

- \* Multicollinearity

- \* Exact multicollinearity

Consider a three-variable regression:

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \varepsilon_i \quad (1)$$

And suppose we have exact collinearity in:

$$x_{i3} = \alpha_1 x_{i1} + \alpha_2 x_{i2} \quad (ii)$$

Then, model (1) is not well-determined or estimable. The reason for this is if we swap  $(\beta_1, \beta_2, \beta_3)$  with any other set  $(\beta_1^*, \beta_2^*, \beta_3^*)$  where

$$\begin{aligned} \beta_1 + \alpha_1 \beta_3 &= \beta_1^* + \alpha_1 \beta_3^* \quad \text{and} \\ \beta_2 + \alpha_2 \beta_3 &= \beta_2^* + \alpha_2 \beta_3^* \end{aligned} \quad (A)$$

then model (1) is unchanged!

- \* In such situations,  $X^T X$  is a singular matrix and  $\hat{\beta}_{OLS}$  does not exist.

- \* Approximate multicollinearity

In most cases, we don't have perfect linear relationships like (ii), but instead have approximate linear dependence.

- In this case,  $X^T X$  will be ill-conditioned - its inverse is numerically unstable.
- It follows that some of the parameter estimates, or ~~at least~~ some combinations of them, will have large variances.

Implication: large perturbations of the parameters will have only a small effect on the fit of the model.

## \* Detecting Multicollinearity

### \* I. Variance Inflation Factors

- For all columns of  $X$  except the intercept, we must first standardize the columns to have mean 0 and variance 1.

R code: use `scale(X)`

- The  $j^{\text{th}}$  variance inflation factor is the  $j^{\text{th}}$  diagonal entry of  $(X^T X)^{-1}$  once  $X$  has been standardized as above. Notation:  $V_j$
- \* Idea: if there were no multicollinearity, the cross-products of each pair of variables would be 0 and  $X^T X$  would be the identity matrix  $\Rightarrow V_j \equiv 1 \quad \forall j$

- \* There are no formal tests for large VIF's, but it is accepted that significant multicollinearity is present if

$$V_j > 10 \quad \text{for any } j = 1, \dots, p.$$

- \* Note: this does not indicate where collinearity occurs, only that it is present

## II. Singular Value Criteria for MC

- \* We start with the singular value decomposition of the standardized version of  $X$ :

$$X = U D V^T, \quad \text{where} \quad (\ddot{u})$$

$U$  is  $n \times p$ ,  $V$  is  $p \times p$  and  $U^T U = V^T V = I_p$  and  $D$  is a  $p \times p$  ~~not~~ diagonal matrix with entries  $\mu_1, \dots, \mu_p$ , which are the singular values of  $X$ .

( $\ddot{u}$ ) implies

$$X^T X = V D^2 V^T \Rightarrow (X^T X)^{-1} = V D^{-2} V^T \quad (**)$$

(\*\*)  $\Rightarrow$  small singular values imply large entries in  $(X^T X)^{-1}$  (and hence large variance),



\* Define

$$\tau_k = \frac{\mu_{\max}}{\mu_k}$$

as the  $k^{\text{th}}$  condition index of  $X$ , where

$$\mu_{\max} = \max \{ \mu_1, \dots, \mu_p \}.$$

Large values of  $\tau_k$  indicate multicollinearity.  
Rule of thumb:

$\tau_k \in (30, 100) \Rightarrow$  moderate to strong association

$\tau_k < 10 \Rightarrow$  weak association.

\* Again,  $\tau_k$  does not show where association occurs. To measure this, we define

$$\pi_{kj} = \frac{v_{jk}^2 / \mu_k^2}{\sum_{k'} (v_{jk'}^2 / \mu_{k'}^2)}$$

= proportion of variance of the  $j^{\text{th}}$  parameter estimate that is accounted for by the  $k^{\text{th}}$  singular value.

\* Rule of thumb: Suspect variables are those w/ high  $\tau_k$  and at least two large values of  $\pi_{kj}$ .

Monday  $\rightarrow$  review

## Remedies for MC

- I. Ridge regression  $\rightarrow$  reduces the effects of MC through regularization
- II. Principal components regression:  
Main idea:
  - ① Identify  $k$  principal components of the columns of  $X$ . These represent the  $k$  directions of most variability in  $X$ .
  - ② Regress  $Y$  on the principal components.  
Since the directions are orthogonal, we have no issues w/ MC but we generally lose interpretability.
  - ③ Removal of suspicious variables from MC diagnostics.
  - ④ Partial Least Squares (not covering here, but it is an option).