



Multiple Linear Regression

- Target (Dependent variable) : Y
- Feature (Independent variable) : X
- 목표 : $Y = F(X)$. Find F

• linear relationship find (Y, X)

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \underbrace{\epsilon}_{\text{noise}}$$

$$\xrightarrow{\text{Find}} \hat{y} = \underbrace{\hat{\beta}_0} + \underbrace{\hat{\beta}_1} x_1 + \underbrace{\hat{\beta}_2} x_2 + \dots + \underbrace{\hat{\beta}_n} x_n$$

- OLS (Ordinary least square)

⇒ Minimize the squared difference between Y, \hat{Y} .

- OLS solve.

$X: n \times (d+1)$ matrix \Rightarrow data.

$y: n \times 1$ vector \Rightarrow actual value

$\hat{\beta}: (d+1) \times 1$ vector \Rightarrow parameter

$$\min \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \min \frac{1}{2} (y - X\hat{\beta})^T (y - X\hat{\beta})$$

$$\Rightarrow (y - X\hat{\beta})' (y - X\hat{\beta}) = y'y - \hat{\beta}' X' y - y' X \hat{\beta} + \hat{\beta}' X' X \hat{\beta}$$

$$= y'y - 2\hat{\beta}' X' y + \hat{\beta}' X' X \hat{\beta}$$

$$\Rightarrow \frac{\partial E(X)}{\partial \hat{\beta}} = -2X'y + X'X\hat{\beta} + \hat{\beta}' X' X$$
$$= -2X'y + 2X'X\hat{\beta}$$

$$\Rightarrow -2X'y + 2X'X\hat{\beta} = 0$$

$$\Rightarrow 2X'X\hat{\beta} = 2X'y$$

$$\Rightarrow (X'X)\hat{\beta} = X'y$$

$$\Rightarrow \hat{\beta} = (X'X)^{-1} X'y$$

↳ Unique, explicit solution

★ 4가지 조건 (OLS의 핵심가정)

① ε 는 normal Distribution

② Homoskedasticity.

③ linear relationship is correct.

④ independent.

check 1

→ QQ - plot. Residual

check 2.

→ Residual plot.

· 평가 방법!

R^2

Sum-of-Squares Decomposition

$$\sum_{j=1}^n (y_j - \bar{y})^2 = \sum_{j=1}^n (\hat{y}_j - \bar{y})^2 + \sum_{j=1}^n \hat{\varepsilon}_j^2$$

(total sum of squares about mean) = (regression sum of squares) + (residual (error) sum of squares)

SST SSR SSE

Total variance Explained variance

$$R^2 = 1 - \frac{\text{SSE}}{\text{SST}}$$

↑ SST

↑ SSE

⇒ 1 - 전체 데이터의 변동성과 회귀식에서 설명할 수 있는 변동성의 비율

Adjusted R^2

$$R^2_{Adj} = 1 - \left[\frac{n-1}{n - (p+1)} \right] \frac{SSE}{SST} \leq R^2$$

\downarrow
변수개수

\Rightarrow 변수가 늘어나면

무조건 커지는 R^2 의 한계를 극복하기
위해 (손을 많은 변수 제외 시켜)

* 다중 공선성.

$$C = 1 - A - B$$

* p-value.

$$H_0 : \alpha = 0 \quad (\text{비무가설})$$

$$H_1 : \alpha \neq 0 \quad (\text{대립가설})$$

\leftarrow 비무가설의 불부

