## Chapter 10. Validation of Regression Models

SKIP!

### Chapter 11. Multicollinearity

$$\mathbf{y} = X \boldsymbol{eta} + \boldsymbol{\epsilon}, \quad \hat{\boldsymbol{eta}} = (X'X)^{-1}X'\mathbf{y}, ext{ if } (X'X) ext{ is non-singular}.$$

**Q**: When is X'X singular?

<u>Ans</u>: When at least one column of X is *linearly dependent* on the other columns. We say that **collinearity** (or **multicollinearity**) exists among the columns of X.

<u>**Def**</u>: The data are said to be **ill-conditioned** if there is a "linear dependency" in the columns of X.

i.e.  $\left. \begin{array}{c} \det(X'X) \approx 0 \\ Var(\hat{\beta}_i) & \text{large} \end{array} \right\}$  Inadequate data.

Read Section 11.2 for sources of multicollinearity.

Centering and Scaling Data: Consider

$$\mathbf{Y}\sqrt{S_{yy}} = \beta_1 \sqrt{S_{11}} \mathbf{X}_1 + \dots + \beta_k \sqrt{S_{kk}} \mathbf{X}_k + \boldsymbol{\epsilon}',$$

where 
$$Y_i = (y_i - \bar{y})/\sqrt{S_{yy}}$$
,  $X_{ji} = (x_{ji} - \bar{x}_j)/\sqrt{S_{jj}}$ ,  
 $j = 1, \dots, k; i = 1, \dots, n$ . Thus, we have

$$\mathbf{Y} = \alpha_1 \mathbf{X}_1 + \dots + \alpha_k \mathbf{X}_k + \epsilon'', \text{ where } \alpha_j = \beta_j (\frac{S_{jj}}{S_{yy}})^{1/2},$$

白 ト イヨト イヨト

 $j=1,\ldots,k.$ 

 $\therefore \boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_k)'$  can be estimated through the normal equations

$$\left(\begin{array}{ccc} 1 & \cdots & r_{ij} \\ \vdots & 1 & \vdots \\ r_{ij} & \cdots & 1 \end{array}\right) \left(\begin{array}{c} \hat{\alpha}_1 \\ \vdots \\ \hat{\alpha}_k \end{array}\right) = \left(\begin{array}{c} r_{1y} \\ \vdots \\ r_{ky} \end{array}\right), \text{ in correlation form.}$$

Here

$$r_{ij} = rac{S_{ij}}{\sqrt{S_{ii}S_{jj}}} = corr(x_i, x_j), ext{ and } r_{iy} = rac{S_{iy}}{\sqrt{S_{ii}S_{yy}}} = corr(x_i, y).$$

▲口▼▲□▼▲目▼▲目▼ 目 ろくの

Then 
$$\hat{\beta}_i = \hat{\alpha}_i (S_{yy}/S_{ii})^{1/2}, i = 1, \dots, k$$
 and  
 $\hat{\beta}_0 = \bar{y} - \hat{\alpha}_1 \bar{x}_1 - \dots - \hat{\alpha}_k \bar{x}_k.$ 

**<u>Note</u>**: People in some fields argue that  $\hat{\alpha}_1, \ldots, \hat{\alpha}_k$  are more meaningful in interpreting the regression.

**<u>Ex.</u>** k = 2,  $y = \beta_1 x_1 + \beta_2 x_2 + \epsilon$  (centered and scaled).

$$X'X = \left(\begin{array}{rrr} 1 & r_{12} \\ r_{12} & 1 \end{array}\right)$$

and

$$C = (X'X)^{-1} = \begin{pmatrix} (1 - r_{12}^2)^{-1} & -r_{12}/(1 - r_{12}^2) \\ (1 - r_{12}^2)^{-1} \end{pmatrix}$$
  
$$\therefore \quad \hat{\beta}_1 = \frac{r_{1y} - r_{12}r_{2y}}{1 - r_{12}^2}, \quad \hat{\beta}_2 = \frac{r_{2y} - r_{12}r_{1y}}{1 - r_{12}^2}.$$
  
$$Var(\hat{\beta}_1) = Var(\hat{\beta}_2) = \sigma^2/(1 - r_{12}^2),$$

and  $Cov(\hat{\beta}_1, \hat{\beta}_2) = -\sigma^2 r_{12}/(1 - r_{12}^2).$ 

If  $|r_{12}| \rightarrow 1$  ( $x_1, x_2$  possess strong multicollinearity), then the LSE of  $\beta$  have large variances and covariance. For k > 2,

$$C_{jj} = \frac{1}{1 - R_j^2}, j = 1, \dots, k,$$
 (Exercise)

where  $R_j^2$  is the coefficient of multiple determination from the regression of  $x_i$  on the remaining k - 1 regressors.

:.  $R_j^2 \rightarrow 1$  if  $x_j$  has strong linear relationship with some subsets of the other regressors.

Thus,  $Var(\hat{\beta}_j) \rightarrow large.$ 

Consider

 $E(\hat{\beta} - \beta)'(\hat{\beta} - \beta) = \text{expected squared distance} = \sum_{j=1}^{k} E(\hat{\beta}_j - \beta_j)^2$  $= \sum_{j=1}^{k} Var(\hat{\beta}_j) = \sigma^2 tr(X'X)^{-1} = \sigma^2 \sum_{j=1}^{k} \frac{1}{\lambda_j},$ 

where  $\lambda_j$ 's are the eigenvalues of X'X.

Strong multicollinearity causes at least one  $\lambda_j$  to be small, thus,  $E(\hat{\beta} - \beta)'(\hat{\beta} - \beta)$  will be large. Moreover,

$$E(\hat{\beta}'\hat{\beta}) = \beta'\beta + \sigma^2 tr(X'X)^{-1} \ge \beta'\beta.$$

Hence, the vector  $\hat{\beta}$  is generally longer then  $\beta$ .

- Scatter plots of (x<sub>i</sub>, x<sub>j</sub>). Examine the correlation matrix only pairs of regressors.
- By experience or prior knowledge.
  - i)  $\hat{\beta}_i$  has different sign from anticipated.
  - ii) Important explanatory variable yields small *t*-statistic.
  - iii) Sensitive of deletion of a row or a column from X.

2. Check VIF's (Variance Inflation Factors).

If  $\mathbf{x}_j$  is orthogonal to all other columns of X,  $\forall j$ ,

$$X'X = \left(egin{array}{cc} S_{11} & 0 \ & S_{jj} \ & 0 & S_{kk} \end{array}
ight)$$
 or  $R = I.$ 

:.  $Var(\hat{\beta}_j) = \sigma^2 / S_{jj}$ . Consider

$$VIF_j = Var(\hat{\beta}_j) / \left(\frac{\sigma^2}{S_{jj}}\right), j = 1, \dots, k.$$

- **<u>Note</u>**: 1.  $VIF_j$  is a measure of how much  $\sigma^2/S_{jj}$  is inflated by columns of X to  $\mathbf{x}_j$ . It's the combined effect of the dependencies among the relationship of other regressors on the variance of  $\hat{\beta}_j$ .
  - 2. For centered and scaled data,  $S_{jj} = 1$  (in R),

$$VIF_j = C_{jj} = \frac{1}{1 - R_j^2}, j = 1, \dots, k.$$

 $C_{jj} 
ightarrow 1(R_j^2 
ightarrow 0)$  orthogonal.

 A large VIF (> 5 or 10) indicates a strong multicollinearity.

- 3. Eigensystem analysis of X'X.
  - Let λ<sub>1</sub>, λ<sub>2</sub>,..., λ<sub>p</sub> be the eigenvalues of X'X. Then one or more of the λ<sub>j</sub>'s will be ≈ 0, if there is one or more near-linear dependencies in the data. Let

$$\kappa = \frac{\lambda_{\max}}{\lambda_{\min}} = \text{ condition number of } X'X$$
  

$$\kappa_j = \frac{\lambda_{\max}}{\lambda_j} = \text{ condition indices of } X'X \text{ for } j = 1, \dots, k.$$

- **Note**: 1.  $\kappa$  measures the spread in the eigenvalue spectrum of X'X.
  - 2.  $\kappa \leq$  100, okay;

 $100 < \kappa < 1000$ , moderate to strong multicollinearity;  $\kappa \geq 1000$ , severe multicollinearity.

If λ<sub>j</sub> ≈ 0 then the elements of the associated eigenvector
 t<sub>j</sub> = (t<sub>j1</sub>,..., t<sub>jk</sub>)' describe the nature of the linear dependence. i.e. ∑<sub>i=1</sub><sup>k</sup> t<sub>ji</sub>x<sub>i</sub> ≈ 0, some of t<sub>ji</sub> ≈ 0. Find the linear relationship, e.g. x<sub>1</sub> ≈ t<sub>2</sub><sup>\*</sup>x<sub>2</sub> + t<sub>4</sub><sup>\*</sup>x<sub>4</sub>.

(2) (Belsley's Method.)

1. Decompose X'X such that  $X'X = T\Lambda T'$  where  $T'T = TT' = I_k$  and  $\Lambda = diag(\lambda_1, \dots, \lambda_k)$ .

<u>Note</u>: 1.  $T = (\mathbf{t}_1, \dots, \mathbf{t}_k)$  and  $\mathbf{t}_j$  is the eigenvector associated with  $\lambda_j$ . 2. Singular value decomposition of X:

$$X_{n\times k}=U_{n\times k}DT',$$

where  $T'T = I_k$  and U is an  $n \times k$  matrix whose columns are the *eigenvectors* associated with the k nonzero eigenvalues of XX' such that  $U'U = I_k$ ,  $D = diag(\mu_1, \ldots, \mu_k)$ ,  $\mu_j$ 's are called the *singular values* of X.

<u>Note</u>: 1.  $X'X = TDU'UDT' = TD^2T'$ ,  $\therefore D^2 = \Lambda$ ,  $\mu_j = \sqrt{\lambda_j}$ .

2. There will be one small singular value for *each* near-linear dependence.

Let

$$\eta_j = \frac{\mu_{max}}{\mu_j}, j = 1, \dots, k, \text{ condition indices of } X;$$
  
$$\eta = \frac{\mu_{max}}{\mu_{min}} = \text{ condition number of } X.$$

<u>Note</u>: 1. This approach deals directly with the data matrix X.
2. The algorithms for singular-value decomposition are more stable numerically than those for eigensystem analysis.

Recall  $Var(\hat{\beta}) = \sigma^2 (X'X)^{-1} = \sigma^2 T \Lambda^{-1} T'$  or

$$Var(\hat{\beta}_j) = \sigma^2 \sum_{i=1}^k \frac{t_{ji}^2}{\mu_i^2} = \sigma^2 \sum_{i=1}^k \frac{t_{ji}^2}{\lambda_i} = \sigma^2 C_{jj}.$$

:. 
$$VIF_j = \sum_{i=1}^k \frac{t_{ji}^2}{\mu_i^2} = \sum_{i=1}^k \frac{t_{ji}^2}{\lambda_i}.$$

Thus, one or more small  $\mu_i^2$  (or  $\lambda_i$ ) can inflate  $Var(\hat{\beta}_j)$  dramatically.

4. Variance-decomposition proportions.

$$\pi_{ij}=rac{t_{ji}^2/\mu_i^2}{VIF_j}, \hspace{0.2cm} i=1,\ldots,k, \hspace{0.2cm} ext{for each } j=1,\ldots,k.$$

**<u>Note</u>**: 1.  $\pi_{ij}$  measures the multicollinearity = proportion of  $Var(\hat{\beta}_j)$  contributed by the  $i^{th}$  singular value.

2.  $\eta_j > 30$  and  $\pi_{ij} > 0.5$  are recommended guidelines for detecting multicollinearity. It indicates that the corresponding regressors are of possible multicollinearity.

T

		<i>x</i> <sub>1</sub>	<i>x</i> <sub>2</sub>	 	x <sub>k</sub>
$\lambda_1$	$\eta_1$	$\pi_{11}$	$\pi_{12}$	 	$\pi_{1k}$
÷	÷	$\pi_{21}$	$\pi_{22}$	 	
$\geq$	$\leq$			÷	
÷	÷	:		÷	
÷	30	÷			
$\lambda_k$	$\eta_k$	$\pi_{k1}$	$\pi_{k2}$	 	$\pi_{kk}$
$\downarrow$	$\downarrow$				
0	$\infty$	$\sum = 1$			

6.  $|X'X| \to 0. \ (0 \le |X'X| \le 1)$  or  $100[|X'X|^{-1/2} - 1]$  large.

<u>Note</u>:  $|X'X|^{-1/2}$  is the size of the confidence region enlarged due to multicollinearity. (|X'X| = 1 if X is orthogonal.) When multicollinearity occurs, one can

1) Add prior information.

e.g. In  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$ . 'Known'  $(x_1, x_2)$  are highly correlated, i.e.  $\beta_1 = c\beta_2$ .  $\implies$  Let  $z = x_1 + cx_2$  and consider  $y = \beta_0 + \beta_1 z + \epsilon$ . 2) Combine models.  $y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \epsilon_i$ ,  $(x_{i1}, x_{i2})$  collinear.

i) Use *previous* information of (y<sub>i</sub><sup>\*</sup>, x<sub>i</sub><sup>\*</sup>) to fit y<sub>i</sub><sup>\*</sup> = β<sub>0</sub><sup>\*</sup> + β<sub>2</sub><sup>\*</sup>x<sub>i</sub><sup>\*</sup> + ϵ<sub>i</sub><sup>\*</sup> ⇒ get β̂<sub>0</sub><sup>\*</sup>, β̂<sub>2</sub><sup>\*</sup>.
ii) Let ỹ<sub>i</sub> = y<sub>i</sub> - β̂<sub>0</sub><sup>\*</sup> - β̂<sub>2</sub><sup>\*</sup>x<sub>i</sub><sup>2</sup> and consider ỹ<sub>i</sub> = β'<sub>0</sub> + β<sub>1</sub>x<sub>i</sub> + ϵ'<sub>i</sub>.

- 3) Delete collinear variables. Must check for model adequacy, etc.
- 4) Principle component analysis.
- 5) Make transformation on x's to eliminate correlations.
- 6) Take more data.
- 7) Bayesian approach.

# Homework 10: (Page 365) 11.2, 11.3, 11.4, 11.5, 11.12. Due: Jan. 16, 2009.

3

Intended to overcome "ill-conditioned" solutions.

 $\implies$  Find out how the ill-conditioning occurs, and add specific additional information to the problem to remove the ill-conditioning.

<u>Motivation</u>: Find a "biased" estimator of  $\beta$  that has a much smaller variance than the LSE.

Let  $\theta > 0$  (usually  $\theta \in (0, 1)$ ). Define the ridge regression estimate  $\hat{\beta}_R$  of  $\beta$  as the solution to

$$(X'X+ heta I)\hat{eta}_R=X'\mathbf{y}$$
 or  $\hat{eta}_R=(X'X+ heta I)^{-1}X'\mathbf{y}.$ 

Note: 1. Here the data are centered and scaled.

2. 
$$\hat{\boldsymbol{\beta}}_{R} = (X'X + \theta I)^{-1}(X'X)(X'X)^{-1}X'\mathbf{y} = Z_{\theta}\hat{\boldsymbol{\beta}}.$$
  $(\hat{\boldsymbol{\beta}}_{R} = \hat{\boldsymbol{\beta}}$   
when  $\theta = 0.$ )

3.  $\theta$  is called the **biasing parameter** since  $E(\hat{\beta}_R) = Z_{\theta}\beta$  is biased.

4. 
$$Var(\hat{\beta}_R) = \sigma^2 (X'X + \theta I)^{-1} X' X (X'X + \theta I)^{-1}.$$

<u>**Result</u></u>: There exists a \theta^\* > 0 such that MSE\_{\theta^\*}(\hat{\beta}\_R) < MSE(\hat{\beta}) if \beta'\beta is bounded.</u>** 

**Q**:  $\theta^*$  depends on  $\sigma^2$  and  $\beta$ .  $\theta^* = ???$ 

#### Solutions:

1. Ridge trace: plot of  $\hat{\beta}_{jR}(\theta)$  versus  $\theta$ , for each j = 1, ..., k.



At a certain value of  $\theta$ , the system will stabilize.

Note that

$$SS_{Res}(\hat{\beta}_R) = (\mathbf{y} - X\hat{\beta}_R)'(\mathbf{y} - X\hat{\beta}_R)$$
  
=  $(\mathbf{y} - X\hat{\beta})'(\mathbf{y} - \hat{\beta}) + (\hat{\beta}_R - \hat{\beta})'X'X(\hat{\beta}_R - \hat{\beta})$   
$$SS_{Res} \uparrow \text{ in } \theta.$$

Large  $\theta \implies$  ridge trace gets stable, but with large  $SS_{Res}$ . So  $\theta$  can't be too large.

 $\therefore$  Select a reasonable *small*  $\theta$  at which  $\hat{\beta}_R$  are stable.

<u>Note</u>: Usually, the ridge estimates do a better job of predicting future observations than LSE.

2. 
$$\theta^{\star} = p\hat{\sigma^2}/(\hat{\beta}'\hat{\beta})$$
, where  $\hat{\sigma^2} = MS_{Res}$  based on LSE.

Note that



sensible from a Bayesian point of view (as  $\beta \sim N(0, \sigma_{\beta}^2 I))$ .

・ロト ・聞 ト ・言 ト ・言 ト ・ ほ ・ うんの

3. Iterative procedure. Let  $\theta_0^{\star} = p\hat{\sigma^2}/(\hat{\beta}'\hat{\beta}) \Longrightarrow \hat{\beta}_R(\theta_0^{\star})$ . Compute  $\theta_{j+1}^{\star} = p\hat{\sigma^2}/(\hat{\beta}'_R(\theta_j^{\star})\hat{\beta}_R(\theta_j^{\star})), \quad j = 0, 1, 2, \dots$ until 'convergence'!

i.e. Stop until 
$$|\theta_{j+1}^{\star} - \theta_{j}^{\star}|/\theta_{j}^{\star} \leq \delta$$
.  
e.g.  $\delta = 20(tr(X'X)^{-1}/k)^{-1.3}$  is suggested.

#### Justification:

- Bayesian regression analysis:  $eta \sim \pi(eta)$  prior distribution,  $\hat{eta} = ?$
- 2 LSE with restriction: Known  $\beta'\beta \leq c^2$ ,  $\hat{\beta} = ?$