

# P9. Least Square Method and its Geometric Meaning

$$X = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{Np} \end{pmatrix}_{N \times p} \quad Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}_{N \times 1} \quad w = \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_p \end{pmatrix}_{p \times 1} \quad f(w) = w^T X \Rightarrow \text{each data}$$

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

$$L(w) = (w^T x_1 - y_1, \dots, w^T x_N - y_N) \begin{pmatrix} w^T x_1 - y_1 \\ \vdots \\ w^T x_N - y_N \end{pmatrix} = [w^T (x_1, x_2, \dots, x_N) - (y_1, y_2, \dots, y_N)] \Delta^T$$

$\therefore L(w) = (w^T X^T - Y^T) (Xw - Y) = w^T X^T Xw - w^T X^T Y - Y^T Xw + Y^T Y$   $w^T X^T Y = (Y^T Xw)^T$  are constants, so they are equal

$$\therefore L(w) = w^T X^T Xw - 2w^T X^T Y + Y^T Y \quad \frac{d}{dw} L(w) = 2X^T Xw - 2X^T Y = 0$$

$2X^T Xw = 2X^T Y \quad w = (X^T X)^{-1} X^T Y$  ( $\because X^T = X^{-T}, X[(X^T X)^{-1} X^T]X = X \therefore (X^T X)^{-1} X^T$  is called pseudoinverse or generalized inverse (伪逆或广义逆))

Geometric Meaning: ① Splitting the error to each data. LS minimizes the sum of the distance between observed values and theoretical values

② Splitting the error to each dimensions of attributes  $f(w) = w^T x = Xw$

All dimensions form a vector space:  $f(w) = \sum_{i=1}^p w_i x_i$

$Y$  doesn't lie in this vector space unless all data fit. LS is to find the closest  $f(w)$  to  $Y$  in the vector space

The  $f(w)$  that minimizes  $Y - f(w)$   $Y - f(w)$  is perpendicular to all  $X \therefore X^T (Y - Xw) = 0 \therefore X^T Y = X^T Xw \quad w = (X^T X)^{-1} X^T Y$  ← analytic solution

numerical solution  $\rightarrow GD$

## P10. Probabilistic Perspective LS

Least Square Estimation  $\longleftrightarrow$  Maximum Likelihood Estimation

When the noise obeys Gaussian distribution  $LSE \longleftrightarrow MLE \quad LSE: \hat{w} = \arg \min L(w) \quad L(w) = \sum_{i=1}^N \|w^T x_i - y_i\|^2$

$$MLE: y = w^T x + \epsilon \quad \epsilon \sim N(0, \sigma^2) \quad y|x, w \sim N(w^T x, \sigma^2) \quad P(y|x, w) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(y - w^T x)^2}{2\sigma^2}\right\}$$

$$L(w) = \log P(y|x, w) = \sum_{i=1}^N \left[ \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} (y_i - w^T x_i)^2 \right] \quad \hat{w} = \arg \max_w L(w)$$

$$\hat{w} = \arg \min_w (y - w^T x)^2 = LSE$$

## P11. Regularization - Ridge Regression - Frequentists (频率学派)

$X_{N \times p}$ ,  $N$  samples,  $p$  attributes Generally  $N \gg p$ , otherwise it may be over-fitted

The  $(X^T X)^{-1}$  in  $\hat{w} = (X^T X)^{-1} X^T Y$  may not exist ① Add data ② Feature Selection/Extraction

③ Regularization:  $\hat{w} = \arg \min_w [L(w) + \lambda P(w)] \rightarrow$  Penalty (惩罚)

$L_1$ : Lasso,  $P(w) = \|w\|_1$   $L_2$ : Ridge, 岭回归,  $P(w) = \|w\|^2$

$$J(w) = (w^T X^T - Y^T)(Xw - Y) - \lambda w^T w = w^T (X^T X + \lambda I)w - 2w^T X^T Y + Y^T Y$$

$$\hat{w} = \arg \min_w J(w) \quad \frac{\partial}{\partial w} J(w) = (X^T X + \lambda I)w - 2X^T Y = 0 \quad \hat{w} = (X^T X + \lambda I)^{-1} X^T Y$$

## P12. Regularization - Ridge Regression - Bayesians (贝叶斯学派)

Regularization: As above ①

Bayesian Perspective:  $y = w^T x + \epsilon$  Let  $\epsilon \sim N(0, \sigma^2)$ ,  $y \sim N(w^T x, \sigma^2)$ ,  $w \sim N(0, \sigma_0^2)$

$$P(w|y) = \frac{P(y|w)P(w)}{P(y)} \quad \text{Maximum a Posteriori Estimation, MAP: } \hat{w} = \arg \max_w [P(y|w)P(w)]$$

$$P(y|w) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(y_i - w^T x_i)^2}{2\sigma^2}\right\} \quad P(w) = \frac{1}{\sqrt{2\pi}\sigma_0} \exp\left\{-\frac{\|w\|^2}{2\sigma_0^2}\right\}$$

$$\hat{w} = \arg \max_w \left\{ \log [P(y|w)P(w)] \right\} = \arg \max_w \left\{ \log \frac{1}{\sqrt{2\pi}\sigma} \frac{1}{\sqrt{2\pi}\sigma_0} - \left[ \sum_{i=1}^N \frac{(y_i - w^T x_i)^2}{2\sigma^2} + \frac{\|w\|^2}{2\sigma_0^2} \right] \right\}$$

$$= \arg \min_w \left[ \sum_{i=1}^N (y_i - w^T x_i)^2 + \frac{\sigma_0^2}{\sigma^2} \|w\|^2 \right] \iff ①$$

Loss  $\downarrow$   $\lambda$  penalty

## Supplement to P11, P12

$LSE \longleftrightarrow MLE$  (When the noise obeys Gaussian distribution)

Regularized (Ridge)  $LSE \longleftrightarrow MAP$  (When the noise and the prior obey the Gaussian distribution)