

Linear Regression with Ordinary Least Square Method

Zengchang Qin (PhD)

2018年4月15日
0:00

LEAST SQUARE

1. Given a training data $D = \{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$
If the model we hope to learn is a linear model:

$$h_{\theta}(x^{(i)}) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

The function we need to learn can map $x^{(i)} \rightarrow y^{(i)}$, we need to find the best parameter θ to satisfy the following equations:

$$\begin{cases} \theta_0 + \theta_1 x_1^{(1)} + \theta_2 x_2^{(1)} + \dots + \theta_n x_n^{(1)} = y^{(1)} \\ \theta_0 + \theta_1 x_1^{(2)} + \theta_2 x_2^{(2)} + \dots + \theta_n x_n^{(2)} = y^{(2)} \\ \vdots \\ \theta_0 + \theta_1 x_1^{(N)} + \theta_2 x_2^{(N)} + \dots + \theta_n x_n^{(N)} = y^{(N)} \end{cases}$$

They can be written into the form of matrices:

$$\underbrace{N \left\{ \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_1^{(N)} & x_2^{(N)} & \dots & x_n^{(N)} \end{bmatrix} \right\}}_{n+1} \cdot \underbrace{\begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}}_{n+1} \approx \underbrace{\begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{bmatrix}}_N$$

In order to find the best approximation, we can solve the following optimization problem:

$$X\theta \approx Y \Rightarrow \min \|\theta - \hat{\theta}\|_2^2 \text{ (Square of } L_2 \text{ norm)}$$

To solve that, we set:

$$S = \|X\theta - Y\|_2^2 = 0$$

$$\begin{aligned}\|X\theta - Y\|_2^2 &= (X\theta - Y)^T (X\theta - Y) \\ &= (\theta^T X^T - Y^T) (X\theta - Y) \\ &= \theta^T X^T X \theta - \theta^T X^T Y - Y^T X \theta + Y^T Y \\ &= \theta^T X^T X \theta - 2\theta^T X^T Y + Y^T Y\end{aligned}$$

Because $\theta^T X^T Y$ is a scalar, so does $Y^T X \theta$

$$1 \left\{ \underbrace{\hspace{2cm}}_{n+1} \right\} \cdot n+1 \left\{ \underbrace{\hspace{2cm}}_N \right\} \cdot N \left\{ \underbrace{\hspace{1cm}}_1 \right\} = \underbrace{\hspace{1cm}}_1 \left\{ \hspace{1cm} \right\} 1$$

$$\frac{\partial S}{\partial \theta} = 0 \Rightarrow \frac{\partial (\theta^T X^T X \theta)}{\partial \theta} - 2 X^T Y = 0$$

Since $\frac{d(u^T v)}{dx} = \frac{du^T}{dx} v + \frac{dv^T}{dx} u$, then:

$$\frac{d(\theta^T \theta)}{d\theta} = \frac{d\theta^T}{d\theta} \theta + \frac{d\theta^T}{d\theta} \theta = 2\theta$$

Given A is a square matrix

$$\begin{aligned}\frac{d(\theta^T A \theta)}{d\theta} &= \frac{d\theta^T}{d\theta} A \theta + \frac{d(\theta^T A^T)}{d\theta} \theta \\ &= A \theta + A^T \theta = (A + A^T) \theta\end{aligned}$$

$X^T X$ is $(n+1) \times (n+1)$ square matrix

$$\frac{d(\theta^T X^T X \theta)}{d\theta} = X^T X \theta + (X^T X)^T \theta = 2X^T X \theta$$

$$\begin{aligned}\frac{\partial S}{\partial \theta} = 0 &\Rightarrow 2X^T X \theta - 2X^T Y = 0 \\ &\Rightarrow X^T X \theta = X^T Y \\ &\Rightarrow \theta = (X^T X)^{-1} X^T Y\end{aligned}$$