

R Lab for Regression Analysis

Young-geun Kim

Department of Statistics, SKKU

dudrms33@g.skku.edu

24 Mar, 2019

Contents

1	Linear Regression Analysis	5
1.1	Relation	5
2	Simple Linear Regression	7
2.1	Model	7
2.2	Least Squares Estimation	8
2.3	Maximum Likelihood Estimation	14
	References	17

Chapter 1

Linear Regression Analysis

```
data(BioOxyDemand, package = "MPV")
(BioOxyDemand <-
  BioOxyDemand %>%
  tbl_df())
```

```
# A tibble: 14 x 2
```

	x	y
	<int>	<int>
1	3	4
2	8	7
3	10	8
4	11	8
5	13	10
6	16	11
7	27	16
8	30	26
9	35	21
10	37	9
11	38	31
12	44	30
13	103	75
14	142	90

1.1 Relation

We wonder how x affects y , especially linearly.

- Functional relation: mathematical equation,

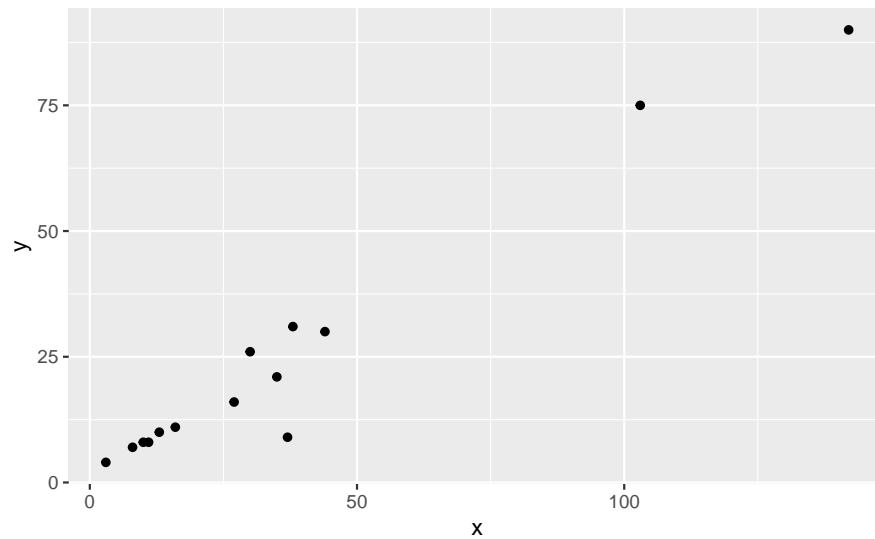
$$y = \beta_0 + \beta_1 x$$

- Statistical relation: embedded with noise

So we try to estimate

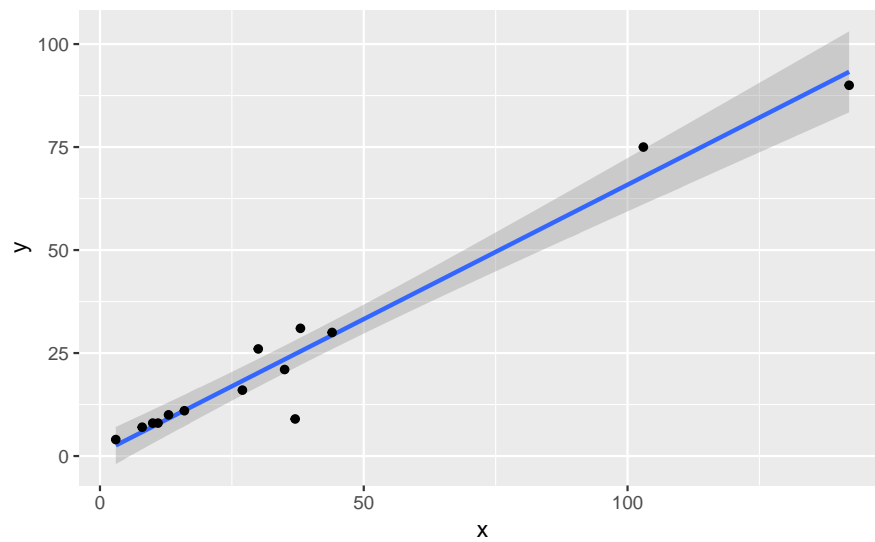
$$y = \beta_0 + \beta_1 x + \epsilon$$

```
BioOxyDemand %>%
  ggplot(aes(x, y)) +
  geom_point()
```



Looking just with the eyes, we can see the linear relationship. Regression analysis estimates the relationship statistically.

```
BioOxyDemand %>%
  ggplot(aes(x, y)) +
  geom_smooth(method = "lm") +
  geom_point()
```



Chapter 2

Simple Linear Regression

2.1 Model

```
delv <- MPV::p2.9 %>% tbl_df()
```

```
delv %>%  
  ggplot(aes(x = x, y = y)) +  
  geom_point() +  
  labs(  
    x = "Number of Cases",  
    y = "Delivery Time"  
  )
```

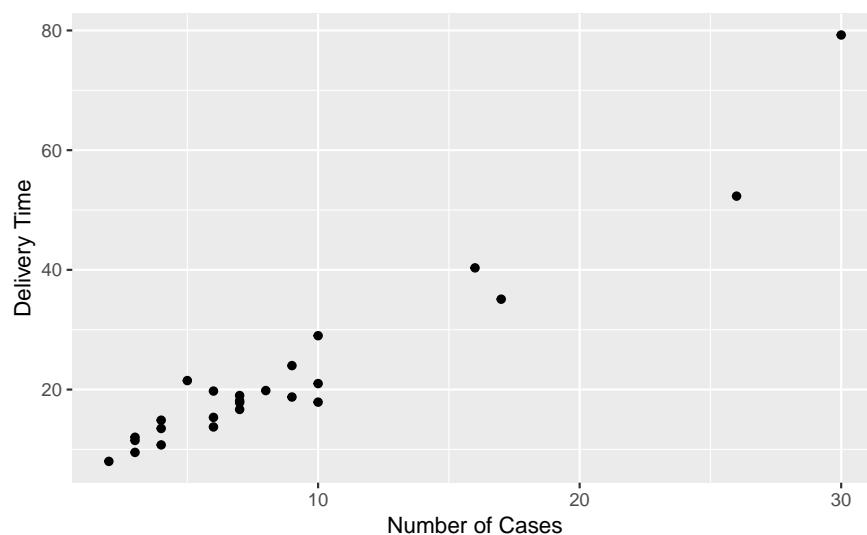


Figure 2.1: The Delivery Time Data

Given data $(x_1, Y_1), \dots, (x_n, Y_n)$, we try to fit linear model

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

Here ϵ_i is a error term, which is a random variable.

$$\epsilon \stackrel{iid}{\sim} (0, \sigma^2)$$

It gives the problem of estimating three parameters $(\beta_0, \beta_1, \sigma^2)$. Before estimating these, we set some assumptions.

1. linear relationship
2. ϵ_i s are independent
3. ϵ_i s are identically distributed, i.e. *constant variance*
4. In some setting, $\epsilon_i \sim N$

2.2 Least Squares Estimation

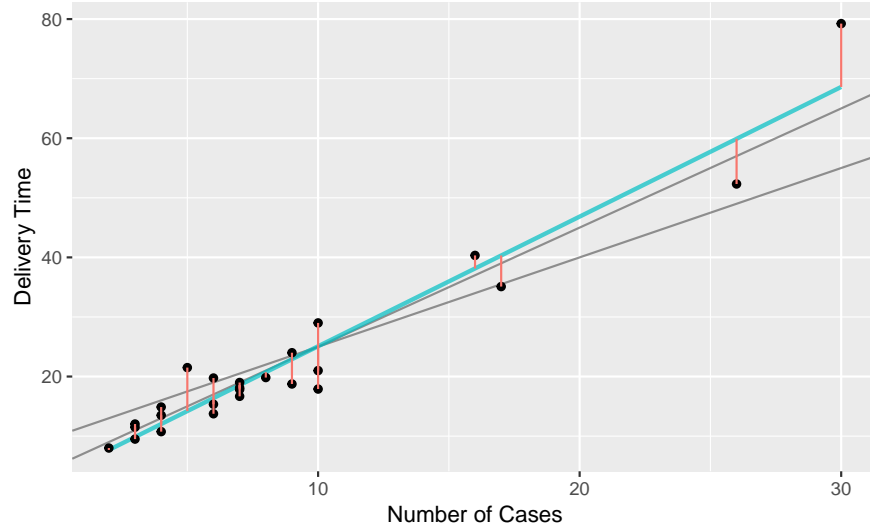


Figure 2.2: Idea of the least square estimation

We try to find β_0 and β_1 that minimize the sum of squares of the vertical distances, i.e.

$$(\beta_0, \beta_1) = \arg \min \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 x_i)^2 \quad (2.1)$$

2.2.1 Normal equations

Denote that Equation (2.1) is quadratic. Then we can find its minimum by find the zero point of the first derivative. Set

$$Q(\beta_0, \beta_1) := \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 x_i)^2$$

Then we have

$$\frac{\partial Q}{\partial \beta_0} = -2 \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 x_i) = 0 \quad (2.2)$$

and

$$\frac{\partial Q}{\partial \beta_1} = -2 \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 x_i) x_i = 0 \quad (2.3)$$

From (2.2),

$$\sum_{i=1}^n Y_i - n\hat{\beta}_0 - \hat{\beta}_1 \sum_{i=1}^n x_i = 0$$

Thus,

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x}$$

(2.3) gives

$$\sum_{i=1}^n x_i (Y_i - \bar{Y} + \hat{\beta}_1 \bar{x} - \hat{\beta}_1 x_i) = \sum_{i=1}^n x_i (Y_i - \bar{Y}) - \hat{\beta}_1 \sum_{i=1}^n x_i (x_i - \bar{x}) = 0$$

Thus,

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i (Y_i - \bar{Y})}{\sum_{i=1}^n x_i (x_i - \bar{x})}$$

Remark.

$$\hat{\beta}_1 = \frac{S_{XY}}{S_{XX}}$$

where $S_{XX} := \sum_{i=1}^n (x_i - \bar{x})^2$ and $S_{XY} := \sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})$

Proof. Note that $\bar{x}^2 = \frac{1}{n^2} \left(\sum_{i=1}^n x_i \right)^2$. Then we have

$$\begin{aligned} S_{XX} &= \sum_{i=1}^n (x_i - \bar{x})^2 \\ &= \sum_{i=1}^n x_i^2 - 2 \sum_{i=1}^n x_i \bar{x} + \sum_{i=1}^n \bar{x}^2 \\ &= \sum_{i=1}^n x_i^2 - \frac{2}{n} \left(\sum_{i=1}^n x_i \right)^2 + \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \\ &= \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \end{aligned} \quad (2.4)$$

It follows that

$$\begin{aligned}
\hat{\beta}_1 &= \frac{\sum x_i(Y_i - \bar{Y})}{\sum x_i(x_i - \bar{x})} \\
&= \frac{\sum x_i(Y_i - \bar{Y}) - \bar{x} \sum (Y_i - \bar{Y})}{\sum x_i^2 - \frac{1}{n}(\sum x_i)^2} \quad \because \sum (Y_i - \bar{Y}) = 0 \\
&= \frac{\sum (x_i - \bar{x})(Y_i - \bar{Y})}{\sum x_i^2 - \frac{1}{n}(\sum x_i)^2} \\
&= \frac{S_{XY}}{S_{XX}}
\end{aligned}$$

□

```
lm(y ~ x, data = delv)
```

Call:

```
lm(formula = y ~ x, data = delv)
```

Coefficients:

(Intercept)	x
3.32	2.18

2.2.2 Prediction and Mean response

“Essentially, all models are wrong, but some are useful.”

—George Box

Recall that we have assumed the **linear assumption** between the predictor and the response variables, i.e. the true model. Estimating β_0 and β_1 is same as estimating the *assumed true model*.

Definition 2.1 (Mean response).

$$E(Y \mid X = x) = \beta_0 + \beta_1 x$$

We can estimate this mean response by

$$\widehat{E(Y \mid x)} = \hat{\beta}_0 + \hat{\beta}_1 x \quad (2.5)$$

However, in practice, the model might not be true, which is included in ϵ term.

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

Our real problem is predicting individual Y , not the mean. The *prediction* of response can be done by

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \quad (2.6)$$

Observe that the values of Equation (2.5) and (2.6) are same. However, due to the **error term in the prediction**, it has larger standard error.

2.2.3 Properties of LSE

Parameters β_0 and β_1 have some properties related to the expectation and variance. We can notice that these lse's are **unbiased linear estimator**. In fact, these are the *best unbiased linear estimator*. This will be covered in the Gauss-Markov theorem.

Lemma 2.1.

$$S_{XX} = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 = \sum_{i=1}^n x_i (x_i - \bar{x})$$

$$S_{XY} = \sum_{i=1}^n x_i Y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n Y_i \right) = \sum_{i=1}^n Y_i (x_i - \bar{x})$$

Proof. We already proven the first part of S_{XX} . See the Equation (2.4). The second part is tivial. Since $\sum (x_i - \bar{x}) = 0$,

$$S_{XX} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n (x_i - \bar{x}) x_i$$

For the first part of S_{XY} ,

$$\begin{aligned} S_{XY} &= \sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y}) \\ &= \sum_{i=1}^n x_i Y_i - \bar{x} \sum_{i=1}^n Y_i - \bar{Y} \sum_{i=1}^n x_i + n \bar{x} \bar{Y} \\ &= \sum_{i=1}^n x_i Y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n Y_i \right) \end{aligned}$$

Second part of S_{XY} also can be proven from the definition.

$$\begin{aligned} S_{XY} &= \sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y}) \\ &= \sum_{i=1}^n Y_i (x_i - \bar{x}) - \bar{Y} \sum_{i=1}^n (x_i - \bar{x}) \\ &= \sum_{i=1}^n Y_i (x_i - \bar{x}) \quad \because \sum_{i=1}^n (x_i - \bar{x}) = 0 \end{aligned}$$

□

Lemma 2.2 (Linearity). *Each coefficient is a linear estimator.*

$$\hat{\beta}_1 = \sum_{i=1}^n \frac{(x_i - \bar{x})}{S_{XX}} Y_i$$

$$\hat{\beta}_0 = \sum_{i=1}^n \left(\frac{1}{n} - \frac{(x_i - \bar{x})}{S_{XX}} \right) Y_i$$

Proof. From lemma 2.1,

$$\begin{aligned}\hat{\beta}_1 &= \frac{S_{XY}}{S_{XX}} \\ &= \frac{1}{S_{XX}} \sum_{i=1}^n (x_i - \bar{x}) Y_i\end{aligned}$$

It gives that

$$\begin{aligned}\hat{\beta}_0 &= \bar{Y} - \hat{\beta}_1 \bar{x} \\ &= \frac{1}{n} \sum_{i=1}^n Y_i - \bar{x} \sum_{i=1}^n \frac{(x_i - \bar{x})}{S_{XX}} Y_i \\ &= \sum_{i=1}^n \left(\frac{1}{n} - \frac{(x_i - \bar{x}) \bar{x}}{S_{XX}} \right) Y_i\end{aligned}$$

□

Proposition 2.1 (Unbiasedness). *Both coefficients are unbiased.*

(a) $E\hat{\beta}_1 = \beta_1$

(b) $E\hat{\beta}_0 = \beta_0$

From the model, $Y_1, \dots, Y_n \stackrel{indep}{\sim} (\beta_0 + \beta_1 x_i, \sigma^2)$.

Proof. From lemma 2.1,

$$\begin{aligned}E\hat{\beta}_1 &= \sum_{i=1}^n \left[\frac{(x_i - \bar{x})}{S_{XX}} E(Y_i) \right] \\ &= \sum_{i=1}^n \frac{(x_i - \bar{x})}{S_{XX}} (\beta_0 + \beta_1 x_i) \\ &= \frac{\beta_1 \sum (x_i - \bar{x}) x_i}{\sum (x_i - \bar{x}) x_i} \quad \because \sum (x_i - \bar{x}) = 0 \\ &= \beta_1\end{aligned}$$

It follows that

$$\begin{aligned}E\hat{\beta}_0 &= E(\bar{Y} - \hat{\beta}_1 \bar{x}) \\ &= E(\bar{Y}) - \bar{x} E(\hat{\beta}_1) \\ &= E(\beta_0 + \beta_1 \bar{x} + \bar{\epsilon}) - \beta_1 \bar{x} \\ &= \beta_0 + \beta_1 \bar{x} - \beta_1 \bar{x} \\ &= \beta_0\end{aligned}$$

□

Proposition 2.2 (Variances). *Variances and covariance of coefficients*

$$(a) \text{Var} \hat{\beta}_1 = \frac{\sigma^2}{S_{XX}}$$

$$(b) \text{Var} \hat{\beta}_0 = \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{XX}} \right) \sigma^2$$

$$(c) \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) = -\frac{\bar{x}}{S_{XX}} \sigma^2$$

Proof. Proving is just arithmetic.

(a)

$$\begin{aligned} \text{Var} \hat{\beta}_1 &= \frac{1}{S_{XX}^2} \sum_{i=1}^n \left[(x_i - \bar{x})^2 \text{Var}(Y_i) \right] + \frac{1}{S_{XX}^2} \sum_{j \neq k}^n \left[(x_j - \bar{x})(x_k - \bar{x}) \text{Cov}(Y_j, Y_k) \right] \\ &= \frac{\sigma^2}{S_{XX}} \quad \because \text{Cov}(Y_j, Y_k) = 0 \text{ if } j \neq k \end{aligned}$$

(b)

$$\begin{aligned} \text{Var} \hat{\beta}_0 &= \sum_{i=1}^n \left(\frac{1}{n} - \frac{(x_i - \bar{x})\bar{x}}{S_{XX}} \right)^2 \text{Var}(Y_i) + \sum_{j \neq k} \left(\frac{1}{n} - \frac{(x_j - \bar{x})\bar{x}}{S_{XX}} \right) \left(\frac{1}{n} - \frac{(x_k - \bar{x})\bar{x}}{S_{XX}} \right) \text{Cov}(Y_j, Y_k) \\ &= \frac{\sigma^2}{n} - 2\sigma^2 \frac{\bar{x}}{S_{XX}} \sum_{i=1}^n (x_i - \bar{x}) + \frac{\sigma^2 \bar{x}^2 \sum_{i=1}^n (x_i - \bar{x})^2}{S_{XX}^2} \\ &= \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{XX}} \right) \sigma^2 \quad \because \sum (x_i - \bar{x}) = 0 \end{aligned}$$

(c)

$$\begin{aligned} \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) &= \text{Cov}(\bar{Y} - \hat{\beta}_1 \bar{x}, \hat{\beta}_1) \\ &= -\bar{x} \text{Var} \hat{\beta}_1 \\ &= -\frac{\bar{x}}{S_{XX}} \sigma^2 \end{aligned}$$

□

2.2.4 Gauss-Markov Theorem

Chapter 2.2.3 shows that the β_0^{LSE} and β_1^{LSE} are the **linear unbiased estimators**. Are these good? Good compared to *what estimators*? Here we consider *linear unbiased estimator*. If variances in the proposition 2.2 are lower than any parameters in this parameter family, β_0^{LSE} and β_1^{LSE} are the **best linear unbiased estimators**.

Theorem 2.1 (Gauss Markov Theorem). *$\hat{\beta}_0$ and $\hat{\beta}_1$ are BLUE, i.e. the best linear unbiased estimator.*

$$\text{Var}(\hat{\beta}_0) \leq \text{Var} \left(\sum_{i=1}^n a_i Y_i \right) \forall a_i \in \mathbb{R} \text{ s.t. } E \left(\sum_{i=1}^n a_i Y_i \right) = \beta_0$$

$$\text{Var}(\hat{\beta}_1) \leq \text{Var} \left(\sum_{i=1}^n b_i Y_i \right) \forall b_i \in \mathbb{R} \text{ s.t. } E \left(\sum_{i=1}^n b_i Y_i \right) = \beta_1$$

2.3 Maximum Likelihood Estimation

In this section, we add an assumption to an random errors ϵ_i .

$$\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

Example 2.1 (Gaussian Likelihood). Note that $Y_i \stackrel{indep}{\sim} N(\beta_0 + \beta_1 x_i, \sigma^2)$. Then the likelihood function is

$$L(\beta_0, \beta_1, \sigma^2) = \prod_{i=1}^n \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{(Y_i - \beta_0 - \beta_1 x_i)^2}{2\sigma^2} \right) \right)$$

and so the log-likelihood function can be computed as

$$l(\beta_0, \beta_1, \sigma^2) = -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 x_i)^2$$

2.3.1 Likelihood equations

Definition 2.2 (Maximum Likelihood Estimator).

$$(\hat{\beta}_0^{MLE}, \hat{\beta}_1^{MLE}, \hat{\sigma}^{2MLE}) := \arg \sup L(\beta_0, \beta_1, \sigma^2)$$

Since $l(\cdot) = \ln L(\cdot)$ is monotone,

Remark.

$$(\hat{\beta}_0^{MLE}, \hat{\beta}_1^{MLE}, \hat{\sigma}^{2MLE}) = \arg \sup l(\beta_0, \beta_1, \sigma^2)$$

We can find the maximum of this *quadratic* function by making first derivative.

$$\frac{\partial l}{\partial \beta_0} = \frac{1}{\sigma^2} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 x_i) = 0 \quad (2.7)$$

$$\frac{\partial l}{\partial \beta_1} = \frac{1}{\sigma^2} \sum_{i=1}^n x_i (Y_i - \beta_0 - \beta_1 x_i) = 0 \quad (2.8)$$

$$\frac{\partial l}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 x_i)^2 = 0 \quad (2.9)$$

Denote that Equations (2.7) and (2.8) given $\hat{\sigma}^2$ are equivalent to the normal equations. Thus,

$$\hat{\beta}_0^{MLE} = \hat{\beta}_0^{LSE}, \quad \hat{\beta}_1^{MLE} = \hat{\beta}_1^{LSE}$$

From (2.9),

$$\hat{\sigma}^{2MLE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 x_i)^2 = \frac{n-2}{n} \hat{\sigma}^{2LSE}$$

Recall that $\hat{\sigma}^{2LSE}$ is an unbiased, i.e. this *MLE is not an unbiased estimator*. Since $\hat{\sigma}^{2MLE} \approx \hat{\sigma}^{2LSE}$ for large n , however, it is *asymptotically unbiased*.

Theorem 2.2 (Rao-Cramer Lower Bound, univariate case). *Let $X_1, \dots, X_n \stackrel{iid}{\sim} f(x; \theta)$. If $\hat{\theta}$ is an unbiased estimator of θ ,*

$$\text{Var}(\hat{\theta}) \geq \frac{1}{I_n(\theta)}$$

$$\text{where } I_n(\theta) = -E\left(\frac{\partial^2 l(\theta)}{\partial \theta^2}\right)$$

To apply this theorem @(\thm:rclb) in the simple linear regression setting, i.e. (β_0, β_1) , we need to look at the *bivariate case*.

Theorem 2.3 (Rao-Cramer Lower Bound, bivariate case). *Let $X_1, \dots, X_n \stackrel{iid}{\sim} f(x; \theta_1, \theta_2)$ and let $\boldsymbol{\theta} = (\theta_1, \theta_2)^T$. If each $\hat{\theta}_1, \hat{\theta}_2$ is an unbiased estimator of θ_1 and θ_2 , then*

$$\text{Var}(\boldsymbol{\theta}) := \begin{bmatrix} \text{Var}(\hat{\theta}_1) & \text{Cov}(\hat{\theta}_1, \hat{\theta}_2) \\ \text{Cov}(\hat{\theta}_1, \hat{\theta}_2) & \text{Var}(\hat{\theta}_2) \end{bmatrix} \geq I_n^{-1}(\theta_1, \theta_2)$$

where

$$I_n(\theta_1, \theta_2) = - \begin{bmatrix} E\left(\frac{\partial^2 l(\theta_1, \theta_2)}{\partial \theta_1^2}\right) & E\left(\frac{\partial^2 l(\theta_1, \theta_2)}{\partial \theta_1 \partial \theta_2}\right) \\ E\left(\frac{\partial^2 l(\theta_1, \theta_2)}{\partial \theta_1 \partial \theta_2}\right) & E\left(\frac{\partial^2 l(\theta_1, \theta_2)}{\partial \theta_2^2}\right) \end{bmatrix}$$

References