

## 线性回归：最小二乘法

$$D = \{(x_i, y_i)\} \quad x_i \in \mathbb{R}^p \text{ } p\text{-dimensional vector} \quad y_i \in \mathbb{R} \quad i=1, 2, \dots, N \quad N \text{ samples}$$

数据集  $D$  可以表示为

表示为

$$X = (x_1 \ x_2 \ \cdots \ x_N)^T = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{pmatrix} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1P} \\ x_{21} & x_{22} & \cdots & x_{2P} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{NP} \end{pmatrix}_{N \times P}$$

$P$ 维向量

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

$$= \begin{pmatrix} W^T x_1 - y_1 \\ W^T x_2 - y_2 \\ \vdots \\ W^T x_N - y_N \end{pmatrix} \begin{matrix} \text{行向量} \\ \text{列向量} \end{matrix}$$

$$= W^T (x_1 \ x_2 \ \dots \ x_N) - (y_1 \ y_2 \ \dots \ y_N) \quad (31)$$

$$= (W^T X^T - Y^T)(XW - Y)$$

$$= W^T X^T X W - W^T X Y - Y X W^T + Y^T Y$$

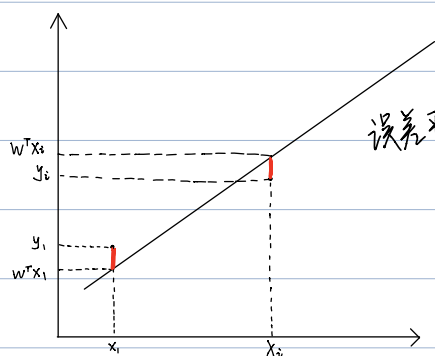
$$= W^T X^T X W - 2W^T X^T Y + Y^T Y$$

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

求导数  $\frac{\partial}{\partial W} L(W) = 2X^T X W - 2X^T Y$

$$\text{令 } 2X^T X W - 2X^T Y = 0 \text{ 得到 } W = (X^T X)^{-1} X^T Y = X^+ Y$$

## 最小二乘法几何意义



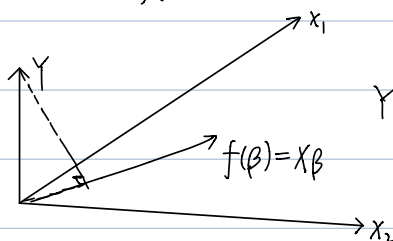
误差平均分散在  $N$  个样本上

由  $W = X^+ Y$  得到  $W \in \mathbb{R}^{p \times 1}$

$$X^T \in \mathbb{R}^{p \times N} \quad Y \in \mathbb{R}^{N \times 1} \quad \beta \in \mathbb{R}^{p \times 1}$$

误差分散在  $p$  个维度上

再假设  $f(w) = w^T x = x\beta$



$Y$  独立于  $X$  构成的  $p$  维空间

最小二乘法: 找到真值  $Y$  在  $p$  维子空间的投影

$x_1, x_2, \dots, x_N$  的线性组合

$$X^T(Y - X\beta) = 0 \quad X^T Y = X^T X \beta \quad \text{所以 } \beta = (X^T X)^{-1} X^T Y = X^+ Y$$

相当于法向量

## 最小二乘法概率视角

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

$p$  维向量

最小二乘法估计损失函数  $L(W) = \sum_{i=1}^N \|W^T x_i - y_i\|_2^2$        $\hat{W} = (X^T X^{-1}) X^T Y$

假设  $y = f(w) + \varepsilon$  (量化模型与真值的误差  $\varepsilon$ , 且要求  $\varepsilon \sim N(0, \sigma^2)$ )

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

$$\text{定义 } L(w) = \log P(Y|X; w) = \log \prod_{i=1}^N P(y_i|x_i; w) = \sum_{i=1}^N \log P(y_i|x_i; w)$$

↑  
log-likelihood

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

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

$$\hat{w} = \arg\max_w L(w) = \arg\min_w \sum_{i=1}^N (y_i - w^T x_i)^2 \longrightarrow \text{极大似然估计} \equiv \text{最小二乘}$$

结论: least square estimation  $\equiv$  maximum likelihood estimation  
当且仅当 noise 服从高斯分布

## 线性回归: 正则化

线性回归引入正则化的原因:

①  $\hat{w} = (X^T X)^{-1} X^T Y$  难以得到解析解. if  $N \gg p$

② 过拟合. (样本数太少)  $\longrightarrow$  { 添加样本  
特征选择  
特征提取 (降维 e.g.)  
正则化 (对参数空间  $w$  进行约束)

正则化  $\arg\min_w \underbrace{L(w)}_{\text{loss}} + \lambda \underbrace{P(w)}_{\text{penalty}}$

L1: LASSO,  $P(w) = \|w\|$

L2: Ridge,  $P(w) = w^T w = \|w\|^2$   
 $\hookrightarrow$  权值衰减

$$\begin{aligned}
L(w) + \lambda P(w) &= \sum_{i=1}^N (y_i - w^T x_i)^2 + \lambda \|w\|^2 \\
&= \begin{pmatrix} w^T x_1 - y_1 & w^T x_2 - y_2 & \dots & w^T x_N - y_N \end{pmatrix} \begin{pmatrix} w^T x_1 - y_1 \\ \vdots \\ w^T x_N - y_N \end{pmatrix} + \lambda \|w\|^2 \\
&= \begin{pmatrix} w^T (x_1 \ x_2 \ \dots \ x_N) - (y_1 \ y_2 \ \dots \ y_N) \end{pmatrix} \begin{pmatrix} w^T x_1 - y_1 \\ \vdots \\ w^T x_N - y_N \end{pmatrix} + \lambda \|w\|^2 \\
&= (w^T X - Y^T) (Xw - Y) + \lambda w^T w \\
&= w^T X^T X w - 2w^T X^T Y + Y^T Y + \lambda w^T w \\
&= w^T (X^T X + \lambda I) w - 2w^T X^T Y + Y^T Y
\end{aligned}$$

$$\frac{\partial}{\partial w} L(w) + \lambda P(w) = \frac{\partial}{\partial w} J(w) = 2(X^T X + \lambda I) w - 2X^T Y$$

$$\text{令导数为0, } \hat{w} = \underbrace{(X^T X + \lambda I)^{-1}}_{\text{半定阵 对角阵}} X^T Y$$

假设  $w \sim p(w) = N(0, \sigma_0^2)$ ,  $y = f(w) + \varepsilon$

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

$$\hat{w} = \arg \max_w p(w|y) = \arg \max_w \frac{P(y|w)P(w)}{\sqrt{2\pi\sigma_0^2} \exp(-\frac{(y-w^T x)^2}{2\sigma_0^2})} \frac{1}{\sqrt{2\pi}\sigma_0} \exp(-\frac{w^2}{2\sigma_0^2})$$

最大后验估计

$$= \arg \max_w \frac{1}{2\pi} \frac{1}{\sigma_0^2} \exp\left(-\frac{(y-w^T x)^2}{2\sigma^2} - \frac{w^2}{2\sigma_0^2}\right)$$

$$= \arg \max_w -\frac{(y-w^T x)^2}{2\sigma^2} - \frac{w^2}{2\sigma_0^2}$$

$$= \arg \min_w \underline{(y-w^T x)^2 + \frac{\sigma^2}{\sigma_0^2} w^2} \quad \text{Ridge Regression!}$$

结论: ① LSE  $\Leftrightarrow$  MLE 极大似然估计  
 ② Regularized LSE  $\Leftrightarrow$  MAP 最大后验估计

} noise is Gaussian