

Advanced Econometrics II

TA Session Problems No. 4

Agnieszka Borowska

27.01.2015

Note: this is only a draft of the problems discussed on Tuesday and might contain some typos or more or less imprecise statements. If you find some, please let me know.

1. ML Basic Concepts
2. Asymptotic Efficiency of the ML Estimator

ML Basic Concepts

Recall

$f(y, \theta) = \prod_{t=1}^n f_t(y_t, \theta),$	(joint PDF)
$\ell(y, \theta) \equiv \log f(y, \theta) = \sum_{t=1}^n \underbrace{\ell_t(y_t, \theta)}_{\text{contribution}},$	(loglikelihood)
$g(y, \theta) = (g_i(y, \theta))_{i=1, \dots, k},$	(gradient/score vector)
$g_i(y, \theta) \equiv \frac{\partial \ell(y, \theta)}{\partial \theta_i} = \sum_{t=1}^n \frac{\partial \ell_t(y_t, \theta)}{\partial \theta_i},$	(typical element of $g(y, \theta)$)
$\mathbf{H}(\theta) = \left(\frac{\partial^2 \ell(\theta)}{\partial \theta_i \partial \theta_j} \right)_{i, j=1, \dots, k},$	(Hessian matrix)
$G(y, \theta) = (G_{ti}(y^t, \theta))_{i=1, \dots, k, t=1, \dots, n},$	(matrix of contributions to the gradient)
$G_{ti}(y^t, \theta) \equiv \frac{\partial \ell(y^t, \theta)}{\partial \theta_i} = \sum_{t=1}^n \frac{\partial \ell_t(y_t, \theta)}{\partial \theta_i},$	(typical element of $G(y, \theta)$)
$g_i(y, \theta) = \sum_{t=1}^n G_{ti}(y^t, \theta),$	(10.27)
$\mathbf{I}(\theta) \equiv \sum_{t=1}^n \underbrace{\mathbf{I}_t(\theta)}_{\text{contribution}}$	(information matrix)
$\mathbf{I}_t(\theta) = \left(\mathbb{E}_\theta (G_{ti}(y^t, \theta) G_{tj}(y^t, \theta)) \right)_{i, j=1, \dots, k},$	(covariance matrix of $G_t(y^t, \theta)$)
$\mathbf{I}_t(\theta) \equiv \mathbb{E}_\theta (g(y, \theta) g^T(y, \theta))$	($G_t(y^t, \theta) - t^{th}$ row of $G(y, \theta)$)
$\mathcal{I}(\theta) \equiv \text{plim}_{n \rightarrow \infty} \frac{1}{n} \mathbf{I}(\theta),$	(covariance matrix of the score vector) (★)
$\mathcal{H}(\theta) \equiv \text{plim}_{n \rightarrow \infty} \frac{1}{n} \mathbf{H}(\theta),$	(asymptotic information matrix)
$\mathcal{I}(\theta) = -\mathcal{H}(\theta),$	(asymptotic Hessian matrix)
	(information matrix equality)

Below, we will prove (★).

Covariance matrix of the gradient vector

DM, 10.5

Prove that the definition

$$\mathbf{I}(\theta) \equiv \sum_{t=1}^n \mathbf{I}_t(\theta) = \sum_{t=1}^n \mathbb{E}_\theta(G_t^T(y, \theta)G_t(y, \theta)) \quad (10.31)$$

of the information matrix is equivalent to the definition

$$\mathbf{I}(\theta) = \mathbb{E}_\theta(g(y, \theta)g^T(y, \theta)).$$

Hint: Use the result

$$\begin{aligned} \mathbb{E}_\theta(G_{ti}^T(y^t, \theta)G_{sj}(y^s, \theta)) &= \mathbb{E}_\theta(\mathbb{E}_\theta(G_{ti}(y^t, \theta)G_{sj}(y^s, \theta)|y^t)) \\ &= \mathbb{E}_\theta(\mathbb{E}_\theta(G_{ti}(y^t, \theta)G_{sj}(y^s, \theta)|y^t)) = 0. \end{aligned} \quad (10.30)$$

The task is to show that the following relation holds

$$\mathbb{E}_\theta(g(y, \theta)g^T(y, \theta)) = \sum_{t=1}^n \mathbb{E}_\theta(G_t^T(y, \theta)G_t(y, \theta)).$$

Since by (10.27) each element of the gradient vector is the sum of the elements of one of the columns of the matrix of contributions to the gradient,

$$g_i(y, \theta) = \sum_{t=1}^n G_{ti}(y^t, \theta),$$

we can also write

$$g(y, \theta) = \sum_{t=1}^n G_t^T(y^t, \theta).$$

Then we easily obtain the required result by writing

$$\begin{aligned} \mathbb{E}_\theta(g(y, \theta)g^T(y, \theta)) &= \mathbb{E}_\theta \left(\left(\sum_{t=1}^n G_t^T(y^t, \theta) \right) \left(\sum_{s=1}^n G_s(y^s, \theta) \right) \right) \\ &= \mathbb{E}_\theta \left(\sum_{t=1}^n \sum_{s=1}^n G_t^T(y^t, \theta)G_s(y^s, \theta) \right) \\ &\stackrel{(*)}{=} \mathbb{E}_\theta \left(\sum_{t=1}^n G_t^T(y^t, \theta)G_t(y^t, \theta) \right) \\ &= \sum_{t=1}^n \mathbb{E}_\theta(G_t^T(y, \theta)G_t(y, \theta)), \end{aligned}$$

where in (*) we used that by (10.30) $\forall t \neq s$

$$\mathbb{E}_\theta \left(\sum_{t=1}^n \sum_{s=1}^n G_t^T(y^t, \theta)G_s(y^s, \theta) \right) = 0.$$

Asymptotic Efficiency of the ML Estimator

Below we will prove the ML estimator asymptotically achieves the **Cramér-Rao lower bound** – which is one of many of its attractive features. Notice, however, that it happens, in general, only **asymptotically**, as, in general, ML estimators are **not unbiased** (but are asymptotically unbiased).

To start with, consider any *other* root- n consistent and asymptotically unbiased estimator, which we will denote $\tilde{\theta}$ ($\hat{\theta}$ will stand for the MLE). It can be shown that¹

$$\text{plim}_{n \rightarrow \infty} \sqrt{n}(\tilde{\theta} - \theta_0) = \text{plim}_{n \rightarrow \infty} \sqrt{n}(\hat{\theta} - \theta_0) + v,$$

where v is a random, zero mean k -vector, uncorrelated with $\text{plim}_{n \rightarrow \infty} \sqrt{n}(\hat{\theta} - \theta_0)$. Hence, taking Var's on the both sides of the above expression gives us the following result

$$\text{Var}\left(\text{plim}_{n \rightarrow \infty} \sqrt{n}(\tilde{\theta} - \theta_0)\right) = \text{Var}\left(\text{plim}_{n \rightarrow \infty} \sqrt{n}(\hat{\theta} - \theta_0)\right) + \text{Var}(v).$$

Because $\text{Var}(v)$ is PSD, it follows that the asymptotic covariance matrix of the *other* estimator $\tilde{\theta}$ must be larger than the one of the MLE $\hat{\theta}$. This asymptotic efficiency result is an asymptotic version of Cramér-Rao lower bound (which applies to any **unbiased** estimator), stating that the covariance matrix of an unbiased estimator cannot be smaller than \mathbf{I}^{-1} . In case of the MLE the latter is asymptotically equal to its covariance matrix. So we can say that the MLE *attains* the Cramér-Rao lower bound.

DM, 10.12

Let $\hat{\theta}$ denote any unbiased estimator of the k parameters of a parametric model fully specified by the loglikelihood function $\ell(\theta)$. The unbiasedness property can be expressed as the following identity:

$$\mathbb{E}_{\theta} \tilde{\theta} = \int L(y, \theta) \tilde{\theta} dy = \theta \quad (1)$$

By using the relationship between $L(y, \theta)$ and $\ell(y, \theta)$ and differentiating this identity with respect to the components of θ , show that

$$\text{Cov}_{\theta}(g(\theta), (\tilde{\theta} - \theta)) = \mathbb{I}$$

where \mathbb{I} is a $k \times k$ identity matrix, and the notation Cov_{θ} indicates that the covariance is to be calculated under the DGP characterized by θ .

Let V denote the $2k \times 2k$ covariance matrix of the $2k$ -vector obtained by stacking the k components of $g(\theta)$ above the k components of $\tilde{\theta} - \theta$. Partition this matrix into 4 $k \times k$ blocks as follows:

$$V = \begin{bmatrix} V_1 & C \\ C^T & V_2 \end{bmatrix} \quad (2)$$

where V_1 and V_2 are, respectively, the covariance matrices of the vectors $g(\theta)$ and $\tilde{\theta} - \theta$ under the DGP characterized by θ . Then use the fact that V is positive semidefinite to show that the difference between V_2 and $\mathbb{I}^{-1}(\theta)$, where $\mathbb{I}(\theta)$ is the (finite-sample) information matrix for the model, is a positive semidefinite matrix.

First, notice that in (1) the RHS is simply θ , which differentiated wrt θ will be simply a $k \times k$ identity matrix \mathbb{I} . For the LHS, we had

$$\ell(y, \theta) = \log L(y, \theta) \quad \Leftrightarrow \quad L(y, \theta) = \exp(\ell(y, \theta)),$$

so

$$\frac{\partial L(y, \theta)}{\partial \theta} = L(y, \theta) \underbrace{\frac{\partial \ell(y, \theta)}{\partial \theta}}_{g(y, \theta)}.$$

Hence, differentiation of (1) wrt θ gives

$$\int L(y, \theta) \tilde{\theta} g(y, \theta) dy = \mathbb{I}. \quad (3)$$

¹Cf. Section 10.4 in DM.

Next, notice that $L(\cdot, \theta)$, i.e. a function of y with θ fixed, is the PDF of y , which means that the integral of two functions wrt $dL(y, \theta)$ (with θ fixed) is just the covariance of these functions wrt the distribution characterized by θ . Hence, (3) describes the covariance matrix of $\tilde{\theta}$ and $g(y, \theta) := g(\theta)$, i.e.

$$\text{Cov}_\theta(g(\theta), \tilde{\theta}) = \int g(y, \theta) \tilde{\theta} L(y, \theta) dy = \mathbb{I}.$$

However, we know that² $\mathbb{E}_\theta g(\theta) = 0$, and by the unbiasedness assumption $\mathbb{E}_\theta \tilde{\theta} = \theta$, so

$$\begin{aligned} \mathbb{I} &= \text{Cov}_\theta(g(\theta), \tilde{\theta}) \\ &\stackrel{\text{def}}{=} \mathbb{E}_\theta \left((g(\theta) - \mathbb{E}_\theta g(\theta)) (\tilde{\theta} - \mathbb{E}_\theta \tilde{\theta}) \right) \\ &= \mathbb{E}_\theta \left(g(\theta) (\tilde{\theta} - \theta) \right) \\ &= \int g(y, \theta) (\tilde{\theta} - \theta) L(y, \theta) dy \\ &= \text{Cov}_\theta(g(\theta), (\tilde{\theta} - \theta)), \end{aligned}$$

which is the first of the required results.

Second, consider V as in (2). By definition of V , its off-diagonal blocks C and C^T are given by the covariance matrix of $g(\theta)$ and $\tilde{\theta} - \theta$, which we have just shown is simply the identity matrix. As far as V_1 is concerned, it is the covariance matrix of the gradient vector $g(\theta)$ and we know from the previous exercise that it is equal to $\mathbf{I}(\theta)$, the information matrix. Finally, V as a covariance matrix needs to be PSD. Hence, we arrive at

$$V = \begin{bmatrix} \mathbf{I}(\theta) & \mathbb{I} \\ \mathbb{I} & V_2 \end{bmatrix} \geq 0,$$

which implies that also

$$V^{-1} \geq 0$$

as well as each of the diagonal blocks of V^{-1} must also be PSD. To complete the exercise we need to show that

$$V_2 - \mathbf{I}^{-1}(\theta) \geq 0.$$

Recall the discussion about the asymptotic distribution of the Wald statistics, where we needed the inverse of the covariance matrix of the IV estimator corresponding to part of this vector³. We used there the following fact

$$A := \begin{bmatrix} A_1 & A_{12} \\ A_{21} & A_2 \end{bmatrix} \Rightarrow A^{-1} = \begin{bmatrix} \cdot & \cdot \\ \cdot & (A_{22} - A_{21} A_{11}^{-1} A_{12})^{-1} \end{bmatrix},$$

which applied to the current case yields

$$V^{-1} = \begin{bmatrix} \cdot & \cdot \\ \cdot & (V_2 - \mathbb{I} \mathbf{I}^{-1}(\theta) \mathbb{I})^{-1} \end{bmatrix}.$$

So the lower diagonal block of the inverse of V , which, as we stated, is PSD, has the form

$$(V_2 - \mathbb{I} \mathbf{I}^{-1}(\theta) \mathbb{I})^{-1} = (V_2 - \mathbf{I}^{-1}(\theta))^{-1}.$$

This means that its inverse,

$$V_2 - \mathbf{I}^{-1}(\theta) \geq 0,$$

i.e. is PSD – which is the second of the required results.

²Recall that a crucial property of the matrix $G(y, \theta)$ is that **if y is generated by the GDP characterized by θ** , then the expectations of all the elements of the matrix, **evaluated at θ** , are zero – which is a consequence of the fact that all densities integrate to 1. Then, summing the expectations of the elements in each column of $G(y, \theta)$ yields that $\mathbb{E}_\theta g(y, \theta) = 0$.

³Week 1, slide 48.