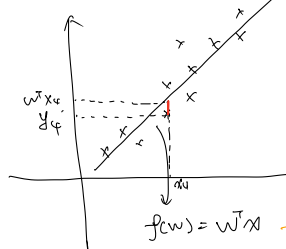


最小二乘法 (矩阵表示, 几何意义)



$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$$

$$x_i \in \mathbb{R}^p, y_i \in \mathbb{R} \quad i=1, 2, \dots, N \quad \text{数据点为向量}$$

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

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

矩阵 $X^T N \times p$

$$\text{最小二乘法} \quad L(w) = \sum_{i=1}^N \|w^T x_i - y_i\|^2$$

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

$$\begin{aligned} &= (w^T x_1 - y_1, w^T x_2 - y_2, \dots, w^T x_N - y_N) \begin{pmatrix} w^T x_1 - y_1 \\ w^T x_2 - y_2 \\ \vdots \\ w^T x_N - y_N \end{pmatrix} \\ &= (w^T x_1, w^T x_2, \dots, w^T x_N) (y_1, y_2, \dots, y_N) \\ &= w^T (x_1, x_2, \dots, x_N) (y_1, y_2, \dots, y_N) \\ &= w^T X^T - Y^T = (w^T X^T - Y^T)^T \\ &= X w - Y \end{aligned}$$

$$= (w^T X^T - Y^T) (X w - Y)$$

$$= w^T X^T X w - w^T X^T Y - Y^T X w + Y^T Y \quad \text{每一项都是标量所以可以交换 } w^T X^T Y = Y^T X w$$

$$= w^T X^T X w - 2 w^T X^T Y + Y^T Y$$

$$\hat{w} = \arg \min L(w)$$

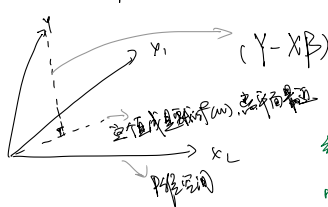
$$\frac{\partial L(w)}{\partial w} = 2 X^T X w - 2 X^T Y = 0 \quad \text{矩阵/向量等式}$$

$$X^T X w = X^T Y$$

$$w = (X^T X)^{-1} X^T Y$$

一般也会把 $(X^T X)^{-1} X^T$ 记作 X^+ 伪逆

$$f(w) = w^T x = x^T \beta \rightarrow \text{把数据点拟合到 } p \text{ 个特征上}$$



$$X = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{Np} \end{pmatrix} \quad (X^T \text{ 的秩为 } p, p \text{ 向量 (特征)})$$

$$\text{结论若 } a \perp b \text{ 则 } a^T \cdot b = 0$$

$$X^T \cdot (Y - X\beta) = 0 \rightarrow \text{向量 } 0 \text{ 是 } p \times 1 \text{ 向量 } \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$X^T Y = X^T X \beta$$

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

N 维高维: 最小二乘法 \Leftrightarrow noise 为 Gaussian 的 ML

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\} \quad x_i \in \mathbb{R}^p, y_i \in \mathbb{R} \quad i=1, 2, \dots, N$$

$$\text{最小二乘法} \quad L(w) = \sum_{i=1}^N \|w^T x_i - y_i\|^2$$

$$\hat{w} = \arg \min L(w)$$

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

$$\text{噪声 } \epsilon \sim N(0, \sigma^2)$$

$$y = f(w) + \epsilon = w^T x + \epsilon$$

$$y | x; w \sim N(w^T x, \sigma^2) \rightarrow P(y | x; w) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{(y - w^T x)^2}{2\sigma^2} \right\}$$

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

$$Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}$$

MLE

独立同分布(i.i.d.)高斯噪声

$$\begin{aligned} \ell(w) &= \log P(y|x;w) = \log \prod_{i=1}^N P(y_i|x_i;w) = \sum_{i=1}^N \log P(y_i|x_i;w) \\ \log \text{ likelihood} &= \sum_{i=1}^N \left(\log \frac{1}{\sqrt{2\pi}} + \log \exp \left\{ -\frac{(y_i - w^T x_i)^2}{2\sigma^2} \right\} \right) \\ &= \sum_{i=1}^N \left(\log \frac{1}{\sqrt{2\pi}} - \frac{1}{2\sigma^2} (y_i - w^T x_i)^2 \right) \end{aligned}$$

$$\hat{w} = \underset{w}{\operatorname{argmax}} \ell(w)$$

$$= \underset{w}{\operatorname{argmax}} \sum_{i=1}^N -\frac{1}{2\sigma^2} (y_i - w^T x_i)^2 \quad \sum \log \frac{1}{\sqrt{2\pi}} \text{ 与 } w \text{ 无关 (不合 } w)$$

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

$$\text{最小二乘估计为 } L(w) = \sum_{i=1}^N \|w^T x_i - y_i\|^2 \quad \hat{w} = \underset{w}{\operatorname{argmin}} L(w) \text{ 两者完全一致}$$

ridge

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

$$w^T = (w^T x)^T w^T y$$

$$\text{过拟合} \rightarrow \begin{cases} \text{① 正则化} \\ \text{② 降维 (特征选择/特征提取)} \\ \text{③ 正则化} \end{cases} \quad \text{PCA}$$

$$\text{正则化效果 } \underset{w}{\operatorname{argmin}} [L(w) + \lambda p(w)]$$

loss penalty

$$L1: \text{LASSO} \quad p(w) = \|w\|$$

$$L2: \text{Ridge} \quad p(w) = \|w\|^2 = w^T w$$

权重衰减

$$J(w) = \sum_{i=1}^N \|w^T x_i - y_i\|^2 + \lambda w^T w$$

拟合误差

$$= (w^T x^T - y^T)(xw - y) + \lambda w^T w$$

$$= (w^T x^T xw - 2w^T x^T y + y^T y) + \lambda w^T w$$

$$= w^T (x^T x + \lambda I) w - 2w^T x^T y + y^T y$$

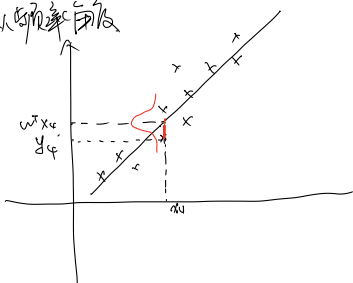
$$\hat{w} = \underset{w}{\operatorname{argmin}} J(w)$$

$$\frac{\partial J(w)}{\partial w} = 2(x^T x + \lambda I) w - 2x^T y = 0$$

$$\boxed{\hat{w} = (x^T x + \lambda I)^{-1} x^T y}$$

从概率角度看 Ridge Regression (岭回归)

之前从频率角度



这是从贝叶斯角度来解.

$w \sim N(0, \sigma_0^2)$ 参数 w 的先验概率

$$p(w|y) = \frac{p(y|w)p(w)}{p(y)}$$

MAP: 最大后验估计

$$\hat{w} = \arg \max_w p(w|y)$$

$$= \arg \max_w p(y|w)p(w)$$

$$= \arg \max_w \log[p(y|w)p(w)]$$

$$= \arg \max_w \log\left(\frac{1}{\sqrt{2\pi}\sigma_0} \cdot \frac{1}{\sqrt{2\pi}\sigma_0} \cdot \exp\left\{-\frac{(y-\tilde{w}x)^2}{2\sigma^2} - \frac{\|w\|^2}{2\sigma_0^2}\right\}\right)$$

$$= \arg \min_w \left(\frac{(y-\tilde{w}x)^2}{2\sigma^2} + \frac{\|w\|^2}{2\sigma_0^2} \right)$$

$$= \arg \min_w \left((y-\tilde{w}x)^2 + \frac{\sigma^2}{\sigma_0^2} \|w\|^2 \right)$$

$$= \arg \min_w \sum_{i=1}^N \underbrace{(y_i - \tilde{w}x_i)^2}_{\text{loss function}} + \underbrace{\left(\frac{\sigma^2}{\sigma_0^2}\right) \|w\|^2}_{\text{penalty}}$$

Regularized LSE \Leftrightarrow MAP (noise Gaussian Distribution, 先验是 GP)

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

$$p(w) = \frac{1}{\sqrt{2\pi}\sigma_0} \exp\left(-\frac{\|w\|^2}{2\sigma_0^2}\right)$$

$$p(y|w) \cdot p(w) = \frac{1}{\sqrt{2\pi}\sigma} \cdot \frac{1}{\sqrt{2\pi}\sigma_0} \cdot \exp\left\{-\frac{(y-\tilde{w}x)^2}{2\sigma^2} - \frac{\|w\|^2}{2\sigma_0^2}\right\}$$

其他机器学习方法

Linear Regression.
 $f(w,b) = w^T x + b$
 $x \in \mathbb{R}$

- ① 线性 \xrightarrow{x}
 - 局部非线性: 特征工程 (多项式回归)
 - 全局非线性: 线性分类 (激活函数是非线性的)
 - 全局非线性: 神经网络.
- ② 全局性 \xrightarrow{x} 线性模型回归. 决策树. (对样本空间的分割)
- ③ 数据降维 \xrightarrow{x} PCA. 流形