $\begin{cases} X_j^{(i)} : \text{the value of feature } j \text{ in the } i^{\text{th}} \text{ training example} \\ X^{(i)} : \text{the input (features) of the } i^{\text{th}} \text{ training examples} \\ m : \text{the number of training examples} \\ n : \text{the number of features} \end{cases}$

eg. $X = \begin{bmatrix} 1 & 34 & 78 \\ 1 & 30 & 43 \\ \vdots & \vdots & \vdots \\ 1 & 94 & 89 \end{bmatrix}$ $m \leftarrow \text{training examples}$

$\begin{matrix} \uparrow & \uparrow & \uparrow \\ X_0 & X_1 & X_2 \end{matrix} \leftarrow \text{features}$

assume: $X_0^{(i)} = 1$

$$h_{\theta}(X) = \theta_0 + \theta_1 X_1 + \theta_2 X_2 + \dots + \theta_n X_n$$

$$= [\theta_0 \ \theta_1 \ \theta_2 \ \dots \ \theta_n] \begin{bmatrix} X_0 \\ X_1 \\ \vdots \\ X_n \end{bmatrix} = \theta^T X$$

2. gradient descent algorithm

① $n=1$:

Repeat $\left\{ \begin{array}{l} \theta_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta) \\ \theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta) \end{array} \right.$

★ update θ_0, θ_1 at the same time

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(X^{(i)}) - y^{(i)})$$
$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(X^{(i)}) - y^{(i)}) X^{(i)}$$

② $n \geq 1$:

Repeat { *simultaneously update θ_j for $j=0, \dots, n$*

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

}

$$\Rightarrow \begin{cases} \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h\theta(x^{(i)}) - y^{(i)}) \underline{x_0^{(i)}} = 1 \\ \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h\theta(x^{(i)}) - y^{(i)}) x_1^{(i)} \\ \dots \end{cases}$$

3. feature scaling

$$x_i := \frac{x_i - \mu_i}{s_i} \rightarrow \text{average}$$

s_i \rightarrow { range of value : max - min
standard deviation

4. learning rate

sd if α is too small : slow convergence

② if α is too large : $J(\theta)$ may not decrease on every iteration
 \Rightarrow may not converge

choose α : try $\dots, 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1, \dots$

3x bigger \rightarrow

4. normal equation

$$X^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1} \Rightarrow X = \begin{bmatrix} \dots (X^{(1)})^T \dots \\ \dots (X^{(2)})^T \dots \end{bmatrix}$$

\swarrow
design matrix

$$\begin{bmatrix} \vdots \\ X_n^{(i)} \end{bmatrix}$$

$$\begin{bmatrix} \vdots \\ \dots (X^{(m)})^T \dots \end{bmatrix}$$

$m \times (n+1)$

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

normal equations (no feature scaling)

5. Compare gradient descent with normal equation

gradient descent:

- ① need to choose α
- ② need many iterations
- ③ $O(n^2)$
- ④ Works well when n is large

normal equation

- ① don't need to choose α
- ② don't need to iterate
- ③ $O(n^2)$, need to calculate inverse of $X^T X$
- ④ slow when n is very large

if $X^T X$ is noninvertible:

- ① Redundant features, where two features are very closely related. (eg. linearly dependent)
- ② Too many features. (eg. $m \leq n$) \rightarrow delete some features

6. Some more information about the normal equation (optional)

(1) matrix derivatives:

$$\textcircled{1} \nabla_A f(A) = \begin{bmatrix} \frac{\partial f}{\partial A_{11}} & \dots & \frac{\partial f}{\partial A_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial A_{m1}} & \dots & \frac{\partial f}{\partial A_{mn}} \end{bmatrix}$$

$$\text{eg } f(A) = \frac{3}{2} A_{11} + 5 A_{12}^2 + A_{21} A_{22}$$

$$\Rightarrow \nabla_A f(A) = \begin{bmatrix} \frac{3}{2} & 10 A_{12} \\ A_{22} & A_{21} \end{bmatrix}$$

$$\textcircled{2} \text{tr} A = \sum_{i=1}^n A_{ii} \longrightarrow \text{trace operator}$$

$$\left\{ \begin{array}{l} \text{tr} AB = \text{tr} BA \\ \text{tr} ABC = \text{tr} CAB = \text{tr} BCA \\ \text{tr} ABCD = \text{tr} DABC = \text{tr} CDAB = \text{tr} DCBA \\ \text{tr} A = \text{tr} A^T \\ \text{tr}(A+B) = \text{tr} A + \text{tr} B \\ \text{tr} \alpha A = \alpha \text{tr} A \\ \nabla_A \text{tr} AB = B^T \\ \nabla_A \text{tr} ABH^T C = CBA + C^T A B^T \\ \nabla_A |A| = |A| (A^T)^T \end{array} \right.$$

(2) normal equation

$$X = \begin{bmatrix} -(X^{(1)})^T \\ -(X^{(2)})^T \\ \vdots \\ -(X^{(m)})^T \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$$\Rightarrow X\theta = \bar{y} \quad \left[(X^{(1)})^T \theta \right] = [y^{(1)}] = [h_\theta(x^{(1)}) - y^{(1)}]$$

$$X^T J = \begin{bmatrix} (X^{(1)})^T \theta \\ (X^{(2)})^T \theta \\ \vdots \\ (X^{(m)})^T \theta \end{bmatrix} \quad \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix} = \begin{bmatrix} h_\theta(X^{(1)}) - y^{(1)} \\ \vdots \\ h_\theta(X^{(m)}) - y^{(m)} \end{bmatrix}$$

$$\Rightarrow \frac{1}{2} (X\theta - \vec{y})^T (X\theta - \vec{y}) = \frac{1}{2} \sum_{i=1}^m (h_\theta(X^{(i)}) - y^{(i)})^2 = J(\theta)$$

$$\because \nabla_A^T \text{tr} ABA^T C = B^T A^T C^T + BA^T C$$

$$\therefore \nabla_\theta J(\theta) = \nabla_\theta \frac{1}{2} (X\theta - \vec{y})^T (X\theta - \vec{y}) = X^T X\theta - X^T \vec{y}$$

$$\therefore \text{minimize } J(\theta)$$

$$\therefore \nabla_\theta J(\theta) = 0$$

$$\Rightarrow \underbrace{X^T X \theta = X^T \vec{y}}_{\text{normal equations}} \Rightarrow \theta = (X^T X)^{-1} X^T \vec{y}$$