

8th Lecture : Maximum Likelihood (ML) 15/12/2022 (1)

- ML estimation of the model

$$y = X\beta + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I_n)$$

- Define the likelihood function

$$\begin{aligned} L(\beta, \sigma^2; y) &= \prod_{i=1}^n N(y_i | x_i' \beta, \sigma^2) = \\ &= (2\pi)^{-n/2} (\sigma^2)^{-n/2} e^{-\frac{1}{2\sigma^2} (y - X\beta)'(y - X\beta)} \end{aligned}$$

- Usually work with log-likelihood:

$$\begin{aligned} \ell(\beta, \sigma^2; y) &= \log L(\beta, \sigma^2; y) \\ &= -\frac{n}{2} \log \pi - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (y - X\beta)'(y - X\beta) \end{aligned}$$

- Use FOCs to derive ~~MLEs~~ MLEs

$$\left. \begin{aligned} \frac{\partial \ell}{\partial \beta} &= 0 \\ \frac{\partial \ell}{\partial \sigma^2} &= 0 \end{aligned} \right\} \Rightarrow \begin{aligned} \hat{\beta}_{ML} &= (X'X)^{-1} X'y \\ \hat{\sigma}_{ML}^2 &= \frac{\hat{\varepsilon}'\hat{\varepsilon}}{n} \end{aligned}$$

(Notice that $\hat{\beta}_{ML} \neq \hat{\sigma}_{ML}^2$, i.e., $\hat{\sigma}_{ML}^2$ is biased)

(2)

Score function: The vector of gradients, i.e., for $\theta = (\beta, \sigma^2)$ the score is the vector $\frac{\partial \log L}{\partial \theta}$.

Proposition: The expected value of the score function is 0, that is $E\left(\frac{\partial \log L}{\partial \theta}\right) = 0$

Proof

$$\begin{aligned}
E\left(\frac{\partial \log L}{\partial \theta}\right) &= E\left(E\left(\frac{\partial \log L}{\partial \theta} \mid X\right)\right) = \\
&= E\left(E\left(\frac{1}{2\sigma^2} X' \epsilon \mid X\right)\right) = \\
&= E\left(\frac{1}{2\sigma^2} X' E(\epsilon \mid X)\right) = 0
\end{aligned}$$

$$\begin{aligned}
E\left(\frac{\partial \log L}{\partial \sigma^2}\right) &= E\left(E\left(-\frac{n}{2\sigma^2} + \frac{1}{\sigma^4} \sum \epsilon_i^2 \mid X\right)\right) = \\
&= -\frac{n}{2\sigma^2} + \frac{1}{\sigma^4} [E(\sum \epsilon_i^2 \mid X)] = \\
&= -\frac{n}{2\sigma^2} + \frac{n\sigma^2}{\sigma^4} = 0
\end{aligned}$$

Information matrix:

(3)

$$I(\theta) \equiv -E \left(\frac{\partial^2 \log L}{\partial \theta \partial \theta'} \right)$$

→ If θ one dimensional

• For $\theta = (\beta, \sigma^2)$ we expect that $I(\theta)$ is 4×4 matrix.

$$i) \frac{\partial \log L}{\partial \beta \partial \beta'} = \frac{\partial}{\partial \beta'} \left(\frac{\partial \log L}{\partial \beta} \right) =$$

$$I(\theta) = \begin{pmatrix} \frac{\partial^2 \log L}{(\partial \beta_1)^2} & \frac{\partial^2 \log L}{\partial \beta_1 \partial \beta_2} \\ \frac{\partial^2 \log L}{\partial \beta_2 \partial \beta_1} & \frac{\partial^2 \log L}{(\partial \beta_2)^2} \end{pmatrix}$$

$$= \frac{\partial}{\partial \beta'} \left(\frac{1}{\sigma^2} X'(y - X\beta) \right)$$

$$= - \frac{X'X}{\sigma^2} \Rightarrow$$

$$\Rightarrow E \left(\frac{\partial \log L}{\partial \beta \partial \beta'} \right) = E \left(E \left(\frac{X'X}{\sigma^2} \mid X \right) \right)$$

$$= \frac{X'X}{\sigma^2}$$

$$ii) \frac{\partial^2 \log L}{\partial \beta \partial \sigma^2} = \frac{\partial}{\partial \sigma^2} \left(\frac{\partial \log L}{\partial \beta} \right) =$$

$$= \frac{\partial}{\partial \sigma^2} \left(\frac{1}{\sigma^2} X'\epsilon \right) =$$

$$= - \frac{X'\epsilon}{\sigma^4}$$

Therefore,

(4)

$$E\left(-\frac{\partial^2 \log L}{\partial \sigma^2} \mid X\right) = \frac{E(X' \varepsilon \mid X)}{\sigma^4} = \frac{X' E(\varepsilon \mid X)}{\sigma^4} = 0$$

~~(iii)~~

$$\begin{aligned} \text{(iii)} \quad \frac{\partial^2 \log L}{(\partial \sigma^2)^2} &= \frac{\partial}{\partial \sigma^2} \left(-\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \varepsilon' \varepsilon \right) = \\ &= -\frac{n}{2\sigma^4} - \frac{1}{\sigma^6} \varepsilon' \varepsilon \end{aligned}$$

$$\begin{aligned} E\left(-\frac{\partial^2 \log L}{(\partial \sigma^2)^2} \mid X\right) &= -E\left(\frac{n}{2\sigma^4} - \frac{1}{\sigma^6} \varepsilon' \varepsilon \mid X\right) \\ &= -\frac{n}{2\sigma^4} + \frac{1}{\sigma^6} E\left(\sum_{i=1}^n \varepsilon_i^2 \mid X\right) \\ &= -\frac{n}{2\sigma^4} + \frac{1}{\sigma^6} \sum_{i=1}^n E(\varepsilon_i^2 \mid X) \\ &= -\frac{n}{2\sigma^4} + \frac{n\sigma^2}{\sigma^6} = \\ &= -\frac{n}{2\sigma^4} + \frac{n}{\sigma^4} \end{aligned}$$

$$(k+1) \times (k+1) = -\frac{n}{2\sigma^4} + \frac{2n}{2\sigma^4} = \frac{n}{2\sigma^4}$$

$$\text{Thus, } I(\theta) = \begin{pmatrix} \frac{X'X}{\sigma^2} & 0_{k \times 1} \\ 0_{1 \times k} & \frac{n}{2\sigma^4} \end{pmatrix}$$

Variance of Score function:

(5)

$$\text{Var}\left(\frac{\partial \log L}{\partial \theta}\right) = E\left[\left(\frac{\partial \log L}{\partial \theta}\right)\left(\frac{\partial \log L}{\partial \theta}\right)'\right] - \left(\begin{array}{l} \star \\ \text{Var}(x) = \\ = E((x - E x)(x - E x)') \\ = E x x' - \\ - E x (E x)' \end{array}\right)$$

$$\ominus - E\left(\frac{\partial \log L}{\partial \theta}\right) E\left(\frac{\partial \log L}{\partial \theta}\right)'$$

$$\text{So } \text{Var} \frac{\partial \log L}{\partial \theta} = - E\left(\frac{\partial^2 \log L}{\partial \theta \partial \theta'}\right) = I(\theta)$$

Remark: We have that $\text{Var}\left(\frac{\partial \log L}{\partial \theta}\right) =$

$$= - E\left(\frac{\partial^2 \log L}{\partial \theta \partial \theta'}\right)$$

because we assume that $E \sim N(0, \sigma^2 I)$
 To check this assumption we ~~can~~ can
 measure the diff. between (*) and
 (**)

- Cramer - Rao Lower Bound (CRLB) ⑥

Theorem (Statistiks): If $E(\hat{\theta}) = \theta$ then
 $\text{Var}(\hat{\theta}) \geq I(\theta)^{-1}$ and $I(\theta)^{-1}$ is called
 the CRLB.

• Remark 1: $I(\theta)^{-1} = \begin{pmatrix} \sigma^2 (X'X)^{-1} & 0_{k \times 1} \\ 0_{1 \times k} & \frac{2\sigma^2}{n} \end{pmatrix}$

• Remark 2:

Since $\hat{\beta}_{GLS} = \hat{\beta}_{OLS}$ we have that ~~the~~ $\hat{\beta}_{OLS}$
 attains the CRLB. This is stronger
 than the GM theorem (BLUE) since
 we remove the linear assumption
 for the estimator. ~~the~~ This is due
 to the normality assumption.

• Properties of the MLEs (7)

* An estimator $\hat{\theta}$ is ~~called~~ asymptotically efficient if it is consistent and asymptotically normal.

• Let $\hat{\theta}$ the MLE of the unknown θ for a well-behaved likelihood function ~~h~~ $h(y; \theta)$ then:

i) $\hat{\theta}$ is consistent $\text{plim}_{n \rightarrow \infty} \hat{\theta} = \theta$

ii) $\hat{\theta}$ is asymptotically normal

$$\sqrt{n} (\hat{\theta} - \theta) \xrightarrow{d} N\left(0, \frac{I(\theta)}{n}\right)^{n-1}$$

where $I(\theta) = -E\left(\frac{\partial^2 \log L}{\partial \theta \partial \theta'}\right)$

which in practice implies the for large n

$$\hat{\theta} \approx N\left(\theta, \left(\frac{I(\theta)}{n}\right)^{-1}\right)$$

iii) $\hat{\theta}$ attains the CRLB and is asympt. efficient

iv) Invariance of MLE: the MLE of $g(\theta)$ is $\tilde{\theta} = g(\hat{\theta})$.

• Hypothesis testing

(8)

Let $\hat{\theta}$ MLE of θ and we wish to test $H_0: g(\theta) = 0$.

In the ML approach this test is conducted by using the Lik. Ratio test:

We test H_0 if the diff

$\log L_U - \log L_R$ is large or not.

(*The analogous in OLS estimation is the F-test based on SSR_U and $SSR_R \rightarrow$ Wald Test)

~~Alternative approach: Lagrange Multiplier or Rao's Score test:~~

~~If H_0 is true then the restricted estimator should be near to the point that max. the log-likelihood~~

→ Theorem: let $H_0: g(\theta) = C$ vs $H_1: g(\theta) \neq C$ ⁹
~~then~~ and set $J = \frac{L_R}{L_U}$ then

$$-2 \log J = 2(\ln L_U - \ln L_R) \xrightarrow{H_0} \chi^2_{(J)} \text{ number.}$$

(Reminder: We have shown the F-test $\sim \chi^2$ of restrictions)

(~~Reminder~~ Some ~~asymptotic~~ result for
 Wald's test ~~can~~ can be shown
 as well as for the Lagrange test)

Generalized Least Squares

Assume the model

$$y = X\beta + \varepsilon \quad \text{but } \varepsilon \neq 0$$

$$E(\varepsilon\varepsilon' | X) = \sigma^2 V(X) \neq \sigma^2 I_n$$

↓
 symmetric
 and
 pos. definite

• Examples for V

i) $V = \begin{pmatrix} h_1 & & \\ & \ddots & \\ & & h_n \end{pmatrix}$ (i.e. only diff. variance \Rightarrow
 \Rightarrow no autocorrelated errors)

ii) $V = \begin{pmatrix} 1 & \rho & \dots & \rho^{n-1} \\ & \ddots & & \\ & & \ddots & \\ \rho & & & 1 \end{pmatrix} \rightarrow$ ~~no~~ autocorrelated errors.

• Then $\hat{\beta}_{OLS}$ is not BLUE:

$$\text{Since } \text{Var}(\hat{\beta}_{OLS} | X) =$$

$$= E\left((\hat{\beta} - E(\hat{\beta} | X)) (\hat{\beta} - E(\hat{\beta} | X))' \right)$$

$$= E\left((\hat{\beta} - \beta) (\hat{\beta} - \beta)' | X \right)$$

$$= E\left((X'X)^{-1} X' \varepsilon \varepsilon' X (X'X)^{-1} | X \right)$$

$$= \sigma^2 (X'X)^{-1} X' V X (X'X)^{-1}$$

(Notice that for $V = I_n \Rightarrow \text{Var}(\hat{\beta} | X) = \sigma^2 (X'X)^{-1}$)

Conseq. of not being BLUE:

We cannot assume that

$$se(\hat{\beta}_j) = \sqrt{\sigma^2 (X'X)^{-1}_{jj}} \Rightarrow$$

~~The~~ The t-test does not follow the t-distribution and the F-test does not follow the F-distribution. as well as for the Wald test.

→ Solution: By noting the GM theorem 11
 requires spherical errors we
~~to~~ make suitable transformations
 to go back to spherical errors.

In particular, let C an $n \times n$ matrix

s.t. $V^{-1} = C' C$. (e.g. Cholesky)

↳ Remark: $CVC' = I_n$ factor

then make the following (linear) transformation:

transformation:

$$y \rightarrow Cy = \tilde{y}$$

$$X \rightarrow CX = \tilde{X}$$

$$\varepsilon \rightarrow C\varepsilon = \tilde{\varepsilon}$$

$$\begin{aligned} C' C^{-1} V^{-1} &= C \\ C' C^{-1} V^{-1} C^{-1} &= I_n \\ (CVC')^{-1} &= I_n \end{aligned}$$

that is we now work with the model

$$\tilde{y} = \tilde{X}\beta + \tilde{\varepsilon}$$

then, $E(\tilde{\varepsilon}\tilde{\varepsilon}') = \sigma^2 CVC' = \sigma^2 I_n$.

(Also $E(\tilde{\varepsilon}|x) = 0$ still true since $E(\varepsilon|x) = C E(\varepsilon|x) = 0$)

• $\hat{\beta}_{GLS}$ is BLUE since the transformed (12)
model follows GLM conditions

$$\text{and } \hat{\beta}_{GLS} = (\tilde{X}'\tilde{X})^{-1}\tilde{X}'\tilde{y} =$$

$$= (X'V^{-1}X)^{-1}X'V^{-1}y$$

We weight
differently each
observation to account
for the fact that
the variance is not
constant across observations

↳ Easier to understand when

$$V = \begin{pmatrix} h_1 & & \\ & \dots & \\ & & h_n \end{pmatrix}$$

Where $\hat{\beta}_{GLS}$ is known as Weighted Least
Squares estimator

$$\begin{aligned} \bullet \text{Var}(\hat{\beta}_{GLS} | X) &= \sigma^2 (\tilde{X}'\tilde{X})^{-1} \\ &= \sigma^2 (X'V^{-1}X)^{-1} \end{aligned}$$

In case we had time:

- FOCs for GLS
- MLE
- Finite Sample Properties
- Asymptotic Properties
- Hypothesis testing
- Explore different forms of Heteroskedasticity.