

Section for Applied Statistics and Data Analysis

TA: Cong Mu

Office Hour: Wednesday 10:00AM - 12:00PM

September 13, 2019

1 Some Statistics

- Regression Formulation
- Least Squares Estimation
- Goodness of Fit

2 Some Programming

- Example with gala Dataset
- Exercise with teengamb Dataset

General Formulation

Suppose we have n observations and p variables.

$$y = X\beta + \epsilon,$$

where

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1} \quad X = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1p} \\ 1 & x_{21} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{np} \end{bmatrix}_{n \times (p+1)} \quad \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}_{(p+1) \times 1}.$$

Estimating Parameters

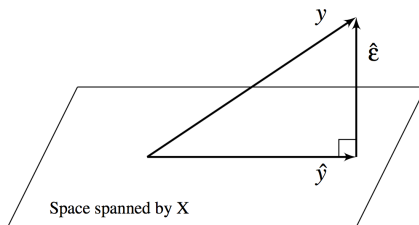


Figure 2.1 Geometrical representation of the estimation β . The data vector Y is projected orthogonally onto the model space spanned by X . The fit is represented by projection $\hat{y} = X\hat{\beta}$ with the difference between the fit and the data represented by the residual vector \hat{e} .

(Figure from Linear Models with R)

Least Squares Formulation

Recall

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \quad \text{where} \quad \boldsymbol{\epsilon} = [\epsilon_1 \quad \epsilon_2 \quad \cdots \quad \epsilon_n]^\top.$$

To minimize the sum of the squared errors

$$L = \sum_{i=1}^n \epsilon_i^2 = \boldsymbol{\epsilon}^\top \boldsymbol{\epsilon} = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}).$$

$$\frac{\partial L}{\partial \boldsymbol{\beta}} = 0 \quad \implies \quad (\mathbf{X}^\top \mathbf{X}) \boldsymbol{\beta} = \mathbf{X}^\top \mathbf{y}.$$

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}, \quad \hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

Some Examples

- $y = \mu + \epsilon$

$$X = \mathbf{1}, \beta = \mu \implies \hat{\beta} = (\mathbf{1}^\top \mathbf{1})^{-1} \mathbf{1}^\top y = \frac{1}{n} \mathbf{1}^\top y = \bar{y}.$$

- $y = \beta_0 + \beta_1 x + \epsilon \iff y = \beta_0 + \beta_1 \bar{x} + \beta_1 (x - \bar{x}) + \epsilon$

$$X = \begin{bmatrix} 1 & x_1 - \bar{x} \\ 1 & x_2 - \bar{x} \\ \vdots & \vdots \\ 1 & x_n - \bar{x} \end{bmatrix}_{n \times 2} \implies X^\top X = \begin{bmatrix} n & 0 \\ 0 & \sum_{i=1}^n (x_i - \bar{x})^2 \end{bmatrix}_{2 \times 2}.$$

$$\hat{\beta} = \begin{bmatrix} \hat{\beta}'_0 \\ \hat{\beta}_1 \end{bmatrix}_{2 \times 1} = \begin{bmatrix} \bar{y} \\ \frac{\sum_{i=1}^n (x_i - \bar{x}) y_i}{\sum_{i=1}^n (x_i - \bar{x})^2} \end{bmatrix}_{2 \times 1}.$$

Some Properties

Gauss-Markov Theorem

In a linear regression model in which the errors are uncorrelated with equal variances and expectation value of zero, the best linear unbiased estimator (BLUE) of the coefficients is given by the ordinary least squares (OLS) estimator. i.e.

$$\mathbb{E} [\hat{\beta}_{\text{OLS}}] = \beta.$$

$$\text{Var} [\hat{\beta}_{\text{OLS}}] \leq \text{Var} [\hat{\beta}] \quad \text{for} \quad \hat{\beta} \in \left\{ \hat{\beta} \mid \mathbb{E} [\hat{\beta}] = \beta \right\}.$$

Coefficient of Determination

- **Sum of Squares Total**

$$SSTO = \sum_{i=1}^n (y_i - \bar{y})^2.$$

- **Sum of Squares Regression**

$$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2.$$

- **Sum of Squares Error**

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2.$$

$$SSTO = SSR + SSE, \quad R^2 = \frac{SSR}{SSTO} = 1 - \frac{SSE}{SSTO}$$

- **Background**

The number of species found on the various Galapagos Islands. There are 30 cases (Islands) and 6 variables in the dataset.

- **Variables**

- Species: the number of species found on the island
- Area: the area of the island (km²)
- Elevation: the highest elevation of the island (m)
- Nearest: the distance from the nearest island (km)
- Scrutz: the distance from Santa Cruz Island (km)
- Adjacent: the area of the adjacent island (km²)

Thanks for listening!