

X. Consequences of a Problematic Regression Space

Usual situation

DATA: $(Y_i, Z_i^T)^T$, $i=1, \dots, n$, $Z_i = (z_{i1}, \dots, z_{ip})^T \in \mathcal{Z}$

$$Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}, \quad X = \begin{pmatrix} X_1^T \\ \vdots \\ X_n^T \end{pmatrix} = (1_n, X^1, \dots, X^{k-1})$$

$$X_i = t_x(Z_i), \quad i=1, \dots, n$$

Full-rank linear model with intercept assumed

$$Y|Z \sim (X\beta, \sigma^2 I_n), \quad \text{rank}(X_{n \times k}) = k < n$$

From now

- $E(Y|Z) \in \mathcal{M}(X)$ (\equiv model is correct)
- $X = (1, \dots)$ (intercept in model)
- $\text{rank}(X) = k$ (full-rank)

10.1 Multicollinearity

2

Principal assumption of linear model $E(Y|Z) \in \text{span}(X)$

Many covariates available \Rightarrow Choose X with many columns (= all available covariates)?

10.1.1 Singular value decomposition of a model matrix

3

$$X_{n \times k} = \sum_{j=0}^{k-1} d_j u_j v_j^T = U D V^T$$

$$U_{n \times k} = (u_0, \dots, u_{k-1})$$

- first k orthonormal eigenvectors of the $n \times n$ matrix $X X^T$

$$V_{k \times k} = (v_0, \dots, v_{k-1})$$

- (all) orthonormal eigenvectors of the $k \times k$ (invertible) matrix $X^T X$

$$D = \text{diag}(d_0, \dots, d_{k-1}), \quad d_j = \sqrt{\lambda_j}, \quad j=0, \dots, k-1$$
$$d_0^2 \geq d_1^2 \geq \dots \geq d_{k-1}^2 > 0$$

• the first k eigenvalues of matrix $X X^T$

• (all) eigenvalues of matrix $X^T X$

$$d_0 \geq d_1 \geq \dots \geq d_{k-1} > 0 : \text{ singular values of matrix } X$$

$\sqrt{\lambda_0} \quad \sqrt{\lambda_1} \quad \dots \quad \sqrt{\lambda_{k-1}}$

2

$$\bullet X^T X = \sum_{j=0}^{k-1} \lambda_j v_j v_j^T = V \Lambda V^T, \quad \Lambda = \text{diag}(\lambda_0, \dots, \lambda_{k-1}) \quad \boxed{3}$$

$$= \sum_{j=0}^{k-1} d_j^2 v_j v_j^T = V D^2 V^T$$

$$\bullet (X^T X)^{-1} = \sum_{j=0}^{k-1} \lambda_j^{-1} v_j v_j^T = V \Lambda^{-1} V^T$$

$$= \sum_{j=0}^{k-1} d_j^{-2} v_j v_j^T$$

Is it inverse?

$$V \Lambda V^T \cdot V \Lambda^{-1} V^T = V \underbrace{\Lambda \Lambda^{-1}}_I V^T = V V^T = I$$

$$\bullet \text{tr}((X^T X)^{-1}) = \sum_{j=0}^{k-1} \lambda_j^{-1} = \sum_{j=0}^{k-1} \frac{1}{d_j^2} \quad \text{yes}$$

REMARK: If $\text{rank}(X) = r < k$, then

$$d_0 \geq \dots \geq d_{r-1} > d_r = \dots = d_{k-1} = 0$$

Moore-Penrose pseudoinverse of $X^T X$ is

$$(X^T X)^+ = \sum_{j=0}^{r-1} \frac{1}{d_j^2} v_j v_j^T$$

What happens if $r=k$ (full-rank) but

$d_{k-1} \rightarrow 0$ then

- $X^T X$ tends to a singular matrix
- columns of X tend to being linearly dependent
- $\text{tr}((X^T X)^{-1}) \rightarrow \infty$

\equiv multicollinearity

10.1.2 Multicollinearity and its impact on precision of the LSE

4

$$\begin{aligned} \text{(i)} \quad \hat{Y} &= (\hat{y}_1, \dots, \hat{y}_n)^T = HY, \quad H = X(X^T X)^{-1} X^T = \underbrace{Q Q^T}_{\substack{\text{orthonormal} \\ \text{basis of } \mathcal{M}(X)}} \\ &= \text{BLUE of } \mu = X\beta = E(Y|Z) \\ \text{var}(\hat{Y}|Z) &= \sigma^2 H \\ \sum_{i=1}^n \text{var}(\hat{y}_i|Z) &= \text{tr}(\sigma^2 H) = \sigma^2 \text{tr}(H) = \sigma^2 \text{tr}(Q Q^T) = \\ &= \sigma^2 \text{tr}(Q^T Q) = \sigma^2 \text{tr}(I_k) = \sigma^2 k \end{aligned}$$

no problem

$$\begin{aligned} \text{(ii)} \quad \hat{\beta} &= (\hat{\beta}_0, \dots, \hat{\beta}_{k-1})^T = (X^T X)^{-1} X^T Y \\ &= \text{BLUE of } \beta \\ \text{var}(\hat{\beta}|Z) &= \sigma^2 (X^T X)^{-1} \\ \sum_{j=0}^{k-1} \text{var}(\hat{\beta}_j|Z) &= \text{tr}(\sigma^2 (X^T X)^{-1}) = \sigma^2 \text{tr}((X^T X)^{-1}) = \\ &= \sigma^2 \sum_{j=0}^{k-1} \frac{1}{d_j^2} \\ &\quad \downarrow d_{k-1} \rightarrow 0 \\ &\quad \infty \end{aligned}$$

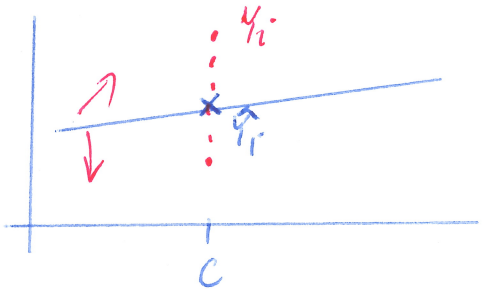
Multicollinearity

- no impact on precision of LSE of $\mu = E(Y|Z)$
- possibly serious inflation of the standard errors of LSE of β

47

Extreme case $X = (1, c1)$, rank $(X_{n \times 2}) = 1$

model $E(Y/Z=z) = \beta_0 + \beta_1 c$



$\hat{Y}_1 = \dots = \hat{Y}_n = \bar{Y}$

"LSE" of $\beta = (\beta_0, \beta_1)^T$

We want to write

$\bar{Y} = \hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 \cdot c$

" X_i (regressor)"

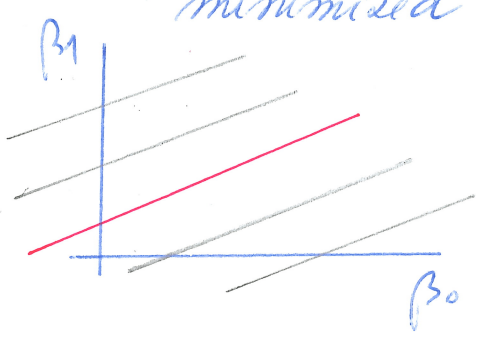
→ infinitely many pairs

$(\hat{\beta}_0, \hat{\beta}_1)^T$ satisfy $\bar{Y} = \hat{\beta}_0 + \hat{\beta}_1 \cdot c$

\equiv "var($\hat{\beta}_0/Z$) = var($\hat{\beta}_1/Z$) = ∞ "

Differently: $SS(\beta) = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 c)^2$ is

minimized for $(\beta_0, \beta_1)^T$ satisfying $\beta_0 + \beta_1 c = \bar{Y}$



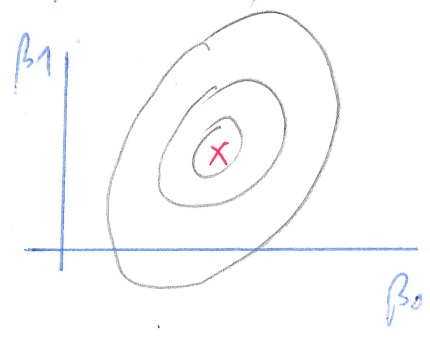
← contours of $SS(\beta)$

\equiv "inverted roof"

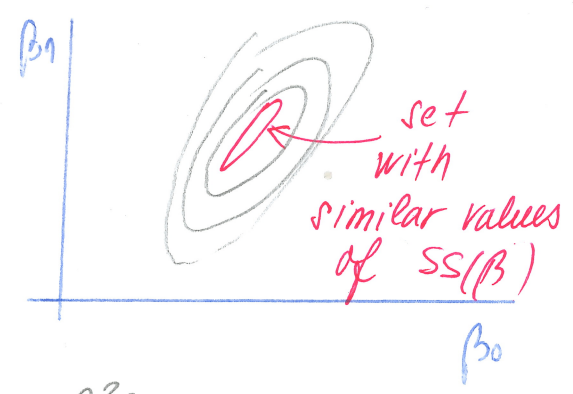
(root of a house after hurricane)

$X = \begin{pmatrix} 1 & X_1 \\ \vdots & \vdots \\ 1 & X_n \end{pmatrix}$, rank $(X) = 2$ with singular values $d_0 \geq d_1 > 0$

$SS(\beta) = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$



$d_1 \rightarrow 0$



also remember that $\frac{\partial^2 SS}{\partial \beta \partial \beta^T} = 2X^T X$

Lemma 10.1 Bias in estimation of the squared norms 5

MOTIVATION: suppose we want to estimate

$$\|E(\hat{Y}|Z)\|^2 = \|X\beta\|^2 \quad \text{or} \quad \|\beta\|^2$$

We know: $E(\hat{Y}|Z) = X\beta$ ($= E(Y|Z)$)

$$E(\hat{\beta}|Z) = \beta$$

$\rightarrow \|\hat{Y}\|^2$ estimator of $\|X\beta\|^2$

$\|\hat{\beta}\|^2$ estimator of $\|\beta\|^2$

- both estimators are biased, but how large is this bias?

\rightarrow let $Y|Z \sim (X\beta, \sigma^2 I_n)$, $\text{rank}(X_n \times n) = k$. Then

$$E(\|\hat{Y}\|^2 - \|X\beta\|^2 | Z) = \sigma^2 k \quad \begin{array}{l} \leftarrow \text{multicollinear.} \\ \text{no problem} \end{array}$$

$$E(\|\hat{\beta}\|^2 - \|\beta\|^2 | Z) = \sigma^2 \text{tr}(X^T X)^{-1} \quad \downarrow \text{problem}$$

Proof: all expectations conditional, given Z :

$$\bullet E\|\hat{Y} - X\beta\|^2 = E \sum_{i=1}^n \underbrace{(\hat{Y}_i - X_i^T \beta)^2}_{\text{var}(\hat{Y}_i)} = \sum_{i=1}^n \text{var}(\hat{Y}_i) = \text{tr}(\sigma^2 H) = \sigma^2 k$$

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

$$= E\|\hat{Y}\|^2 + \|X\beta\|^2 - 2\beta^T X^T E\hat{Y} = E\|\hat{Y}\|^2 - \|X\beta\|^2$$

$$\Rightarrow E\|\hat{Y}\|^2 - \|X\beta\|^2 = \sigma^2 k$$

In the same way since $E\hat{\beta} = \beta$:

$$\bullet E\|\hat{\beta} - \beta\|^2 = \sum_{j=0}^{k-1} \text{var}\hat{\beta}_j = \text{tr}(\text{var}\hat{\beta}) = \sigma^2 \text{tr}(\sigma^2 (X^T X)^{-1}) = \sigma^2 \text{tr}(X^T X)^{-1}$$

$$\bullet E\|\hat{\beta} - \beta\|^2 = E\|\hat{\beta}\|^2 - \|\beta\|^2 \Rightarrow E\|\hat{\beta}\|^2 - \|\beta\|^2 = \sigma^2 \text{tr}(X^T X)^{-1}$$

10.1.3 Variance inflation factor and tolerance

6

most important items from slides:

$$T_y = \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2} = \sqrt{SS_T}$$

$$R^2 = 1 - \frac{sse}{SS_T} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\underbrace{\sum_{i=1}^n (y_i - \bar{y})^2}_{T_y^2}}$$

Consider linear model with

$$X^j = \text{response}, \quad X^{(-j)} = (1, X^1, \dots, X^{j-1}, X^{j+1}, \dots, X^{k-1})$$

7

as model matrix

$$\rightarrow \hat{X}^j = \begin{pmatrix} \hat{X}_{1j} \\ \vdots \\ \hat{X}_{nj} \end{pmatrix} = \text{fitted values}$$

$$\bar{X}^j = \text{sample mean over } X^j$$

$$T_j = \sqrt{\sum_i (X_{ij} - \bar{X}^j)^2} = \sqrt{SS_T(\text{model } M_j: X^j \sim X^{(-j)})}$$

$$R_j^2 = 1 - \frac{\sum_i (X_{ij} - \hat{X}_{ij})^2}{\sum_i (X_{ij} - \bar{X}^j)^2} = 1 - \frac{sse(X^j \sim X^{(-j)})}{SS_T(X^j \sim X^{(-j)})}$$

slide + additional comments

8

$R^2 \equiv$ how closely can Y be expressed as a linear combination of X^1, \dots, X^{k-1}

$R_j^2 \equiv$ how closely can X^j be expressed as a linear combination of $X^1, \dots, X^{j-1}, X^{j+1}, \dots, X^{k-1}$

$R_j^2 = 1 \equiv X^j$ is a perfect linear combination of $X^1, \dots, X^{j-1}, X^{j+1}, \dots, X^{k-1}$

7

Theorem 10.2 Estimated variances of the LSE 9
of the regression coefficients

For a given dataset for which a linear model $Y|Z \sim (X\beta, \sigma^2 I_n)$, $\text{rank}(X_{n \times k}) = k$, $X = (1, X_1, \dots, X_{k-1})$ is applied, diagonal elements of the matrix $\widehat{\text{var}}(\hat{\beta} | Z) = \text{MSE} (X^+ X)^{-1}$, can also be calculated, for $j = 1, \dots, k-1$ as

$$\widehat{\text{var}}(\hat{\beta}_j | Z) = \left(\frac{T_y}{T_j} \right)^2 \cdot \frac{1-R^2}{n-k} \cdot \frac{1}{1-R_j^2}$$

square root of this is standard error of $\hat{\beta}_j$

not related to collinearity of regressors

quantifies strength of collinearity of regressor j with remaining regressors.

Proof: see notes (if interested)

Def 10.1 Variance inflation factor and tolerance 10

For given $j = 1, \dots, k-1$, the variance inflation factor and the tolerance of the j th regressor of the linear model $Y|Z \sim (X\beta, \sigma^2 I_n)$, $\text{rank}(X_{n \times k}) = k$ are values VIF_j and Toler_j , respectively, defined as

$$\text{VIF}_j = \frac{1}{1-R_j^2}, \quad \text{Toler}_j = 1-R_j^2 = \frac{1}{\text{VIF}_j}$$

$$R_j^2 = 0 \quad \equiv \quad \text{VIF}_j = 1, \quad \text{Toler}_j = 1$$

$$R_j^2 \rightarrow 1 \quad \equiv \quad \text{VIF}_j \rightarrow \infty, \quad \text{Toler}_j \rightarrow 0.$$

Interpretation and use of VIF

Remember: $(1-\alpha) 100\%$ confidence interval for β_j (under normality):

$$\hat{\beta}_j \pm t_{n-k} \left(1 - \frac{\alpha}{2}\right) \sqrt{\text{var}(\hat{\beta}_j | Z)}$$

$$\hat{\beta}_j \pm t_{n-k} \left(1 - \frac{\alpha}{2}\right) \sqrt{\text{VIF}_j}$$

not related to collinearity \swarrow
 related to collinearity \nwarrow

$$\rightarrow \text{VIF}_j = \left(\frac{\text{VOL}_j}{\text{VOL}_{0,j}} \right)^2$$

VOL_j = length (volume) of the confidence interval for β_j

$\text{VOL}_{0,j}$ = length (volume) ... if it was $R_j^2 = 0$

Note: In case of categorical covariates their collinearity is viewed through possible inflation of $\hat{\beta}^2 \equiv \text{LSE}$ of regress. parameters based on (pseudo)contrast parameterization

\rightarrow confidence ellipsoid for β^2 : # of levels \downarrow

$$\alpha \beta^2: (\beta^2 - \hat{\beta}^2)^T (\text{MSE} \underline{V}^2)^{-1} (\beta^2 - \hat{\beta}^2) < (G^{-1})_{G, p+k}^{-1} (1-\alpha)$$

\rightarrow volume VOL_Z block of matrix $(X^T X)^{-1}$

$\rightarrow \text{VOL}_{0,Z}$ = volume of above ellipsoid if all columns in X corresponding to (pseudo) contrasts are orthogonal to remaining columns

→ generalized VIF

11

$$gVIF_z = \left(\frac{VOL_z}{VOL_{0,z}} \right)^2$$

- to take into account "dimensionality"

$$\left(\frac{VOL_z}{VOL_{0,z}} \right)^{\frac{1}{m}}$$

is to be calculated and compared,
i.e., calculate and compare values
of $gVIF_z^{\frac{1}{2m}}$, $m = G-1$

→ R function vif
(package car)

10.1.4 Basic treatment of multicollinearity

12

remember:

$$\sum \text{var}(\hat{\beta}_j | Z) = \sigma^2 k$$

$$\sum \text{var}(\hat{\beta}_j | Z) = \sigma^2 \underbrace{\sum_{j=0}^{k-1} \frac{1}{d_j^2}}_{\rightarrow \infty, d_{k-1} \rightarrow 0}$$

$\rightarrow \infty, d_{k-1} \rightarrow 0$

+ comments using slides

13

Example: $Y = IQ$

$Z = (\text{gender}, z_{n7}, z_{n8})$

$z_{n7} \approx z_{n8}$ (corr = 0.95)

$$E(Y | Z) = \beta_0 + \beta_1 \mathbb{I}(\text{boy}) + \underbrace{\beta_2 z_{n7} + \beta_3 z_{n8}}$$

$$\approx \beta_2 z_n + \beta_3 z_n$$

$$= \underbrace{(\beta_2 + \beta_3)}_{\text{this sum identifiable}} \cdot z_n$$

this sum identifiable
from data, how to factorize
it if z_{n7} and z_{n8} both in model?

$\rightarrow \infty$ many ways

15-18

REMEMBER: If multicollinearity not treated
it is very dangerous to claim that
some factor is not associated with
response because it's non-significant
in the model!

11

