

# Linear Regression 01

Chen Gong

12 October 2019

数据集  $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$ , 其中  $x_i \in \mathbb{R}^p$ ,  $y_i \in \mathbb{R}$ ,  $i = 1, 2, \dots, N$ 。  
数据矩阵为: (这样可以保证每一行为一个数据点)

$$X = (x_1, x_2, \dots, x_N)^T = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{pmatrix} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{Np} \end{pmatrix}_{N \times p} \quad (0.1)$$

$$Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}_{N \times 1} \quad (0.2)$$

设拟合的函数为:  $f(w) = w^T x$ 。

## 1 最小二乘估计: 矩阵表示

很简单可以得到损失函数 (Loss function) 为:

$$L(w) = \sum_{i=1}^N \|w^T x_i - y_i\|^2 \quad (1.1)$$

$$= (w^T x_1 - y_1, w^T x_2 - y_2, \dots, w^T x_N - y_N) \begin{pmatrix} w^T x_1 - y_1 \\ w^T x_2 - y_2 \\ \vdots \\ w^T x_N - y_N \end{pmatrix} \quad (1.2)$$

其中:

$$(w^T x_1 - y_1, w^T x_2 - y_2, \dots, w^T x_N - y_N) = [(w^T x_1, w^T x_2, \dots, w^T x_N) - (y_1, y_2, \dots, y_N)] \quad (1.3)$$
$$= w^T X^T - Y^T$$

所以:

$$L(w) = (Xw - Y)^T (Xw - Y) \quad (1.4)$$
$$= w^T X^T X w - 2w^T X^T Y + Y^T Y$$

由于  $X^T X$  是一个半正定矩阵,  $L(w)$  是一个凸函数, 那么我要求的  $w$ , 可记为  $\hat{w} = \arg \min_w L(w)$ 。这是一个无约束优化问题, 可以通过求偏导解决。那么有:

$$\frac{\partial L(w)}{\partial w} = 2X^T X w - 2X^T Y = 0 \quad (1.5)$$

解得:

$$\hat{w} = (X^T X)^{-1} X^T Y \quad (1.6)$$

## 2 最小二乘估计: 几何意义

将  $X$  矩阵从列向量的角度来看, 可以看成是一个  $p$  维的向量空间  $S$ , 为了简便计算, 令  $w^T X = X\beta$ 。可以看成  $Y$  向量到  $S$  的距离最短, 那么将有约束条件:

$$X^T (Y - X\beta) = 0 \quad (2.1)$$

$$X^T Y - X^T X\beta = 0 \quad (2.2)$$

$$\beta = (X^T X)^{-1} X^T Y \quad (2.3)$$

## 3 最小二乘估计: 概率角度

假设一个分布  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ , 那么所有的观测值可看为  $y = w^T x + \varepsilon$ 。因为  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ , 那么  $p(y|x; w) \sim \mathcal{N}(w^T x, \sigma^2)$ 。我们的目的是求  $w$  使,  $y$  出现的概率最大, 在这里可以使用极大似然估计 (MLE) 求解。首先写出  $p(y|x; w)$  的概率密度函数为:

$$p(y|x; w) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y - w^T x)^2}{2\sigma^2}\right) \quad (3.1)$$

对数似然函数为  $\log p(Y|X; w)$ , 使似然函数最大化的过程求解如下:

$$L(w) = \log p(Y|X; w) = \log \prod_{i=1}^N p(y_i|x_i; w) \quad (3.2)$$

$$= \sum_{i=1}^N \log p(y_i|x_i; w) \quad (3.3)$$

$$= \sum_{i=1}^N \left( \log \frac{1}{\sqrt{2\pi}\sigma} + \log \exp\left(-\frac{(y_i - w^T x_i)^2}{2\sigma^2}\right) \right) \quad (3.4)$$

求解目标为  $\hat{w} = \arg \max_w L(w)$ , 因为第一项其中并没有包含  $w$ , 于是可以直接省略, 那么有:

$$\hat{w} = \arg \max_w L(w) \quad (3.5)$$

$$= \arg \max_w \sum_{i=1}^N -\frac{(y_i - w^T x_i)^2}{2\sigma^2}$$

$$= \arg \min_w \sum_{i=1}^N (y_i - w^T x_i)^2 \quad (3.6)$$

那么我可以得到一个结论：最小二乘估计  $\iff$  极大似然估计 (噪声符合高斯分布)。最小二乘估计中隐藏了一个假设条件，那就是噪声符合高斯分布。