

1. Replicate Figures 1,5,6 and 7 in the Matlab exercise. In doing so, make the following modifications to the figures.

- (1) In figure 1, use “plot” instead of “stairs”, and use a dotted line instead of a solid line.
- (2) In figure 5, plot the data from Jan 2004 only and remove all grid lines.
- (3) Put figures 6 and 7 on a common figure frame in  $2 \times 2$  format.

**Solution:**

2. Suppose that the data are generated by the following model

$$y = X\beta + u \quad (1)$$

$$= X_1\beta_1 + X_2\beta_2 + u \quad (2)$$

where  $X = [X_1, X_2]$ . Assume that  $X$  is of full column rank,  $T^{-1}X'X \xrightarrow{p} Q$ , and  $\beta_2 \neq 0$ . Denote the OLS estimate for (1) by  $\hat{\beta}_1$  and  $\hat{\beta}_2$ . Suppose you estimate the regression

$$y = X_1\beta_1 + e \quad (3)$$

by OLS and denote the resulting estimate by  $\hat{b}_1$ .

- a) Show that  $\hat{b}_1$  is inconsistent for  $\beta_1$ , with assuming  $E(X_1'X_2) \neq 0$ .
- b) The inconsistency of  $\hat{b}_1$  is an example of the omitted variable bias. A natural estimate would then be based on an instrumental variable procedure. Show that the OLS estimate  $\hat{\beta}_1$  can indeed be given an IV interpretation.

**Solution:**

- a) The omitted variable bias lingers even when the sample size grows large causing the estimators to be inconsistent. We can show this by simply plugging in the true model of  $y$  to the estimator  $\hat{b}_1$ .

$$\hat{b}_1 = (X_1'X_1)^{-1} X_1'y \quad (4)$$

$$= (X_1'X_1)^{-1} X_1'(X_1\beta_1 + X_2\beta_2 + u) \quad (5)$$

$$= \beta_1 + (X_1'X_1)^{-1} X_1'X_2\beta_2 + (X_1'X_1)^{-1} X_1'u \quad (6)$$

The problem states that  $E(X_1'X_2) \neq 0$ ,  $\beta_2 \neq 0$ , preserving the second term. Note that the only *random variable* in the above equation is  $u$ , converging in probability to zero as  $T \rightarrow \infty$  where  $T$  is the dimension of  $y$  ( $T \times 1$ ). Thus, there is no reason for the second term  $(X_1'X_1)^{-1} X_1'X_2\beta_2$  to disappear.

$$\hat{b}_1 \xrightarrow{p} \beta_1 + \beta_2 (X_1'X_1)^{-1} X_1'X_2 \quad (7)$$

- b)  $X_2$  is by assumption independent of the error term  $u$ . Therefore, using the method of moment, if we have  $E(x_{i2}u_i) = 0$ ,

$$\frac{1}{T} \sum_{i=1}^T x_{i2}u_i = \frac{1}{T} X_2'u = \frac{1}{T} X_2'(y - X_1\beta_1). \quad (8)$$

Solving with respect to  $\beta_1$  yields

$$\hat{\beta}_1 = (X_2'X_1)^{-1} X_2'y \quad (9)$$

which is the IV estimate.

### 3. Consider the linear regression model:

$$y = X\beta + u \quad (10)$$

$$y, u : T \times 1 \quad (11)$$

$$X : T \times k \quad (12)$$

$$\beta : k \times 1 \quad (13)$$

with the  $q$  moment conditions  $E(z_t u_t) = 0$ . Let

$$J_T(\beta, W_T) = g_T(\beta)' W_T g_T(\beta) \quad (14)$$

where  $g_T(\beta) = T^{-1} \sum z_t (y_t - x_t'\beta)$  and  $W_T$  is some weighting matrix. The GMM estimator for  $\beta$  is obtained as the minimizer of  $J_T(\beta, W_T)$ .

- How does your choice of  $W_T$  affect the GMM estimator? Discuss the implications on the consistency, the asymptotic normality and efficiency.
- If the model is exactly identified ( $k = q$ ), explain why the choice of  $W_T$  becomes irrelevant.

**Solution:**

a)

4. Let  $y_t \stackrel{\text{iid}}{\sim} \text{Exp}(\theta)$  for  $t = 1, \dots, T$ .

- a) Derive the score function, Hessian function and information matrix, using the exponential density.
- b) Derive the MLE for  $\theta$ . Sketch the proof that the MLE is asymptotically normal. Be specific with the asymptotic variance.

**Solution:**

a) The likelihood function of  $y_1, \dots, y_T$ ,

$$L(\theta | y_1, \dots, y_T) = \theta^T \exp\left(-\theta \sum_{t=1}^T y_t\right). \quad (15)$$

Taking the logarithm yields

$$\ell(\theta | y_1, \dots, y_T) = T \log \theta - \theta \sum_{t=1}^T y_t. \quad (16)$$

By definition of the score function is the first derivative of the log-likelihood.

$$\frac{d\ell}{d\theta} = \frac{T}{\theta} - \sum_{t=1}^T y_t. \quad (17)$$

The Hessian gets reduced to the second derivative for a univariate function.

$$\frac{d^2\ell}{d\theta^2} = -\frac{T}{\theta^2} - \sum_{t=1}^T y_t. \quad (18)$$

To get the Fisher information,

$$\mathcal{I}_T(\theta) = \frac{1}{T} \text{E} \left( \frac{T}{\theta^2} + \sum_{t=1}^T y_t \right) \quad (19)$$

$$= \frac{1}{\theta^2} + \frac{1}{T} \sum_{t=1}^T \text{E}(y_t) \quad (20)$$

$$= \frac{1}{\theta^2} + \frac{1}{\theta} \quad (21)$$

b) We have obtained the first and second derivatives of the log-likelihood already. Recall:

$$\frac{d\ell}{d\theta} = \frac{T}{\theta} - \sum_{t=1}^T y_t \quad (22)$$

$$\frac{d^2\ell}{d\theta^2} = -\frac{T}{\theta^2} - \sum_{t=1}^T y_t \quad (23)$$

Using the first-order condition, the MLE comes with a closed-form expression.

$$\hat{\theta}^{\text{MLE}} = T / \left( \sum_{t=1}^T y_t \right) \quad (24)$$

The second-order condition validates that the estimator is actually a maximum since

$$\frac{d^2\ell}{d\theta^2} < 0. \quad (25)$$

Now, even as rough a proof as what we will shortly give here takes at least 2 steps: the consistency of MLE and the asymptotic normality of MLE.

- (*Consistency*) Recall that we take the product of every single PDF of  $y_t$  through  $t = 1, \dots, T$  to compute the likelihood function, which also means taking the logarithm will convert the product into summation.

$$\ell(\theta | y_1, \dots, y_T) = \sum_{t=1}^T \ell(\theta | y_t) \quad (26)$$

By the strong law of large numbers, we get the following relation:

$$\frac{1}{T} \sum_{t=1}^T \ell(\theta | y_t) \xrightarrow{a.s.} E_{\theta_0} \ell(\theta | y_1) \quad (27)$$

for some unknown true parameter value  $\theta_0$ . We can then show that the expected log-likelihood function w.r.t. the true parameter is always greater than that of an arbitrary parameter  $\theta$  by *Kullback-Leibler divergence*. The KL divergence is defined as follows.

$$\text{KL}(f(y_1 | \theta_0) \| f(y_1 | \theta)) = E_{\theta_0} \left[ \log \frac{f(y_1 | \theta_0)}{f(y_1 | \theta)} \right] \quad (28)$$

$$= - \int \log \frac{f(y_1 | \theta)}{f(y_1 | \theta_0)} f(y_1 | \theta_0) dy_1 \quad (29)$$

By Jensen's inequality,

$$\underbrace{-\log \int \frac{f(y_1 | \theta)}{f(y_1 | \theta_o)} f(y_1 | \theta_o) dy_1}_{=0} \leq \underbrace{-\int \log \frac{f(y_1 | \theta)}{f(y_1 | \theta_o)} f(y_1 | \theta_o) dy_1}_{=KL(f(y_1 | \theta_o) \| f(y_1 | \theta))}. \quad (30)$$

Therefore, it always follows that the KL divergence is nonnegative. In fact, it is strictly positive if  $f(y_1 | \theta) \neq f(y_1 | \theta_o)$ . This indicates that

$$\theta_o = \sup_{\theta \in \Omega} E_{\theta_o} \ell(\theta | y_1). \quad (31)$$

Recall the following:

$$\hat{\theta}^{MLE} = \sup_{\theta \in \Omega} \frac{1}{T} \sum_{t=1}^T \ell(\theta | y_t). \quad (32)$$

Therefore by 27,  $\hat{\theta}^{MLE} \xrightarrow{P} \theta_o$  for a finite parameter space  $\Omega$ . We can also prove this for a compact parameter space by starting from the lemma that

$$\frac{1}{T} \sum_{t=1}^T \ell(\theta | y_t) \xrightarrow{\text{uniformly convergent}} \int \ell(\theta | y_1) f(y_1 | \theta_o) dy_1 \quad (= E_{\theta_o} \ell(\theta | y_1)) \quad (33)$$

which is equivalent to

$$\Pr \left( \sup_{\theta \in \Omega} \left| \frac{1}{T} \sum_{t=1}^T \ell(\theta | y_t) - E_{\theta_o} \ell(\theta | y_1) \right| > \epsilon \right) \xrightarrow{a.s.} 0, \quad \forall \epsilon > 0. \quad (34)$$

Pointwise convergence is not enough with infinite parameter spaces because the convergence at one parameter of the log-likelihood function as a function of  $y_{1:T}$  does not guarantee that the log-likelihood function with another set of  $y_{1:T}$  generated with a different parameter value will not be close to the true expected log-likelihood. Thus, the convergence at one parameter value does not generalize. Anyway, the proof for the consistency of MLE ends here.

- (Asymptotic Normality) We approximate the score function with its first-order Taylor expansion around the true parameter  $\theta_o$  and apply the mean value theorem.

$$\frac{d\ell}{d\theta} \approx \frac{d\ell}{d\theta} \Big|_{\theta=\theta_o} + \frac{d^2\ell}{d\theta^2} \Big|_{\theta=\bar{\theta}} (\theta - \theta_o) \quad (35)$$

where  $\bar{\theta}$  lies somewhere between  $\theta$  and  $\theta_o$ . Since MLE is the value which sets the first

derivative to zero, we can think of the following identity.

$$\left. \frac{d\ell}{d\theta} \right|_{\theta=\theta_o} + \left. \frac{d^2\ell}{d\theta^2} \right|_{\theta=\bar{\theta}} (\hat{\theta}^{\text{MLE}} - \theta_o) = 0 \quad (36)$$

This, in turn, translates to the following relationship.

$$\hat{\theta}^{\text{MLE}} - \theta_o = - \left( \left. \frac{d\ell}{d\theta} \right|_{\theta=\theta_o} \right) / \left( \left. \frac{d^2\ell}{d\theta^2} \right|_{\theta=\bar{\theta}} \right) \quad (37)$$

with  $\bar{\theta} = s\hat{\theta}^{\text{MLE}} + (1-s)\theta_o$ ,  $s \in [0, 1]$ . Let's slowly examine the RHS of 37. First, the numerator can be expressed as a summation of the log-likelihood of a single observation.

$$\left. \frac{d\ell(\theta | y_{1:T})}{d\theta} \right|_{\theta=\theta_o} = \sum_{t=1}^T \left. \frac{d\ell(\theta | y_t)}{d\theta} \right|_{\theta=\theta_o} \quad (38)$$

By the *Central Limit Theorem*,

$$E \left( \left. \frac{d\ell(\theta | y_1)}{d\theta} \right|_{\theta=\theta_o} \right) = 0 \quad (39)$$

$$\text{Var} \left( \left. \frac{d\ell(\theta | y_1)}{d\theta} \right|_{\theta=\theta_o} \right) = T\mathcal{I}_1(\theta) \quad (40)$$

$$\left. \frac{d\ell(\theta | y_{1:T})}{d\theta} \right|_{\theta=\theta_o} \xrightarrow{d} \mathcal{N}(0, T\mathcal{I}_1(\theta_o)) \quad (41)$$

where  $\mathcal{I}_1$  is the Fisher information for a single observation  $y_1$ . The calculations for the expectation and the variance are given in the end. Now the denominator behaves like the following which makes use of the weak law of large numbers.

$$\left. \frac{d^2\ell(\theta | y_{1:T})}{d\theta^2} \right|_{\theta=\bar{\theta}} = \sum_{t=1}^T \left. \frac{d^2\ell(\theta | y_t)}{d\theta^2} \right|_{\theta=\bar{\theta}} \xrightarrow{p} TE \left( \left. \frac{d^2\ell(\theta | y_1)}{d\theta^2} \right|_{\theta=\bar{\theta}} \right) = -T\mathcal{I}_1(\bar{\theta}) \quad (42)$$

With 39, 42, and the *Slutsky's theorem*, we can conclude that 37 converges in distribution to a normal distribution.

$$- \left( \left. \frac{d\ell}{d\theta} \right|_{\theta=\theta_o} \right) / \left( \left. \frac{d^2\ell}{d\theta^2} \right|_{\theta=\bar{\theta}} \right) \xrightarrow{d} \mathcal{N} \left( 0, \frac{T\mathcal{I}_1(\theta_o)}{(T\mathcal{I}_1(\bar{\theta}))^2} = \frac{\mathcal{I}_T(\theta_o)}{(\mathcal{I}_T(\bar{\theta}))^2} \right) \quad (43)$$

Finally, we know that  $\bar{\theta} \in [\hat{\theta}^{\text{MLE}}, \theta_o]$  if  $\hat{\theta}^{\text{MLE}} < \theta_o$  or  $\bar{\theta} \in [\theta_o, \hat{\theta}^{\text{MLE}}]$  if  $\hat{\theta}^{\text{MLE}} \geq \theta_o$ . As

$T \rightarrow \infty$ ,  $\hat{\theta}^{\text{MLE}} \xrightarrow{p} \theta_0$  which also means  $\bar{\theta} \xrightarrow{p} \theta_0$ . Thus,

$$\hat{\theta}^{\text{MLE}} \xrightarrow{d} \mathcal{N}(\theta_0, (\mathcal{I}_T(\theta_0))^{-1}) \quad (44)$$

There are 3 different ways to compute the Fisher information and all three are equivalent under regularity conditions. The three are

$$\mathcal{I}(\theta) = \text{E} \left( \left( \frac{d\ell}{d\theta} \right)^2 \right) \quad (45)$$

$$= -\text{E} \left( \frac{d^2\ell}{d\theta^2} \right) \quad (46)$$

$$= \text{Var} \left( \frac{d\ell}{d\theta} \right) \quad (47)$$