

Exam 1: mean: 40.2/50 sd \approx 9.6

median 43 . 1st quartile (25%): 37 . 3rd quartile (75%): 46.5

Diagnosis:

§ 1. Collinearity Issue

(§ 3.3.3.6 Intro to Stats learning)

• Assumption: $(X^T X)^{-1}$ exists

• Collinearity refers to the situation when 2 or more predictors are closely correlated with each other.

• The presence of collinearity can pose problems in Regression Analysis as it can be difficult to separate out the individual effects of the collinear variables on Y .

large s.e. ($\hat{\beta}_j$)
↑

In particular, collinearity reduces the accuracy of the estimates of $\hat{\beta}_j$ (if x_j is highly correlated with some other predictors).

To see this, consider a simple linear reg.

$$\widehat{\text{Var}} \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} = \hat{\sigma}^2 (X^T X)^{-1}$$

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

Suppose two columns are highly correlated.

$$\det(X^T X)$$

$$= n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2$$

"relatively small" (compared with n)
↓

Lecture 12:

$$(X^T X)^{-1} = \frac{1}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2} \begin{bmatrix} \sum_{i=1}^n x_i^2 & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & n \end{bmatrix}_{2 \times 2}$$

$$\Rightarrow \text{s.e.}(\hat{\beta}_0) = \sqrt{\widehat{\text{Var}}(\hat{\beta}_0)} = \sqrt{\hat{\sigma}^2 \cdot \frac{\sum_{i=1}^n x_i^2}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2}} \quad \text{large s.e.}$$

$$\text{s.e.}(\hat{\beta}_1) = \dots = \sqrt{\hat{\sigma}^2 \cdot \frac{n}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2}}$$

- How to detect ?

① use the correlation matrix of predictors (large correlation close to 1)
 \Rightarrow collinearity

But this method cannot detect multi-collinearity

(one predictor is highly correlated with a linear combination of multiple predictors)

E.g. $X_1 \approx X_2 + X_3 + X_4$

$$\left. \begin{array}{l} \text{Corr}(X_1, X_2) \\ \text{Corr}(X_1, X_3) \\ \text{Corr}(X_1, X_4) \end{array} \right\} \text{ may not large.}$$

② A recommended method.

- We usually use the R^2 from the regression of

$$X_j \sim \text{all other predictors (including intercept)}$$

denote them by X_{-j}

for $j=1, \dots, p$.

- We denote $R^2_{X_j | X_{-j}}$ to be the R^2 from the above regression

- In practice, if $R^2_{X_j | X_{-j}} > 0.9$ or 0.8 , this indicates a problematic amount of collinearity.

* Equivalently, people also use the "Variance Inflation Factor" (VIF)

For $j=1, \dots, p$.

$$VIF_j = \frac{\widehat{\text{Var}}(\hat{\beta}_j) \text{ from the full regression } Y \sim X_1 + X_2 + \dots + X_p}{\widehat{\text{Var}}(\hat{\beta}_j) \text{ from the simple regression } Y \sim X_j}$$

$$= \frac{1}{1 - R_{X_j|X_{-j}}^2}$$

• $VIF_j > 10$ or 5

§2. ★ Gauss - Markov Theorem.

Under the linear regression model assumptions, the OLS estimator $\hat{\beta}_{(p+1) \times 1}$
 (n, p)
 is the Best Linear Unbiased Estimator of β (true)
 (BLUE)

That is, if $\tilde{\beta}$ is any other unbiased linear estimator of β ,

unbiased: $E(\tilde{\beta} | x) = \beta$
 linear estimator means that we can write
 $\tilde{\beta}$ as a linear combination of Y , i.e.

$$\tilde{\beta} = A Y$$

$$\begin{matrix} (p+1) \times (p+1) & & (p+1) \times n & & n \times 1 \end{matrix}$$

then

$$\text{Var}(\hat{\beta}_{OLS} | x) \leq \text{Var}(\tilde{\beta} | x)$$

$$\begin{matrix} & \uparrow & \\ & (p+1) \times (p+1) & \end{matrix}$$

in the sense that for any $a \in \mathbb{R}^{p+1}$

$$\underbrace{a^T \text{Var}(\hat{\beta} | x) a}_{\|1 \times 1\|} \leq \underbrace{a^T \text{Var}(\tilde{\beta} | x) a}_{\|1 \times 1\|}$$

$$\text{Var}(a^T \hat{\beta} | x) \leq \text{Var}(a^T \tilde{\beta} | x)$$

Proof: Let $\tilde{\beta} = A\gamma$. Then

$$\begin{aligned}\tilde{\beta} - \hat{\beta}_{OLS} &= A\gamma - (X^T X)^{-1} X^T Y \\ &= \underbrace{\left[A - (X^T X)^{-1} X^T \right]}_{C_{(p+1) \times n}} \gamma\end{aligned}$$

$$\text{Then } \tilde{\beta} = \hat{\beta} + \underbrace{C}_{(p+1) \times n} \underbrace{Y}_{n \times 1} = \hat{\beta} + C(X\beta + e)$$

This implies

$$\begin{aligned}E(\tilde{\beta} | x) &= E\left[\hat{\beta}_{OLS} | x\right] + E\left[C(X\beta + e) | x\right] \\ &= \beta + Cx\beta.\end{aligned}$$

Because $\tilde{\beta}$ is unbiased, $E(\tilde{\beta} | x) = \beta$.

$$\Rightarrow Cx\beta = 0 \quad \text{for any } \beta.$$

$$\Rightarrow Cx = 0$$

This gives

$$\text{Var}(\tilde{\beta} | x) = \text{Var}(\tilde{\beta} - \beta | x)$$

$$= \text{Var}\left[\underbrace{C + (X^T X)^{-1} X^T}_{\text{red arrow}} e | x\right]$$

$$\left[\begin{array}{l} \text{This is because } \tilde{\beta} = \hat{\beta} + \underbrace{Cx\beta}_0 + ce = \hat{\beta} + ce \\ \text{and } \tilde{\beta} - \beta = ce + \underbrace{\hat{\beta} - \beta} \\ = ce + \underbrace{(X^T X)^{-1} X^T e} = \left[C + (X^T X)^{-1} X^T \right] e \end{array} \right]$$

For matrix B only depending
on x ,

$$\text{Var}(Be | x) = B \text{Var}(e | x) B^T$$

$$= \underbrace{\left[C + (X^T X)^{-1} X^T \right]}_{6^2 I_n} \underbrace{\text{Var}(e | x)}_{6^2 I_n} \left(C + (X^T X)^{-1} X^T \right)^T$$

$$\begin{aligned}
&= \sigma^2 \left[C + (X^T X)^{-1} X^T \right] \left[C^T + X (X^T X)^{-1} \right] \\
&= \sigma^2 \left[CC^T + \underbrace{C}_{\text{0}} \underbrace{X (X^T X)^{-1}}_{\text{0}} + (X^T X)^{-1} \underbrace{X^T C^T}_{\text{0}} \right. \\
&\quad \left. + (X^T X)^{-1} \underbrace{X^T X}_{\text{0}} (X^T X)^{-1} \right] \\
&= \sigma^2 \left[CC^T + (X^T X)^{-1} \right].
\end{aligned}$$

$$\Rightarrow \text{Var}(\tilde{\beta} | X) = \sigma^2 CC^T + \underbrace{\sigma^2 (X^T X)^{-1}}_{\text{Var}(\hat{\beta} | X)}$$

Since $\underbrace{CC^T}_{(p+1) \times (p+1)} \succeq 0$ (that is, for any $a \in \mathbb{R}^{p+1}$, $a C C^T a^T \geq 0$),

we have

$$\text{Var}(\tilde{\beta} | X) \succeq \underbrace{\text{Var}(\hat{\beta} | X)}_{\text{OLS}}. \quad \#$$

Next lecture: Regression with categorical predictors. (§ 5. Applied linear Reg).
 (factors)