

## LECTURE 14 MAXIMUM LIKELIHOOD ESTIMATION

### Definition

Suppose that the econometrician observes the data  $\{W_i : i = 1, \dots, n\}$ , where  $W_i$  is a random  $p$ -vector. Assume that  $W_i$  are iid with the PDF  $f(w_i, \theta)$ , where  $\theta \in \Theta \subset R^k$  is unknown vector of parameters. The set  $\Theta$  is usually assumed to be compact.

**Example** (Normal regression model). Let  $W_i = (Y_i, X_i')'$ ,  $\theta = (\beta', \sigma^2)'$ , where  $\beta \in R^k$  and  $\sigma^2 \in R$ . Assume that  $Y_i = X_i'\beta + U_i$ , and  $U_i|X_i \sim N(0, \sigma^2)$ . Then,  $f(y_i, x_i, \beta, \sigma^2) = (2\pi\sigma^2)^{-1/2} \exp\left(-\frac{(y_i - x_i'\beta)^2}{2\sigma^2}\right)$ .

Since the observations are iid, the joint PDF of  $W_1, \dots, W_n$  is given by

$$\prod_{i=1}^n f(w_i, \theta).$$

The joint PDF gives us the *likelihood* of our sample given the value of  $\theta$ . Log of joint PDF is

$$\sum_{i=1}^n \log f(w_i, \theta).$$

We define the log-likelihood function as  $1/n$  log of joint PDF evaluated at random sample  $W_1, \dots, W_n$ :

$$\log L_n(\theta) = n^{-1} \sum_{i=1}^n \log f(W_i, \theta).$$

The maximum likelihood (ML) estimator is defined as

$$\hat{\theta}_n^{ML} = \arg \max_{\theta \in \Theta} \log L_n(\theta).$$

Thus, for a fixed set of observations, the ML estimate is the value of  $\theta$  for which we are most likely to observe the values of  $W_1, \dots, W_n$  obtained in the sample.

In the normal regression example,

$$\log L_n(\beta, \sigma^2) = -\frac{1}{2} \log \sigma^2 - \frac{1}{2} \log(2\pi) - \frac{1}{2\sigma^2 n} \sum_{i=1}^n (Y_i - X_i'\beta)^2,$$

and

$$\begin{aligned} \hat{\beta}_{n,ML} &= \left( \sum_{i=1}^n X_i X_i' \right)^{-1} \sum_{i=1}^n X_i Y_i, \\ \hat{\sigma}_{n,ML}^2 &= n^{-1} \sum_{i=1}^n \left( Y_i - X_i' \hat{\beta}_n^{ML} \right)^2. \end{aligned}$$

In this case, the ML estimator of  $\beta$  is identical to the OLS estimator, since maximization of  $L_n$  with respect to  $\beta$  is equivalent to minimization of  $\sum_{i=1}^n (Y_i - X_i'\beta)^2$ .

### Asymptotic properties of the ML estimator

Let  $\theta_0$  be the true value of  $\theta$ .

## Consistency

By the WLLN, we should expect that for each value of  $\theta$ ,

$$\begin{aligned}\log L_n(\theta) &= n^{-1} \sum_{i=1}^n \log f(W_i, \theta) \\ &\rightarrow_p E \log f(W_i, \theta) \\ &= \int (\log f(w, \theta)) f(w, \theta_0) dw.\end{aligned}$$

Next, consider

$$\begin{aligned}& E \log f(W_i, \theta) - E \log f(W_i, \theta_0) \\ &= E \log \frac{f(W_i, \theta)}{f(W_i, \theta_0)} \\ &\leq \log E \frac{f(W_i, \theta)}{f(W_i, \theta_0)} \\ &= \log \int \frac{f(w, \theta)}{f(w, \theta_0)} f(w, \theta_0) dw \\ &= \log \int f(w, \theta) dw \\ &= \log 1 \\ &= 0.\end{aligned}$$

The inequality above follows from the fact that  $\log$  is a concave function and the Jensen's inequality: if  $f$  is concave, then  $E f(X) \leq f(EX)$ . The inequality is in fact strict provided that  $P(f(W_i, \theta_0) \neq f(W_i, \theta)) > 0$  for all  $\theta \neq \theta_0$ . As a result,  $\theta_0$  uniquely maximizes  $E \log f(W_i, \theta)$ , and, under additional technical assumptions, we have that

$$\begin{aligned}\hat{\theta}_n^{ML} &= \arg \max_{\theta \in \Theta} \log L_n(\theta) \\ &\rightarrow_p \arg \max_{\theta \in \Theta} E \log f(W_i, \theta) \\ &= \theta_0.\end{aligned}$$

## Asymptotic normality

The ML estimator solves the first order conditions

$$0 = n^{-1} \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(W_i, \hat{\theta}_n^{ML}).$$

Using the mean value theorem element-by-element, we obtain

$$0 = n^{-1} \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(W_i, \theta_0) + n^{-1} \sum_{i=1}^n \frac{\partial^2}{\partial \theta \partial \theta'} \log f(W_i, \theta_n^*) (\hat{\theta}_n^{ML} - \theta_0), \quad (1)$$

where  $\theta_n^*$  lies between  $\hat{\theta}_n^{ML}$  and  $\theta_0$ . Note that, since  $\hat{\theta}_n^{ML} \rightarrow_p \theta_0$ , we have that  $\theta_n^* \rightarrow_p \theta_0$ . Re-arranging (1) gives

$$n^{1/2} (\hat{\theta}_n^{ML} - \theta_0) = - \left( n^{-1} \sum_{i=1}^n \frac{\partial^2}{\partial \theta \partial \theta'} \log f(W_i, \theta_n^*) \right)^{-1} n^{-1/2} \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(W_i, \theta_0). \quad (2)$$

Under additional technical assumptions,

$$n^{-1} \sum_{i=1}^n \frac{\partial^2}{\partial \theta \partial \theta'} \log f(W_i, \theta_n^*) \rightarrow_p E \frac{\partial^2}{\partial \theta \partial \theta'} \log f(W_i, \theta_0). \quad (3)$$

Next, consider  $\partial \log f(W_i, \theta_0) / \partial \theta$ . Assuming that we can change the order of integration and differentiation,

$$\begin{aligned} E \frac{\partial}{\partial \theta} \log f(W_i, \theta_0) &= E \frac{\partial f(W_i, \theta_0) / \partial \theta}{f(W_i, \theta_0)} \\ &= \int \frac{\partial f(w, \theta_0) / \partial \theta}{f(w, \theta_0)} f(w, \theta_0) dw \\ &= \int \frac{\partial f(w, \theta_0)}{\partial \theta} dw \\ &= \frac{\partial}{\partial \theta} \int f(w, \theta_0) dw \\ &= \frac{\partial}{\partial \theta} 1 \\ &= 0. \end{aligned}$$

Thus, by the CLT, we should expect that

$$n^{-1/2} \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(W_i, \theta_0) \rightarrow_d N \left( 0, E \frac{\partial}{\partial \theta} \log f(W_i, \theta_0) \frac{\partial}{\partial \theta'} \log f(W_i, \theta_0) \right). \quad (4)$$

Combining (2), (3) and (4), we obtain

$$\begin{aligned} n^{1/2} (\hat{\theta}_n^{ML} - \theta_0) &\rightarrow_d - \left( E \frac{\partial^2}{\partial \theta \partial \theta'} \log f(W_i, \theta_0) \right)^{-1} N \left( 0, E \frac{\partial}{\partial \theta} \log f(W_i, \theta_0) \frac{\partial}{\partial \theta'} \log f(W_i, \theta_0) \right) \\ &= N(0, V), \end{aligned}$$

where

$$V = \left( E \frac{\partial^2}{\partial \theta \partial \theta'} \log f(W_i, \theta_0) \right)^{-1} E \frac{\partial}{\partial \theta} \log f(W_i, \theta_0) \frac{\partial}{\partial \theta'} \log f(W_i, \theta_0) \left( E \frac{\partial^2}{\partial \theta \partial \theta'} \log f(W_i, \theta_0) \right)^{-1}.$$

Next,

$$\begin{aligned} E \frac{\partial^2}{\partial \theta \partial \theta'} \log f(W_i, \theta_0) &= E \frac{\partial}{\partial \theta'} \frac{\partial f(W_i, \theta_0) / \partial \theta}{f(W_i, \theta_0)} \\ &= E \frac{\partial^2 f(W_i, \theta_0) / \partial \theta \partial \theta'}{f(W_i, \theta_0)} - E \frac{\partial f(W_i, \theta_0) / \partial \theta}{f(W_i, \theta_0)} \frac{\partial f(W_i, \theta_0) / \partial \theta'}{f(W_i, \theta_0)} \\ &= -E \frac{\partial}{\partial \theta} \log f(W_i, \theta_0) \frac{\partial}{\partial \theta'} \log f(W_i, \theta_0), \end{aligned} \quad (5)$$

since,

$$\begin{aligned} E \frac{\partial^2 f(W_i, \theta_0) / \partial \theta \partial \theta'}{f(W_i, \theta_0)} &= \int \frac{\partial^2 f(w, \theta_0) / \partial \theta \partial \theta'}{f(w, \theta_0)} f(w, \theta_0) dw \\ &= \int \frac{\partial^2 f(w, \theta_0)}{\partial \theta \partial \theta'} dw \\ &= \frac{\partial^2}{\partial \theta \partial \theta'} \int f(w, \theta_0) dw \\ &= \frac{\partial^2}{\partial \theta \partial \theta'} 1 \\ &= 0. \end{aligned}$$

The result in (5) is called the *information equality*. It implies that

$$\begin{aligned} V &= - \left( E \frac{\partial^2}{\partial \theta \partial \theta'} \log f(W_i, \theta_0) \right)^{-1} \\ &= I^{-1}(\theta_0), \end{aligned}$$

where

$$I(\theta_0) = -E \frac{\partial^2}{\partial \theta \partial \theta'} \log f(W_i, \theta_0)$$

is called the *information matrix*. Thus, we have that

$$n^{1/2} \left( \hat{\theta}_n^{ML} - \theta_0 \right) \rightarrow_d N(0, I^{-1}(\theta_0)).$$

## Remarks

- In general, ML estimation requires numerical maximization of the log-likelihood function.
- Hypothesis testing can be performed using the Wald type statistic. Suppose that  $H_0 : h(\theta_0) = 0$ , where  $h : R^k \rightarrow R^q$ . The Wald statistic is given by

$$W_n = nh \left( \hat{\theta}_n^{ML} \right)' \left( \frac{\partial h \left( \hat{\theta}_n^{ML} \right)}{\partial \theta'} I^{-1} \left( \hat{\theta}_n^{ML} \right) \frac{\partial h' \left( \hat{\theta}_n^{ML} \right)}{\partial \theta} \right)^{-1} h \left( \hat{\theta}_n^{ML} \right).$$

One should reject the null if  $W_n > \chi_{q,1-\alpha}^2$ . Alternatively, the null hypothesis  $h(\theta_0) = 0$  can be tested using the *likelihood ratio* statistic

$$LR_n = 2 \left( \log L_n \left( \hat{\theta}_n^{ML} \right) - \log L_n \left( \tilde{\theta}_n^{ML} \right) \right),$$

where  $\tilde{\theta}_n^{ML}$  is the null-restricted ML estimator:

$$\tilde{\theta}_n^{ML} = \arg \max_{\theta \in \Theta, h(\theta)=0} \log L_n(\theta).$$

Under the null,  $LR_n \rightarrow_d \chi_q^2$ , and asymptotically equivalent to the Wald statistic.

- The ML estimator is efficient. Let  $\hat{\theta}_n$  be an estimator such that  $n^{1/2} \left( \hat{\theta}_n - \theta_0 \right) \rightarrow_d N(0, \Sigma(\theta_0))$ , then  $\Sigma(\theta_0) - I^{-1}(\theta_0)$  is positive semi-definite. Thus, the asymptotic variance of the ML estimator is as small as, or smaller than, the variance of any consistent and asymptotically normal estimator.
- The ML estimation relies on very strong assumptions - the true PDF is known up to the value of parameters. If the PDF is misspecified, the estimator is called the quasi-ML estimator. In some cases, the quasi-ML estimator is still consistent. The OLS estimator is an example. It is still consistent even if the data is not normally distributed.