

Multicollinearity

San-Teng Huang, Shang-Chien Ho

National Dong Hwa University

2018/12/19

Outline

Why Collinearity Is a Problem

Matrix-Geometric Perspective on Multicollinearity

Multiple Linear Regression

- The coefficients of the estimates:

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

- The variance of the estimates:

$$\text{Var} [\hat{\beta}] = \sigma^2 (X^T X)^{-1}$$

If that matrix isn't exactly singular, but is close to being non-invertible, the variances will become huge.

Collinearity

There are several equivalent conditions for any square matrix, say u , to be singular or non-invertible:

- The determinant $\det u$ or $|u|$ is 0.
- At least one eigenvalue of u is 0.
- u is rank deficient, meaning that one or more of its columns (or rows) is equal to a linear combination of the other rows.

- Dealing with Collinearity by Deleting Variables
- Diagnosing Collinearity Among Pairs of Variables
- Geometric Perspective
- Why Multicollinearity Is Harder

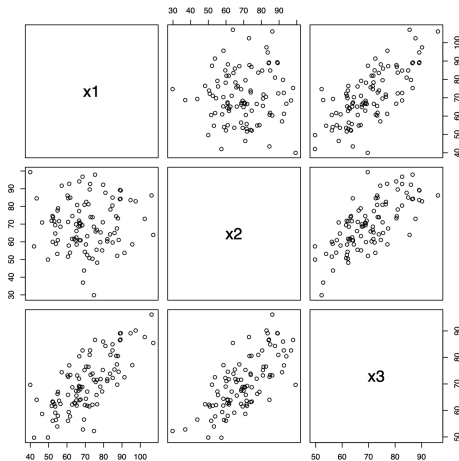


Figure: suppose X_1 and X_2 are independent Gaussians, of equal variance σ^2 , and X_3 is their average, $X_3 = (X_1 + X_2) / 2$

Multicollinearity

Multicollinearity means that

$$c_1X_1 + c_2X_2 + \dots c_pX_p = \sum_{i=1}^p c_iX_i = c_0$$

To simplify this, let's introduce $p \times 1$ matrix $a = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_p \end{bmatrix}$, so we can

write multicollinearity as $a^T X = c_0$, for $a \neq 0$

$$\text{Var}[a^T X] = 0, \quad a \neq 0$$

$$\begin{aligned}\text{Var}[a^T X] &= \text{Var}\left[\sum_{i=1}^p c_i X_i\right] \\ &= \sum_{i=1}^p \sum_{j=1}^p c_i c_j \text{Cov}[X_i, X_j] \\ &= a^T \text{Var}[X] a\end{aligned}$$

- The eigenvectors of $\text{Var}[X]$, such that

$$\text{Var}[X] v_i = \lambda v_i$$

- The eigenvalues are all ≥ 0 .
- Any vector can be re-written as a sum of eigenvectors:

$$a = \sum_{i=1}^p (a^T v_i) v_i$$

- The eigenvectors can be chosen so that they all have length 1 and are orthogonal to each other.
- $\text{Var}[X]$ can be expressed as

$$\text{Var}[X] = V D V^T$$

$$\begin{aligned} \text{Var}[X] a &= \text{Var}[X] \sum_{i=1}^p (a^T v_i) v_i \\ &= \sum_{i=1}^p (a^T v_i) \text{Var}[X] v_i \\ &= \sum_{i=1}^p (a^T v_i) \lambda_i v_i \\ a^T \text{Var}[X] a &= \left(\sum_{i=1}^p (a^T v_i) v_i \right)^T \sum_{i=1}^p (a^T v_i) \lambda_i v_i \\ &= \sum_{i=1}^p \sum_{j=1}^p (a^T v_i) (a^T v_j) v_j^T v_i \\ &= \sum_{i=1}^p (a^T v_i)^2 \lambda_i \end{aligned}$$

- The predictors are multi-collinear if and only if $\text{Var} [X]$ has zero eigenvalues.
- Every multi-collinear combination of the predictors is either an eigenvector of $\text{Var} [X]$ with zero eigenvalue, or a linear combination of such eigenvectors.