

3.5 Maximum likelihood estimation for β and σ^2 when $\varepsilon_{n \times 1} \sim N(0, \sigma^2 \mathbf{I})$

When the error distribution is assumed to be known, namely, $\varepsilon_{n \times 1} \sim N(0, \sigma^2 \mathbf{I})$ and σ^2 is unknown.

The likelihood function is

$$L(\beta, \sigma^2) = \left(\frac{1}{2\pi\sigma^2} \right)^{\frac{n}{2}} e^{-\frac{1}{2\sigma^2} \underbrace{(\mathbf{Y} - \mathbf{X}\beta)^\top (\mathbf{Y} - \mathbf{X}\beta)}_{\vec{\varepsilon}}}.$$

The log-likelihood function is

$$\begin{aligned} \log L(\beta, \sigma^2) &= -\frac{n}{2} \log 2\pi - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (\mathbf{Y} - \mathbf{X}\beta)^\top (\mathbf{Y} - \mathbf{X}\beta) \\ &= -\frac{n}{2} \log 2\pi - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (\mathbf{Y}^\top \mathbf{Y} - 2\mathbf{Y}^\top \mathbf{X}\beta + \beta^\top \mathbf{X}^\top \mathbf{X}\beta); \end{aligned}$$

score function $\left\{ \begin{array}{l} \frac{\partial \log L}{\partial \beta} = \frac{1}{\sigma^2} (\mathbf{X}^\top \mathbf{Y} - \mathbf{X}^\top \mathbf{X}\beta), \\ \frac{\partial \log L}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} (\mathbf{Y} - \mathbf{X}\beta)^\top (\mathbf{Y} - \mathbf{X}\beta). \end{array} \right.$ *derive it by yourself* \square

Setting the two equations to zero, we obtain the MLE of β and σ^2 :

$$\text{MLE}(\beta) = \tilde{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}.$$

$$\begin{aligned} \text{MLE}(\sigma^2) = \tilde{\sigma}^2 &= \frac{1}{n} (\mathbf{Y} - \mathbf{X}\tilde{\beta})^\top (\mathbf{Y} - \mathbf{X}\tilde{\beta}) \\ &= \frac{1}{n} (\mathbf{Y} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y})^\top (\mathbf{Y} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}) \\ &= \frac{1}{n} \mathbf{Y}^\top (\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \mathbf{Y} \\ &= \frac{1}{n} \mathbf{Y}^\top (\mathbf{I} - \mathbf{H}) \mathbf{Y} \\ &= \frac{1}{n} [\mathbf{Y}^\top \mathbf{Y} - \tilde{\beta}^\top \mathbf{X}^\top \mathbf{Y}]. \end{aligned}$$

It can be seen that when $\varepsilon_{n \times 1} \sim N(0, \sigma^2 \mathbf{I})$, the least square estimate is the MLE. In other words, when $\varepsilon_{n \times 1} \sim N(0, \sigma^2 \mathbf{I})$, minimizing the least-square objective function is equivalent to maximizing the likelihood function.

3.5.1 Properties of the MLE.

There are a number of properties of $\tilde{\beta}$ and $\tilde{\sigma}^2$.

1. The MLE of β , $\tilde{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$ follows $\underline{N(\beta, (\mathbf{X}^\top \mathbf{X})^{-1} \sigma^2)}$. This is because $\tilde{\beta}$ is linear in \mathbf{Y}

and

$$\begin{aligned}
 E(\hat{\beta}) &= E[(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}] \\
 &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top E(\mathbf{Y}) \\
 &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \beta_0 \\
 &= \beta_0, \\
 \text{Var}(\hat{\beta}) &= \text{Var}((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}) \\
 &= (\mathbf{X}^\top \mathbf{X})^{-1} \sigma^2.
 \end{aligned}$$

2. The sum of squares of the deviations of the observed Y_i 's from their estimated expected values is usually known as the *residual error sum of squares* or *sum of squares due to error*, denoted as SSE. It is given by

$$SSE = (\mathbf{Y} - \hat{\mathbf{Y}})^\top (\mathbf{Y} - \hat{\mathbf{Y}}) = \mathbf{Y}^\top (\mathbf{I} - \mathbf{H}) \mathbf{Y}.$$

Then, $\tilde{\sigma}^2 = \frac{SSE}{n}$ is the MLE. But we have shown that $\hat{\sigma}^2 = \frac{SSE}{n-r(\mathbf{X})}$ is an unbiased estimator of σ^2 with $r(\mathbf{X}) = \text{rank of } \mathbf{X}$, implying $\tilde{\sigma}^2$ is biased for finite n . However, it is consistent or asymptotically unbiased as $n \rightarrow \infty$.

3. The MLE $\hat{\beta}$ and SSE are independent. To show this, recall that $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$ and $SSE = \mathbf{Y}^\top (\mathbf{I} - \mathbf{H}) \mathbf{Y}$ where $\mathbf{Y} \sim N(\mathbf{X}\beta, \sigma^2 \mathbf{I})$. Since $(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\sigma^2 \mathbf{I}) (\mathbf{I} - \mathbf{H}) = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{I} - \mathbf{H}) = 0$, $\hat{\beta}$ and SSE are independent and thus $\hat{\beta}$ and $\hat{\sigma}^2$ are independent.

4. Distribution of $\hat{\sigma}^2$. Consider $\frac{SSE}{\sigma^2} = \frac{1}{\sigma^2} \mathbf{Y}^\top (\mathbf{I} - \mathbf{H}) \mathbf{Y}$ and $\mathbf{Y} \sim N(\mathbf{X}\beta, \mathbf{I}\sigma^2)$. Since

$$\left(\frac{\mathbf{I} - \mathbf{H}}{\sigma^2}\right)(\mathbf{I}\sigma^2) = \mathbf{I} - \mathbf{H}$$

which is idempotent, then

$$\frac{SSE}{\sigma^2} \sim \chi^2_{r(\frac{\mathbf{I}-\mathbf{H}}{\sigma^2}), \frac{1}{2}(\mathbf{X}\beta)^\top (\frac{\mathbf{I}-\mathbf{H}}{\sigma^2})(\mathbf{X}\beta)}.$$

However,

$$(\mathbf{X}\beta)^\top \left(\frac{\mathbf{I} - \mathbf{H}}{\sigma^2}\right)(\mathbf{X}\beta) = 0.$$

Therefore, the noncentrality parameter is 0. In addition,

$$r\left(\frac{\mathbf{I} - \mathbf{H}}{\sigma^2}\right) = r(\mathbf{I} - \mathbf{H}) = n - r(\mathbf{X}).$$

Consequently,

$$\begin{aligned}
 \frac{SSE}{\sigma^2} &\sim \chi^2_{n-r(\mathbf{X})}, \quad \text{or} \\
 \frac{[n - r(\mathbf{X})]\hat{\sigma}^2}{\sigma^2} &\sim \chi^2_{n-r(\mathbf{X})}.
 \end{aligned}$$

Remark 1. When the error distribution is unknown, that is the density function $f_{\varepsilon}(\cdot)$ is unknown, the linear model with unknown error distribution is a semiparametric model. In semiparametric linear regression or accelerated failure time models, complications in efficient estimation arise from the multiple roots of the efficient score and density estimation. The maximum likelihood estimation or the semiparametric efficient estimation of linear models were studied by Zeng and Lin (2007, JASA) and Lin and Chen (2013, Biometrika).

Example 1. Consider $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$, $i = 1, 2, \dots, n$, where $\varepsilon_i \sim N(0, \sigma^2)$.

In matrix notation, $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, where $\mathbf{Y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$ and $\mathbf{X}^T \mathbf{Y} = \begin{pmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_i y_i \end{pmatrix}$.

$$\begin{aligned} (\mathbf{X}^T \mathbf{X})^{-1} &= \frac{1}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \begin{pmatrix} \sum_{i=1}^n x_i^2 & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & n \end{pmatrix} \\ &= \frac{1}{n \sum_{i=1}^n (x_i - \bar{x})^2} \begin{pmatrix} \sum_{i=1}^n x_i^2 & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & n \end{pmatrix}. \end{aligned}$$

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= \frac{1}{n \sum_{i=1}^n (x_i - \bar{x})^2} \begin{pmatrix} \sum_{i=1}^n y_i \sum_{i=1}^n x_i^2 - \sum_{i=1}^n x_i \sum_{i=1}^n x_i y_i \\ n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n y_i \sum_{i=1}^n x_i \end{pmatrix} \\ &= \begin{pmatrix} \bar{y} - \hat{\beta}_1 \bar{x} \\ \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \end{pmatrix} \\ &= \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix}. \end{aligned}$$

$$\begin{aligned}
\hat{\sigma}^2 &= \frac{1}{n-2}(\mathbf{Y}^\top \mathbf{Y} - \hat{\beta}^\top \mathbf{X}^\top \mathbf{Y}) \\
&= \frac{1}{n-2}[\sum_{i=1}^n y_i^2 - \hat{\beta}_0 \sum_{i=1}^n y_i - \hat{\beta}_1 \sum_{i=1}^n x_i y_i] \\
&= \frac{1}{n-2}[\sum_{i=1}^n y_i^2 - \bar{y} \sum_{i=1}^n y_i + \hat{\beta}_1 \bar{x} \sum_{i=1}^n y_i - \hat{\beta}_1 \sum_{i=1}^n x_i y_i] \\
&= \frac{1}{n-2}[\sum_{i=1}^n (y_i - \bar{y})^2 - \hat{\beta}_1 (\sum_{i=1}^n x_i y_i - n \sum_{i=1}^n y_i \sum_{i=1}^n x_i)] \\
&= \frac{1}{n-2}[\sum_{i=1}^n (y_i - \bar{y})^2 - \frac{[\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})]^2}{\sum_{i=1}^n (x_i - \bar{x})^2}].
\end{aligned}$$

Example 1(cont'd)

	parameter	UMVU Estimator
1)	β_1	$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$
2)	β_0	$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$
3)	σ^2	$\hat{\sigma}^2$
4)	$2\beta_1 - 3\beta_0$	$2\hat{\beta}_1 - 3\hat{\beta}_0$
5)	$5\sigma^2 + 8\beta_1$	$5\hat{\sigma}^2 + 8\hat{\beta}_1$

6) $\beta_0 + 1.94\sigma$

Since

$$\begin{aligned}
E\left[\frac{(n-2)\hat{\sigma}^2}{\sigma^2}\right] &= n-2 \quad \text{where } \frac{(n-2)\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-2}^2, \\
\Rightarrow E\left[\frac{\sqrt{n-2}\hat{\sigma}}{\sigma}\right] &= \frac{\sqrt{2}\Gamma(\frac{n-2}{2} + \frac{1}{2})}{\Gamma(\frac{n-2}{2})} \\
\Rightarrow E\left[\frac{\sqrt{n-2}\Gamma(\frac{n-2}{2})}{\sqrt{2}\Gamma(\frac{n-1}{2})}\hat{\sigma}\right] &= \sigma \\
\Rightarrow \text{The UMVUE is } \hat{\beta}_0 + 1.94\left[\frac{\sqrt{n-2}\Gamma(\frac{n-2}{2})}{\sqrt{2}\Gamma(\frac{n-1}{2})}\hat{\sigma}\right].
\end{aligned}$$

7) $\frac{\beta_0}{\sigma^2}$

Since $\hat{\beta}_0$ and $\hat{\sigma}^2$ are independent and

$$\begin{aligned}
E\left[\frac{\sigma^2}{(n-2)\hat{\sigma}^2}\right] &= \frac{\Gamma(\frac{n-2}{2} - 1)}{\Gamma(\frac{n-2}{2})} 2^{-1} = \frac{1}{n-4} \\
\Rightarrow E\left[\frac{n-4}{(n-2)\hat{\sigma}^2}\right] &= \frac{1}{\sigma^2} \\
\Rightarrow \text{The UMVUE is } \hat{\beta}_0 \left(\frac{n-4}{(n-2)\hat{\sigma}^2}\right).
\end{aligned}$$

3.5.2 Deviations from Means

The following lemma is useful to find the inverse matrix of a partitioned full rank symmetric matrix.

Lemma 1.

If

Don't spend time memorizing such complex formulas.

$$\begin{aligned} M &= \begin{bmatrix} \mathbf{X}^\top \\ \mathbf{Z}^\top \end{bmatrix} [\mathbf{X} \quad \mathbf{Z}] \\ &= \begin{bmatrix} \mathbf{X}^\top \mathbf{X} & \mathbf{X}^\top \mathbf{Z} \\ \mathbf{Z}^\top \mathbf{X} & \mathbf{Z}^\top \mathbf{Z} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^\top & \mathbf{D} \end{bmatrix}, \end{aligned}$$

and put

$$\begin{aligned} \mathbf{W} &= (\mathbf{D} - \mathbf{B}^\top \mathbf{A}^{-1} \mathbf{B})^{-1} \\ &= [\mathbf{Z}^\top \mathbf{Z} - \mathbf{Z}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Z}]^{-1}, \end{aligned}$$

then,

$$\begin{aligned} \mathbf{M}^{-1} &= \begin{bmatrix} \mathbf{A}^{-1} + \mathbf{A}^{-1} \mathbf{B} \mathbf{W} \mathbf{B}^\top \mathbf{A}^{-1} & -\mathbf{A}^{-1} \mathbf{B} \mathbf{W} \\ \mathbf{W} \mathbf{B}^\top \mathbf{A}^{-1} & \mathbf{W} \end{bmatrix} \leftarrow \text{often used formula.} \\ &= \begin{bmatrix} -\mathbf{A}^{-1} \mathbf{B} \\ \mathbf{I} \end{bmatrix} \mathbf{W} [-\mathbf{B}^\top \mathbf{A}^{-1} \quad \mathbf{I}] + \begin{bmatrix} \mathbf{A}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}. \end{aligned}$$

In this section, we still consider the linear model in matrix form

$$\mathbf{Y}_{n \times 1} = \mathbf{X} \boldsymbol{\beta}_{(k+1) \times 1} + \underbrace{\boldsymbol{\varepsilon}_{n \times 1}}_{\sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)}.$$

Recall that the least square estimate $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$ and $\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^\top \mathbf{X})^{-1} \sigma^2$.

Write $\mathbf{X} = (\mathbf{1}, \quad \mathbf{X}_1)$ and $\boldsymbol{\beta}^\top = (\beta_0, \mathbf{b}^\top)$, where $\mathbf{X}_1 = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1k} \\ x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \cdots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nk} \end{bmatrix}$ and $\mathbf{b}^\top = (\beta_1, \beta_2, \dots, \beta_k)$.

Let $\bar{\mathbf{X}}^\top = (\bar{x}_{.1}, \bar{x}_{.2}, \dots, \bar{x}_{.k})$ with $\bar{x}_{.i} = (1/n) \sum_{j=1}^n x_{j,i}$, $i = 1, \dots, k$.

Note that

$$\begin{aligned} \mathbf{1}^\top \mathbf{1} &= n \\ \mathbf{1}^\top \mathbf{Y} &= n\bar{y} \\ \mathbf{1}^\top \mathbf{X}_1 &= n\bar{\mathbf{X}}^\top, \end{aligned}$$

where $\mathbf{1}_{n \times 1} = (1, 1, \dots, 1)^\top$. We then rewrite $\hat{\beta}$ as follows:

$$\begin{aligned}\hat{\beta} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} \\ &= \left[\begin{pmatrix} \mathbf{1}^\top \\ \mathbf{X}_1^\top \end{pmatrix} (\mathbf{1} \quad \mathbf{X}_1) \right]^{-1} \begin{pmatrix} \mathbf{1}^\top \\ \mathbf{X}_1^\top \end{pmatrix} \mathbf{Y} \\ &= \begin{bmatrix} n & n\bar{\mathbf{X}}^\top \\ n\bar{\mathbf{X}} & \mathbf{X}_1^\top \mathbf{X}_1 \end{bmatrix}^{-1} \begin{bmatrix} n\bar{y} \\ \mathbf{X}_1^\top \mathbf{Y} \end{bmatrix}.\end{aligned}$$

By Lemma 1, we have

$$\hat{\beta} = \begin{bmatrix} \frac{1}{n} + \bar{\mathbf{X}}^\top \mathbf{S}^{-1} \bar{\mathbf{X}} & -\bar{\mathbf{X}}^\top \mathbf{S}^{-1} \\ -\mathbf{S}^{-1} \bar{\mathbf{X}} & \mathbf{S}^{-1} \end{bmatrix} \begin{bmatrix} n\bar{y} \\ \mathbf{X}_1^\top \mathbf{Y} \end{bmatrix},$$

where $\mathbf{S} = \mathbf{X}_1^\top \mathbf{X}_1 - n\bar{\mathbf{X}}\bar{\mathbf{X}}^\top = \mathbf{Z}^\top \mathbf{Z}$ and $\mathbf{Z} = \mathbf{X}_1 - \mathbf{1}\bar{\mathbf{X}}^\top$. This implies

$$\begin{aligned}\Rightarrow \begin{bmatrix} \hat{\beta}_0 \\ \hat{\mathbf{b}} \end{bmatrix} &= \begin{bmatrix} \bar{y} - \bar{\mathbf{X}}^\top \mathbf{S}^{-1} (\mathbf{X}_1^\top \mathbf{Y} - n\bar{y}\bar{\mathbf{X}}) \\ \mathbf{S}^{-1} (\mathbf{X}_1^\top \mathbf{Y} - n\bar{y}\bar{\mathbf{X}}) \end{bmatrix} \hat{\mathbf{b}} \\ \Rightarrow \hat{\beta}_0 &= \bar{y} - \bar{\mathbf{X}}^\top \hat{\mathbf{b}} = \frac{1}{n} \mathbf{1}^\top \mathbf{Y} - \bar{\mathbf{X}}^\top \hat{\mathbf{b}} \\ \hat{\mathbf{b}} &= \mathbf{S}^{-1} (\mathbf{X}_1^\top \mathbf{Y} - n\bar{y}\bar{\mathbf{X}}) \\ &= (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{Y}.\end{aligned}$$

Similarly, it follows directly that

$$\begin{aligned}\text{var}(\hat{\beta}) &= \text{var} \begin{pmatrix} \hat{\beta}_0 \\ \hat{\mathbf{b}} \end{pmatrix} = (\mathbf{X}^\top \mathbf{X})^{-1} \sigma^2 \\ &= \begin{bmatrix} \frac{1}{n} + \bar{\mathbf{X}}^\top (\mathbf{Z}^\top \mathbf{Z})^{-1} \bar{\mathbf{X}} & -\bar{\mathbf{X}}^\top (\mathbf{Z}^\top \mathbf{Z})^{-1} \\ -(\mathbf{Z}^\top \mathbf{Z})^{-1} \bar{\mathbf{X}} & (\mathbf{Z}^\top \mathbf{Z})^{-1} \end{bmatrix} \sigma^2, \\ \Rightarrow \text{var}(\hat{\mathbf{b}}) &= (\mathbf{Z}^\top \mathbf{Z})^{-1} \sigma^2,\end{aligned}$$

$$\begin{aligned}\text{var}(\hat{\beta}_0) &= \frac{\sigma^2}{n} + \bar{\mathbf{X}}^\top (\mathbf{Z}^\top \mathbf{Z})^{-1} \bar{\mathbf{X}} \sigma^2 \\ &= \frac{\sigma^2}{n} + \bar{\mathbf{X}}^\top \text{var}(\hat{\mathbf{b}}) \bar{\mathbf{X}},\end{aligned}$$

$$\begin{aligned}\text{cov}(\hat{\beta}_0, \hat{\mathbf{b}}^\top) &= -\bar{\mathbf{X}}^\top (\mathbf{Z}^\top \mathbf{Z})^{-1} \sigma^2 \\ &= -\bar{\mathbf{X}}^\top \text{var}(\hat{\mathbf{b}}).\end{aligned}$$

intercept & coefficients
are not independent! Take care.