

Cov matrix는 안 안다!

## \* Multicollinearity in multiple regression data.

- Correlation matrix

$$X = (\underline{1}, \underline{x}_1, \underline{x}_2, \dots, \underline{x}_k) \quad n \times p \text{ matrix}$$

$$\downarrow \text{centering \& scaling} \quad x_{ij} \Rightarrow x_{ij}^* = \frac{x_{ij} - \bar{x}_j}{\sqrt{\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2}}$$

$$X = \begin{pmatrix} \frac{1}{\sqrt{n}} & x_{11}^* & x_{12}^* & \dots & x_{1k}^* \\ \frac{1}{\sqrt{n}} & x_{21}^* & x_{22}^* & \dots & x_{2k}^* \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{1}{\sqrt{n}} & x_{n1}^* & x_{n2}^* & \dots & x_{nk}^* \end{pmatrix} = \left( \frac{1}{\sqrt{n}} \underline{1}, \underline{x}_1^*, \underline{x}_2^*, \dots, \underline{x}_k^* \right)$$

↑  
y는 centering 하지 않았기 때문에 상수항의 변인이 됨

$$\begin{aligned} y &= \beta_0 + \beta_1 x_1^* + \beta_2 x_2^* + \dots + \beta_k x_k^* + \varepsilon \\ \bar{y} &\neq 0 \quad \text{But} \quad \bar{x}_1^* = \bar{x}_2^* = \dots = \bar{x}_k^* = 0 \\ \therefore \hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x}_1^* - \dots - \hat{\beta}_k \bar{x}_k^* \\ &= \bar{y} \end{aligned}$$

symmetric

$$X^* X^* = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & r_{12} & \dots & r_{k1} \\ \vdots & r_{12} & 1 & \dots & r_{k2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & r_{k1} & r_{k2} & \dots & 1 \end{pmatrix}$$

$X^* X^*$   
:  $x$ 들의 correlation matrix

$$\begin{aligned} (\text{eg}) (X^* X^*)_{23} &= \underline{x}_1^{*'} \underline{x}_2^* = \sum_{i=1}^n \frac{x_{i1} - \bar{x}_1}{\sqrt{\sum_{i=1}^n (x_{i1} - \bar{x}_1)^2}} \cdot \frac{x_{i2} - \bar{x}_2}{\sqrt{\sum_{i=1}^n (x_{i2} - \bar{x}_2)^2}} \\ &= \widehat{\text{Corr}}(x_1, x_2) \\ &= r_{12} \end{aligned}$$

$$\therefore (X^* X^*)_{ij} = \underline{x}_{i-1}^{*'} \underline{x}_{j-1}^* = r_{i-1, j-1}$$

## - Problem of Multicollinearity

:  $x$ 들의 correlation이 높으면, coefficient 등의 variance에 inflation이 생겨난다.

$$\text{if } \begin{bmatrix} 1 & r_{12} \\ r_{12} & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0.99915 \\ 0.99915 & 1 \end{bmatrix} \xrightarrow{\text{Inverse}} \begin{bmatrix} 63.94 & -63.44 \\ -63.44 & 63.94 \end{bmatrix} \quad \hat{\beta} \sim (\underline{1}, (X^* X^*)^{-1} \sigma^2) \quad \text{이니까!}$$

$(X^* X^* = I)$  일때 대비 약 64배!!

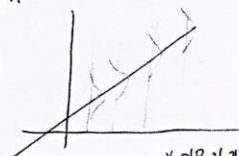
$\Rightarrow$  따라서 리커계수 자체에 관성이 있다면, (y와  $x$ 사이의 관계). 치명적인 문제

But.  $y$ 에 관성이 있다면, 크게 문제가 되지 않음.

$\Rightarrow$  이런 Multicollinearity의 정도를 나타내는 것이 VIF (Variance Inflation Factor)

\* 가설가가  
불관용해로  
꼭지점의 수는  
변화가 거의 X

(그러니까  
리커계수의 추정이  
안타깝지 않음.)



# - VIF $\nearrow x_j$

: 한 변수가 나머지 다른 변수들과 선형 관계를 가지는 것 때문에  $\hat{\beta}_j$ 의 variance가 얼마나 증가한 것인가를 보여줌.

$\leftrightarrow$  j-th 변수가 나머지 변수들과 linearly independent 또는 orthogonal 한때

$\text{Var}(\hat{\beta}_j) = 1$  이라면, VIF<sub>j</sub>는 그것들의 선형 관계에 의해 몇배나 증가하는가.

$$\text{VIF}_j = \frac{1}{1-R_j^2} \rightarrow x_j^* = \beta_1 x_1^* + \dots + \beta_{j-1} x_{j-1}^* + \beta_{j+1} x_{j+1}^* + \dots + \beta_k x_k^* + \varepsilon$$

에서의  $R^2$  (모든  $x^*$ 들은 centering 했기 때문에 상수항이 없다)

## \* VIF 구하기 - 1

i)  $X^* = (x_1^*, x_2^*, \dots, x_k^*) : n \times k \text{ matrix}$

$= (x_1^* \quad x_1^*) \rightarrow x_1$ 은 제외한 나머지  $k-1$ 개의 독립변수

(cf)  $C = \begin{pmatrix} A & B \\ B' & D \end{pmatrix}$  symmetric 이냐

$$X^{*'} X^* = \begin{pmatrix} x_1^{*'} \\ x_1^{*'} \end{pmatrix} (x_1^*, x_1^*) = \begin{pmatrix} x_1^{*'} x_1^* & x_1^{*'} x_1^* \\ x_1^{*'} x_1^* & x_1^{*'} x_1^* \end{pmatrix}$$

$$C^{-1} = \frac{1}{AD-BC} \begin{pmatrix} D & -B \\ -C & A \end{pmatrix}$$

$$\therefore C^{-1}_{11} = (A - BD'B)^{-1}$$

ii)  $(X^{*'} X^*)^{-1} = (x_1^{*'} x_1^* - x_1^{*'} x_1^* (x_1^{*'} x_1^*)^{-1} x_1^{*'} x_1^*)^{-1} \quad \therefore H = X^* (X^{*'} X^*)^{-1} X^{*'} \quad \leftarrow$

$$= (x_1^{*'} (I-H) x_1^*)^{-1}$$

$$= \left( \frac{x_1^{*'} (I-H) x_1^*}{x_1^{*'} x_1^*} \right)^{-1} \quad \because x_1^{*'} x_1^* = 1 \text{ (scaled! 제곱합 1)}$$

$$= \left( \frac{\text{SSE}}{\text{SST}} \right)^{-1}$$

$$= \frac{1}{1-R_j^2}$$

$x_1^* = x_1 x_1 + \varepsilon$  이 모델의 SSE

$\therefore \text{SST} = \sum (x_1^* - \bar{x}_1^*)^2 = \sum x_1^{*2} = x_1^{*'} x_1^* \quad \begin{matrix} \uparrow \\ \text{centered} \end{matrix} = 1 \therefore \text{scaled}$

iii)  $\text{Var}(\hat{\beta}_1) = (X'X)^{-1}_{11} \sigma^2$   
 $= (X^{*'} X^*)^{-1}_{11} \sigma^2$   
 $= \text{VIF}_1 \cdot \sigma^2$

$\therefore$  VIF의 뜻이  $\text{Var}(\hat{\beta}_j) = \sigma^2$  일때 대비 몇배 증가하는가!

$\text{S.E.}(\hat{\beta}_1)^2 / \text{MSE} = \text{VIF}_1 ?$

## <Remark>

- 모든 변수에 대해서 VIF를 계산

만약 변수  $x_j^*$ 가 다른 변수들에 의해 설명이 많이된다면,  $\text{Var}(\hat{\beta}_j)$ 의 값이 커짐.

- VIF > 10 이면 Multicollinearity의 문제가 존재한다고 한다.