This lecture:  estimation of $\beta_{(p+1) \times 1} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}$

Goal: to find $\hat{\beta} = \begin{pmatrix} \hat{\beta}_0 \\ \vdots \\ \hat{\beta}_p \end{pmatrix}_{(p+1) \times 1}$ such that it minimizes the RSS function

$$\underset{\underset{J(\beta)}{\downarrow}}{RSS} = \sum_{i=1}^{n} (y_i - x_i^T \beta)^2.$$

$J(\beta)$

From calculus, the minimizer of $J(\beta)$ is achieved when

$$\nabla_\beta J(\beta) = 0$$

that is,  $\nabla_\beta J(\beta) \Big|_{\beta = \hat{\beta}} = 0.$  ✱ ($\leftarrow$ (p+1) equations)

Method 1:

$$\nabla_\beta J(\beta) = \nabla_\beta \sum_{i=1}^{n} (y_i - x_i^T \beta)^2. \qquad x_i : (p+1) \times 1$$

(p+1)×1 dim

$$= \sum_{i=1}^{n} \nabla_\beta (y_i - x_i^T \beta)^2 \qquad \overset{\shortparallel}{\begin{pmatrix} 1 \\ x_{i,1} \\ \vdots \\ x_{i,p} \end{pmatrix}}$$

$$= \sum_{i=1}^{n} 2 \cdot (y_i - x_i^T \beta) \cdot \underbrace{\nabla_\beta (y_i - x_i^T \beta)}_{\overset{\shortparallel}{-x_i}}$$

$$= -2 \sum_{i=1}^{n} (y_i - x_i^T \beta) x_i. \qquad \leftarrow (p+1) \times 1$$

We know from the above discussion that $\hat{\beta}$ satisfies

$$\cancel{-2} \sum_{i=1}^{n} (y_i - x_i^T \hat{\beta}) x_i = 0$$

Use the matrix notation, we have

$$\sum_{i=1}^{n} (y_i - x_i^T \hat{\beta}) x_i = \underbrace{[x_1, x_2, \dots, x_n]}_{} \begin{bmatrix} y_1 - x_1^T \hat{\beta} \\ y_2 - x_2^T \hat{\beta} \\ \vdots \\ y_n - x_n^T \hat{\beta} \end{bmatrix}_{n \times 1}$$

$(p+1) \times 1 \qquad \overset{(p+1)}{\underset{n}{\times}}$

$$X^T \quad \{ \text{Recall that}$$

$$X_{n\times(p+1)} = \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1p} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n,1} & \cdots & x_{np} \end{pmatrix}$$

$$= X^T \cdot \left[ \underbrace{\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}}_{Y_{n\times 1}} - \underbrace{\begin{pmatrix} x_1^T \hat\beta \\ \vdots \\ x_n^T \hat\beta \end{pmatrix}}_{\parallel} \right]$$

$$\begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \hat\beta = \underset{n\times(p+1)}{X}\, \underset{(p+1)\times 1}{\hat\beta}$$

$$= \underbrace{X^T (Y - X\hat\beta) = O_{(p+1)\times 1}}_{\text{normal equations.}}$$

So $\quad X^T y - X^T X \hat\beta = 0$

$$\Rightarrow \underset{(p+1)\times(p+1)}{(X^T X)}\; \underset{(p+1)\times 1}{\hat\beta} = \underset{(p+1)\times 1}{X^T y}$$

$$\Rightarrow \quad \hat\beta = (X^T X)^{-1} X^T y \quad \text{if } \underline{(X^T X)^{-1} \text{ exists}}$$

OLS estimator

$\bullet$ some $\quad$ rank$(X^T X) = (p+1)$.

E.g. ~~two~~ columns of $X_{n\times(p+1)}$

"colinearity" issue. $\longrightarrow$ are the same

$\Rightarrow$ rank$(X^T X) < (p+1)$

$\bullet$ or a column of $X$ is a linear combination of other columns.

E.g: $\quad \underline{y_i = x_i} \rightarrow \underset{\sim}{\beta_0 = 0.\ \beta_1 = 1}$

$$y_i = \beta_0 + \beta_1 \underset{\text{1dim}}{x_i} + \varepsilon_i \quad i = 1, \ldots, n$$

$$\underset{n\times 2}{X} = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$$

$\underset{\text{not linearly dependent.}}{\uparrow \quad \uparrow} \quad \text{if } \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \neq \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$

high-dimension data $\left\{ \begin{array}{l} \bullet \text{ If } p > n \\ \quad \text{rank}(X^T X) \le n < p \end{array} \right.$

# Method 2

We first rewrite $J(\beta)$ using matrix notation.

$$J(\beta) = \sum_{i=1}^{n} (y_i - x_i^T \beta)^2$$

$$= \underbrace{(y_1 - x_1^T\beta, \ \cdots\cdots, \ y_n - x_n^T\beta)}_{(Y - X\beta)^T} \quad \begin{pmatrix} y_1 - x_1^T\beta \\ \vdots \\ \vdots \\ y_n - x_n^T\beta \end{pmatrix}_{n\times 1}$$

$$1\times n$$

↑ as shown in Method 1

$$\underset{n\times 1}{Y} - \underset{n\times(p+1)}{X} \ \underset{(p+1)\times 1}{\beta}$$

$$= (Y - X\beta)^T (Y - X\beta)$$

$Tr(AA^T) = Tr(A^TA)$

$$\nabla_\beta J(\beta) = \nabla_\beta \ Tr\{(Y - X\beta)(Y - X\beta)^T\}$$

$$= \nabla_\beta \left[ Tr\{(Y - X\beta)^T (Y - X\beta)\} \right]$$

"B" in ②.

$$\beta^T \underbrace{X^T X} \beta$$

$$\underbrace{Y^T X \beta}_{\|}$$

$$= \nabla_\beta \left[ Tr\left( Y^TY - \underbrace{(X\beta)^T Y}_{\text{scalar } 1\times 1} - \underbrace{Y^T X \beta}_{\text{"A" in ①}} + \underbrace{(X\beta)^T (X\beta)}_{} \right) \right]$$

$$= 0 - 2\cdot (Y^T X)^T$$

$$+ 2(X^T X)\beta$$

$$= -2 X^T Y + 2(X^T X)\beta$$

Therefore, $\hat{\beta}$ satisfies

$$-\cancel{2} X^T Y + \cancel{2}(X^T X)\hat{\beta} = 0$$

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

if $(X^T X)^{-1}$ exists.

---

Recall that from linear alg file
(page 23-26 § 4.3)

① $\nabla_\beta Tr(A\beta) = A^T$ for any A matrix

② $\nabla_\beta Tr(\beta^T B \beta) = 2B\beta$
for symmetric matrix B.

These results hold generally for
$\beta$ being vectors or matrix
↑
(HW1. Q5)

---

Hints: • Q5. HW. multivariate response $Y_i \in R^m$. [OLS. m=1]

$$\underset{1\times m}{\underset{\bigcirc}{Y_i^T}} = \underset{1\times d}{\underset{\bigcirc}{X_i^T}} \underset{d\times m}{\underline{B}} + error \qquad (d = p+1)$$

(a) $RSS(B) = \sum_{i=1}^{n} \underbrace{(Y_i^T - X_i^T B)}_{} \underbrace{(Y_i^T - X_i^T B)^T}_{}$

Hint: For $A = \begin{bmatrix} -a_1- \\ \vdots \\ -a_n- \end{bmatrix}$ $\quad Tr(AA^T) = \sum_{i=1}^{n} \underbrace{a_i a_i^T}_{}$.

$$RSS(B)$$
$$\|$$
$$Tr(AA^T)$$
with $a_i = Y_i^T - X_i^T B$

(b). Find $\hat{B}$ s.t $\nabla_B RSS(B) = 0$

$$\nabla_B Tr\left[(Y - XB)(Y - XB)^T\right] = 0$$

$$A = \begin{bmatrix} Y_1^T - X_1^T B \\ \vdots \\ Y_n^T - X_n^T B \end{bmatrix}$$

$$= \begin{bmatrix} Y_1^T \\ \vdots \\ Y_n^T \end{bmatrix} - \begin{bmatrix} X_1^T \\ \vdots \\ X_n^T \end{bmatrix} B$$

$$= Y_{n \times m} - X_{n \times d} B_{d \times m}$$

---

Example  Simple linear Regression  $(p = 1)$   note that now $x_1 \cdots x_n \in R$
(Q4 Hw1)

$$y_{n \times 1} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \qquad X_{n \times 2} = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \qquad \beta_{2 \times 1} = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}$$

To find OLS $\hat{\beta}$.

Following Method 1, we can have the normal equations.

$$\sum_{i=1}^{n} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)\begin{pmatrix} 1 \\ x_i \end{pmatrix} = 0$$

$$\Rightarrow \begin{cases} \sum_{i=1}^{n} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0 \cdots \text{ⓐ} \quad \leftarrow \quad \dfrac{\partial J(\beta)}{\partial \beta_0} = 0 \\[4mm] \sum_{i=1}^{n} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i = 0 \cdots \text{ⓑ} \quad \leftarrow \quad \dfrac{\partial J(\beta)}{\partial \beta_1} = 0 \end{cases}$$

(practice)

Eq ⓐ $\Rightarrow$ $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x}$

(practice)

ⓐ + ⓑ $\Rightarrow$ $\hat{\beta}_1 = \dfrac{\sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{Y})}{\sum_{i=1}^{n} (x_i - \bar{x})^2}$

where $\bar{Y} = \dfrac{1}{n}\sum_{i=1}^{n} y_i$ , $\bar{x} = \dfrac{1}{n}\sum_{i=1}^{n} x_i$ .

Note that $\hat{\beta}_1 = \dfrac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{Y})/n-1 \quad \leftarrow \text{Sample covariance of x and Y}}{\sum_{i=1}^{n}(x_i - \bar{x})^2 /n-1 \quad \leftarrow \text{Sample variance of X}}$

(practice)

$= \text{Sample correlation}(X, Y) \cdot \dfrac{\text{Sample Standard deviation (Y)}}{\text{Sample Standard dev (X)}}$