## 8.1 Example: Mixed Estimation (Theil and Goldberger Model)

A generalization of the restricted OLS $\implies$ Stochastic linear restriction:

$$r = R\beta + v, \qquad \text{E}(v) = 0 \text{ and } \text{V}(v) = \sigma^2\Psi$$

$$y = X\beta + u, \qquad \text{E}(u) = 0 \text{ and } \text{V}(u) = \sigma^2 I_n$$

Using a matrix form,

$$\begin{pmatrix} y \\ r \end{pmatrix} = \begin{pmatrix} X \\ R \end{pmatrix}\beta + \begin{pmatrix} u \\ v \end{pmatrix}, \qquad \text{E}\begin{pmatrix} u \\ v \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \text{ and } \text{V}\begin{pmatrix} u \\ v \end{pmatrix} = \sigma^2 \begin{pmatrix} I_n & 0 \\ 0 & \Psi \end{pmatrix}$$

For estimation, we do not need normality assumption.

Applying GLS, we obtain:

$$
b = \left( (X' \quad R') \begin{pmatrix} I_n & 0 \\ 0 & \Psi \end{pmatrix}^{-1} \begin{pmatrix} X \\ R \end{pmatrix} \right)^{-1} \left( (X' \quad R') \begin{pmatrix} I_n & 0 \\ 0 & \Psi \end{pmatrix}^{-1} \begin{pmatrix} y \\ r \end{pmatrix} \right)
$$

$$
= \left( X'X + R'\Psi^{-1}R \right)^{-1} \left( X'y + R'\Psi^{-1}r \right).
$$

127

Mean and Variance of $b$:        $b$ is rewritten as follows:

$$b = \left( (X' \quad R') \begin{pmatrix} I_n & 0 \\ 0 & \Psi \end{pmatrix}^{-1} \begin{pmatrix} X \\ R \end{pmatrix} \right)^{-1} \left( (X' \quad R') \begin{pmatrix} I_n & 0 \\ 0 & \Psi \end{pmatrix}^{-1} \begin{pmatrix} y \\ r \end{pmatrix} \right)$$

$$= \beta + \left( (X' \quad R') \begin{pmatrix} I_n & 0 \\ 0 & \Psi \end{pmatrix}^{-1} \begin{pmatrix} X \\ R \end{pmatrix} \right)^{-1} \begin{pmatrix} u \\ v \end{pmatrix}$$

Therefore, the mean and variance are given by:

$$\mathrm{E}(b) = \beta \qquad \Longrightarrow \qquad b \text{ is unbiased.}$$

$$\mathrm{V}(b) = \sigma^2 \left( (X' \quad R') \begin{pmatrix} I_n & 0 \\ 0 & \Psi \end{pmatrix}^{-1} \begin{pmatrix} X \\ R \end{pmatrix} \right)^{-1}$$

$$= \sigma^2 \left( X'X + R'\Psi^{-1}R \right)^{-1}$$

# 9  Maximum Likelihood Estimation (MLE,　　　)

    $\longrightarrow$ **Review**

1. The distribution function of $\{X_i\}_{i=1}^n$ is $f(x; \theta)$, where $x = (x_1, x_2, \cdots, x_n)$.

   $\theta$ is a vector or matrix of unknown parameters, e.g., $\theta = (\mu, \Sigma)$, where $\mu = E(X_i)$ and $\Sigma = V(X_i)$.

   Note that $X$ is a vector of random variables and $x$ is a vector of their realizations (i.e., observed data).

   Likelihood function $L(\cdot)$ is defined as $L(\theta; x) = f(x; \theta)$.

   Note that $f(x; \theta) = \prod_{i=1}^n f(x_i; \theta)$ when $X_1, X_2, \cdots, X_n$ are mutually independently and identically distributed.

The maximum likelihood estimate (MLE) of $\theta$ is the $\theta$ such that:

$$\max_{\theta} \ L(\theta; x). \qquad \Longleftrightarrow \qquad \max_{\theta} \ \log L(\theta; x).$$

Thus, MLE satisfies the following two conditions:

(a) $\dfrac{\partial \log L(\theta; x)}{\partial \theta} = 0.$ $\implies$ Solution of $\theta$: $\tilde{\theta} = \tilde{\theta}(x)$

(b) $\dfrac{\partial^2 \log L(\theta; x)}{\partial \theta \partial \theta'}$ is a negative definite matrix.

2. $x = (x_1, x_2, \cdots, x_n)$ are used as the observations (i.e., observed data).

$X = (X_1, X_2, \cdots, X_n)$ denote the random variables associated with the joint distribution $f(x; \theta) = \prod_{i=1}^{n} f(x_i; \theta)$.

3. Replacing *x* by *X*, we otain the maximum likelihood **estimator** (MLE, which is the same word as the maximum likelihood **estimate**).

   That is, MLE of $\theta$ satisfies the following two conditions:

   (a) $\dfrac{\partial \log L(\theta; X)}{\partial \theta} = 0.$ $\implies$ Solution of $\theta$: $\tilde{\theta} = \tilde{\theta}(X)$

   (b) $\dfrac{\partial^2 \log L(\theta; X)}{\partial \theta \partial \theta'}$ is a negative definite matrix.

4. **Fisher's information matrix (　　　　　　　　　　) or simply information matrix**, denoted by $I(\theta)$, is given by:

$$I(\theta) = -\mathrm{E}\Big(\frac{\partial^2 \log L(\theta; X)}{\partial \theta \partial \theta'}\Big),$$

   where we have the following equality:

$$-\mathrm{E}\Big(\frac{\partial^2 \log L(\theta; X)}{\partial \theta \partial \theta'}\Big) = \mathrm{E}\Big(\frac{\partial \log L(\theta; X)}{\partial \theta} \frac{\partial \log L(\theta; X)}{\partial \theta'}\Big) = \mathrm{V}\Big(\frac{\partial \log L(\theta; X)}{\partial \theta}\Big)$$

   Note that $\mathrm{E}(\cdot)$ and $\mathrm{V}(\cdot)$ are expected with respect to *X*.

**Proof of the above equality:**

$$\int L(\theta; x)\mathrm{d}x = 1$$

Take a derivative with respect to $\theta$.

$$\int \frac{\partial L(\theta; x)}{\partial \theta}\mathrm{d}x = 0$$

(We assume that (i) the domain of $x$ does not depend on $\theta$ and (ii) the derivative $\frac{\partial L(\theta; x)}{\partial \theta}$ exists.)

Rewriting the above equation, we obtain:

$$\int \frac{\partial \log L(\theta; x)}{\partial \theta} L(\theta; x)\mathrm{d}x = 0,$$

i.e.,

$$\mathrm{E}\left(\frac{\partial \log L(\theta; X)}{\partial \theta}\right) = 0.$$

132

Again, differentiating the above with respect to $\theta$, we obtain:

$$
\int \frac{\partial^2 \log L(\theta; x)}{\partial\theta\partial\theta'} L(\theta; x)\mathrm{d}x + \int \frac{\partial \log L(\theta; x)}{\partial\theta} \frac{\partial L(\theta; x)}{\partial'\theta}\mathrm{d}x
$$
$$
= \int \frac{\partial^2 \log L(\theta; x)}{\partial\theta\partial\theta'} L(\theta; x)\mathrm{d}x + \int \frac{\partial \log L(\theta; x)}{\partial\theta} \frac{\partial \log L(\theta; x)}{\partial\theta'} L(\theta; x)\mathrm{d}x
$$
$$
= \mathrm{E}\Big(\frac{\partial^2 \log L(\theta; X)}{\partial\theta\partial\theta'}\Big) + \mathrm{E}\Big(\frac{\partial \log L(\theta; X)}{\partial\theta} \frac{\partial \log L(\theta; X)}{\partial\theta'}\Big) = 0.
$$

Therefore, we can derive the following equality:

$$
-\mathrm{E}\left(\frac{\partial^2 \log L(\theta; X)}{\partial\theta\partial\theta'}\right) = \mathrm{E}\left(\frac{\partial \log L(\theta; X)}{\partial\theta} \frac{\partial \log L(\theta; X)}{\partial\theta'}\right) = \mathrm{V}\left(\frac{\partial \log L(\theta; X)}{\partial\theta}\right),
$$

where the second equality utilizes $\mathrm{E}\left(\dfrac{\partial \log L(\theta; X)}{\partial\theta}\right) = 0$.

5. **Cramer-Rao Lower Bound (** ) is given by: $(I(\theta))^{-1}$.

Suppose that an estimator of $\theta$ is given by $s(X)$.

The expectation of $s(X)$ is:

$$\mathrm{E}(s(X)) = \int s(x)L(\theta; x)\mathrm{d}x.$$

Differentiating the above with respect to $\theta$,

$$\frac{\partial \mathrm{E}(s(X))}{\partial \theta} = \int s(x)\frac{\partial L(\theta; x)}{\partial \theta}\mathrm{d}x = \int s(x)\frac{\partial \log L(\theta; x)}{\partial \theta}L(\theta; x)\mathrm{d}x$$
$$= \mathrm{Cov}\left(s(X), \frac{\partial \log L(\theta; X)}{\partial \theta}\right)$$

For simplicity, let $s(X)$ and $\theta$ be scalars.

Then,

$$\left(\frac{\partial \mathrm{E}(s(X))}{\partial \theta}\right)^2 = \left(\mathrm{Cov}\left(s(X), \frac{\partial \log L(\theta; X)}{\partial \theta}\right)\right)^2 = \rho^2 \mathrm{V}\left(s(X)\right) \mathrm{V}\left(\frac{\partial \log L(\theta; X)}{\partial \theta}\right)$$
$$\leq \mathrm{V}\left(s(X)\right) \mathrm{V}\left(\frac{\partial \log L(\theta; X)}{\partial \theta}\right),$$

134

where $\rho$ denotes the correlation coefficient between $s(X)$ and $\dfrac{\partial \log L(\theta; X)}{\partial \theta}$, i.e.,

$$\rho = \frac{\text{Cov}\left(s(X), \dfrac{\partial \log L(\theta; X)}{\partial \theta}\right)}{\sqrt{\text{V}\left(s(X)\right)}\sqrt{\text{V}\left(\dfrac{\partial \log L(\theta; X)}{\partial \theta}\right)}}.$$

Note that $|\rho| \leq 1$.

Therefore, we have the following inequality:

$$\left(\frac{\partial \text{E}(s(X))}{\partial \theta}\right)^2 \leq \text{V}(s(X))\, \text{V}\left(\frac{\partial \log L(\theta; X)}{\partial \theta}\right),$$

i.e.,

$$\text{V}(s(X)) \geq \frac{\left(\dfrac{\partial \text{E}(s(X))}{\partial \theta}\right)^2}{\text{V}\left(\dfrac{\partial \log L(\theta; X)}{\partial \theta}\right)}$$

135

Especially, when $E(s(X)) = \theta$, i.e., when $s(X)$ is an unbiased estimator of $\theta$, the numerator of the right-hand side leads to one.

Therefore, we obtain:

$$V(s(X)) \geq \frac{1}{-E\left(\dfrac{\partial^2 \log L(\theta; X)}{\partial \theta^2}\right)} = (I(\theta))^{-1}.$$

Even in the case where $s(X)$ is a vector, the following inequality holds.

$$V(s(X)) \geq (I(\theta))^{-1},$$

where $I(\theta)$ is defined as:

$$\begin{aligned}
I(\theta) &= -E\left(\frac{\partial^2 \log L(\theta; X)}{\partial \theta \partial \theta'}\right) \\
&= E\left(\frac{\partial \log L(\theta; X)}{\partial \theta} \frac{\partial \log L(\theta; X)}{\partial \theta'}\right) = V\left(\frac{\partial \log L(\theta; X)}{\partial \theta}\right).
\end{aligned}$$

The variance of any unbiased estimator of $\theta$ is larger than or equal to $(I(\theta))^{-1}$.

Thus, $(I(\theta))^{-1}$ results in the lower bound of the variance of any unbiased estimator of $\theta$.

6. Asymptotic Normality of MLE:

Let $\tilde{\theta}$ be MLE of $\theta$.

As $n$ goes to infinity, we have the following result:

$$\sqrt{n}(\tilde{\theta} - \theta) \longrightarrow N\left(0, \lim_{n\to\infty}\left(\frac{I(\theta)}{n}\right)^{-1}\right),$$

where it is assumed that $\lim_{n\to\infty}\left(\dfrac{I(\theta)}{n}\right)$ converges.

$\longrightarrow$ The proof will be shown later.

That is, when $n$ is large, $\tilde{\theta}$ is approximately distributed as follows:

$$\tilde{\theta} \sim N\left(\theta, (I(\theta))^{-1}\right).$$

Suppose that $s(X) = \tilde{\theta}$.

When $n$ is large, $V(s(X))$ is approximately equal to $(I(\theta))^{-1}$.

7. **Optimization (      ):**

   MLE of $\theta$ results in the following maximization problem:

   $$\max_{\theta} \quad \log L(\theta; x).$$

   We often have the case where the solution of $\theta$ is not derived in closed form.

   $\Longrightarrow$ Optimization procedure

   $$0 = \frac{\partial \log L(\theta; x)}{\partial \theta} = \frac{\partial \log L(\theta^*; x)}{\partial \theta} + \frac{\partial^2 \log L(\theta^*; x)}{\partial \theta \partial \theta'}(\theta - \theta^*).$$

138

Solving the above equation with respect to $\theta$, we obtain the following:

$$\theta = \theta^* - \left(\frac{\partial^2 \log L(\theta^*; x)}{\partial\theta\partial\theta'}\right)^{-1} \frac{\partial \log L(\theta^*; x)}{\partial\theta}.$$

Replace the variables as follows:

$$\theta \longrightarrow \theta^{(i+1)}$$

$$\theta^* \longrightarrow \theta^{(i)}$$

Then, we have:

$$\theta^{(i+1)} = \theta^{(i)} - \left(\frac{\partial^2 \log L(\theta^{(i)}; x)}{\partial\theta\partial\theta'}\right)^{-1} \frac{\partial \log L(\theta^{(i)}; x)}{\partial\theta}.$$

$\Longrightarrow$ **Newton-Raphson method (** )

Replacing $\dfrac{\partial^2 \log L(\theta^{(i)}; x)}{\partial\theta\partial\theta'}$ by $E\left(\dfrac{\partial^2 \log L(\theta^{(i)}; x)}{\partial\theta\partial\theta'}\right)$, we obtain the following op-

139

timization algorithm:

$$\theta^{(i+1)} = \theta^{(i)} - \left( \mathrm{E} \left( \frac{\partial^2 \log L(\theta^{(i)}; x)}{\partial \theta \partial \theta'} \right) \right)^{-1} \frac{\partial \log L(\theta^{(i)}; x)}{\partial \theta}$$

$$= \theta^{(i)} + \left( I(\theta^{(i)}) \right)^{-1} \frac{\partial \log L(\theta^{(i)}; x)}{\partial \theta}$$

$\Longrightarrow$ **Method of Scoring (          )**