

# Econometrics

## Homework 5

**Problem 1.** Consider the following simple regression without an intercept:

$$Y_i = \beta X_i + U_i.$$

Assume the observations  $(Y_i, X_i)$ ,  $i = 1, 2, \dots, n$  are iid. Assume  $E(X_i U_i) = 0$ ,  $E(X_i^2 U_i^2) < \infty$  and  $0 < E(X_i^2) < \infty$ .

(a) Provide the expression of the OLS estimator  $\hat{\beta}_n$  and show it is a consistent estimator of  $\beta$ .

(b) Show that

$$\sqrt{n}(\hat{\beta}_n - \beta) \rightarrow_d N(0, V), \text{ where } V = \frac{E(X_i^2 U_i^2)}{E(X_i^2)^2}.$$

(c) How to construct a consistent estimator of  $V$ ? You do not need to show your estimator is consistent.

### Solution.

(a) The OLS minimization problem:

$$\min_b \sum_{i=1}^n (Y_i - bX_i)^2.$$

The first-order condition:

$$-2 \sum_{i=1}^n X_i (Y_i - \hat{\beta}_n X_i) = 0 \implies \hat{\beta}_n = \frac{\sum_{i=1}^n X_i Y_i}{\sum_{i=1}^n X_i^2}.$$

Decompose:

$$\begin{aligned} \hat{\beta}_n &= \frac{\sum_{i=1}^n X_i Y_i}{\sum_{i=1}^n X_i^2} \\ &= \frac{\sum_{i=1}^n X_i (\beta X_i + U_i)}{\sum_{i=1}^n X_i^2} \\ &= \beta + \frac{\frac{1}{n} \sum_{i=1}^n X_i U_i}{\frac{1}{n} \sum_{i=1}^n X_i^2}. \end{aligned}$$

By LLN,

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n X_i U_i &\rightarrow_p E(X_i U_i) = 0 \\ \frac{1}{n} \sum_{i=1}^n X_i^2 &\rightarrow_p E(X_i^2) > 0. \end{aligned}$$

By Slutsky,

$$\frac{\frac{1}{n} \sum_{i=1}^n X_i U_i}{\frac{1}{n} \sum_{i=1}^n X_i^2} \xrightarrow{p} \frac{E(X_i U_i)}{E(X_i^2)} = 0.$$

(b)

$$\widehat{\beta}_n - \beta = \frac{\frac{1}{n} \sum_{i=1}^n X_i U_i}{\frac{1}{n} \sum_{i=1}^n X_i^2} \implies \sqrt{n} (\widehat{\beta}_n - \beta) = \frac{\frac{1}{\sqrt{n}} \sum_{i=1}^n X_i U_i}{\frac{1}{n} \sum_{i=1}^n X_i^2}.$$

By CLT,

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n X_i U_i \rightarrow_d N(0, E(X_i^2 U_i^2)),$$

therefore,

$$\sqrt{n} (\widehat{\beta}_n - \beta) \rightarrow_d \frac{1}{E(X_i^2)} N(0, E(X_i^2 U_i^2)) \sim N\left(0, \frac{E(X_i^2 U_i^2)}{E(X_i^2)^2}\right)$$

(c) An estimator:

$$\frac{\frac{1}{n} \sum_{i=1}^n X_i^2 \widehat{U}_i^2}{\left(\frac{1}{n} \sum_{i=1}^n X_i^2\right)^2},$$

where  $\widehat{U}_i = Y_i - \widehat{\beta}_n X_i$ .

**Problem 2.** (a) Consider the following simple regression model:

$$Y_i = \alpha + \beta X_i + U_i.$$

Suppose the observations  $(Y_i, X_i)$ ,  $i = 1, 2, \dots, n$  are iid. Assume  $E|U_i| < \infty$ ,  $E|X_i| < \infty$  and  $E(U_i) = 0$ . Let  $\widetilde{\beta}_n$  be any consistent estimator of  $\beta$  (not necessarily the OLS estimator). Define the following estimator for  $\alpha$ :

$$\widetilde{\alpha}_n = \overline{Y}_n - \widetilde{\beta}_n \overline{X}_n,$$

where  $\overline{Y}_n = n^{-1} \sum_{i=1}^n Y_i$  and  $\overline{X}_n = n^{-1} \sum_{i=1}^n X_i$ . Prove that  $\widetilde{\alpha}_n$  is a consistent estimator of  $\alpha$ . Hint: Show  $\overline{Y}_n = \alpha + \beta \overline{X}_n + \overline{U}_n$ , where  $\overline{U}_n = n^{-1} \sum_{i=1}^n U_i$ .

(b) Consider the following regression model without a regressor:

$$Y_i = \alpha + U_i.$$

Suppose the observations  $Y_i$ ,  $i = 1, 2, \dots, n$  are iid and  $E(Y_i^2) < \infty$ . Assume  $E(U_i) = 0$ . What is the expression of the OLS estimator  $\widehat{\alpha}_n$ ? Show that  $\sqrt{n}(\widehat{\alpha}_n - \alpha) \rightarrow_d N(0, V)$  and find  $V$ .

**Solution.**

(a)

$$Y_i = \alpha + \beta X_i + U_i \implies \overline{Y}_n = \alpha + \beta \overline{X}_n + \overline{U}_n,$$

therefore

$$\begin{aligned} \widetilde{\alpha}_n &= \overline{Y}_n - \widetilde{\beta}_n \overline{X}_n \\ &= \alpha + \beta \overline{X}_n + \overline{U}_n - \widetilde{\beta}_n \overline{X}_n \\ &= \alpha + \overline{X}_n (\beta - \widetilde{\beta}_n) + \overline{U}_n. \end{aligned}$$

By WLLN,

$$\begin{aligned}\bar{U}_n &\rightarrow_p E(U_i) = 0 \\ \bar{X}_n &\rightarrow_p E(X_i).\end{aligned}$$

By assumption,

$$\beta - \tilde{\beta}_n \rightarrow_p 0.$$

Therefore, by Slutsky Lemma,

$$\begin{aligned}\tilde{\alpha}_n - \alpha &= \bar{X}_n (\beta - \tilde{\beta}_n) + \bar{U}_n \\ &\rightarrow_p E(X_i) \times 0 + 0.\end{aligned}$$

(b) The OLS estimator  $\hat{\alpha}_n$  is the solution of

$$\min_a \sum_{i=1}^n (Y_i - a)^2.$$

First-order condition:

$$\begin{aligned}-2 \sum_{i=1}^n (Y_i - \hat{\alpha}_n) &= 0 \implies \hat{\alpha}_n = \frac{1}{n} \sum_{i=1}^n Y_i. \\ E(Y_i) &= \alpha + E(U_i) = \alpha \implies \alpha = E(Y_i).\end{aligned}$$

By CLT,

$$\sqrt{n}(\hat{\alpha}_n - \alpha) \rightarrow_d N(0, \text{Var}(Y_i)).$$

$$\text{Var}(Y_i) = \text{Var}(\alpha + U_i) = \text{Var}(U_i).$$

**Problem 3.** Let  $Y$  be the face number showing when a die is rolled. Define  $X$  as

$$X = \begin{cases} Y & \text{if } Y \text{ is even,} \\ 0 & \text{if } Y \text{ is odd.} \end{cases}$$

Let  $R(Y|X)$  denote the best linear approximation to the conditional expectation  $E(Y|X)$ .  $R(Y|X) = \beta_0 + \beta_1 X$ , where

$$(\beta_0, \beta_1) = \underset{b_0, b_1}{\operatorname{argmin}} E[(E(Y|X) - b_0 - b_1 X)^2].$$

Calculate  $E[(Y - R(Y|X))^2]$  and  $E[(Y - E(Y|X))^2]$ .

**Solution.** First we derive the joint distribution of  $(X, Y)$ :

$$\begin{aligned}\frac{1}{6} &= \Pr((X, Y) = (0, 1)) \\ &= \Pr((X, Y) = (0, 3)) \\ &= \Pr((X, Y) = (0, 5)) \\ &= \Pr((X, Y) = (2, 2)) \\ &= \Pr((X, Y) = (4, 4)) \\ &= \Pr((X, Y) = (6, 6)).\end{aligned}$$

It is easy to derive that the best linear predictor is:

$$R(Y|X) = E(Y) + \frac{\text{Cov}(X, Y)}{\text{Var}(X)} (X - E(X)).$$

We have

$$\begin{aligned} E(Y) &= \sum_{i=1}^6 \frac{i}{6} = \frac{7}{2} \\ E(X) &= \frac{1}{2} \times 0 + \frac{1}{6} (2 + 4 + 6) = 2 \\ \text{Var}(X) &= \frac{1}{2} \times (0 - 2)^2 + \frac{1}{6} [(2 - 2)^2 + (4 - 2)^2 + (6 - 2)^2] = \frac{16}{3} \\ E(XY) &= \frac{1}{6} (2^2 + 4^2 + 6^2) = \frac{28}{3} \\ \text{Cov}(X, Y) &= \frac{28}{3} - 2 \times \frac{28}{3} = \frac{7}{3}. \end{aligned}$$

So

$$R(Y|X) = \frac{21}{8} + \frac{7}{16}X.$$

The conditional mean is:

$$E(Y|X) = \begin{cases} X & \text{if } X = 2, 4, 6 \\ \frac{1}{3}(1 + 3 + 5) & \text{if } X = 0. \end{cases}$$

Calculate:

$$\begin{aligned} E((Y - R(Y|X))^2) &= \frac{1}{6} \left[ \left(1 - \frac{21}{8}\right)^2 + \left(3 - \frac{21}{8}\right)^2 + \left(5 - \frac{21}{8}\right)^2 + \left(2 - \left(\frac{21}{8} + \frac{7}{16} \times 2\right)\right)^2 \right. \\ &\quad \left. + \left(4 - \left(\frac{21}{8} + \frac{7}{16} \times 4\right)\right)^2 + \left(6 - \left(\frac{21}{8} + \frac{7}{16} \times 6\right)\right)^2 \right] \\ &= 1.896 \end{aligned}$$

and

$$E[(Y - E(Y|X))^2] = \frac{1}{6} [(1 - 3)^2 + (3 - 3)^2 + (5 - 3)^2] = 1.333.$$

**Problem 4.** Consider a simple model to estimate the effect of personal computer (PC) ownership on college grade point average for graduating seniors at a large public university:

$$GPA = \beta_0 + \beta_1 PC + u,$$

where  $PC$  is a binary variable indicating PC ownership.

- (i) Why might PC ownership be correlated with  $u$ ?
- (ii) Explain why  $PC$  is likely to be related to parents' annual income. Does this mean parental income is a good IV for  $PC$ ? Why or why not?

- (iii) Suppose that, four years ago, the university gave grants to buy computers to roughly one-half of the incoming students, and the students who received grants were randomly chosen. Carefully explain how you would use this information to construct an instrumental variable for  $PC$ .

**Solution.**

- (i) Personal income/wealth is an omitted variable in this regression, as wealth can affect GPA (the student can hire more tutors), however, wealth is correlated with  $PC$ .
- (ii) While parents' income is correlated with  $PC$  ownership, it does not satisfy the exogeneity assumption as it is correlated with personal income of the student.
- (iii) A dummy variable indicating whether the student received the grant can be used as an IV: it is correlated with  $PC$  and independent of students' income since it was given randomly. However, the econometrician will have to restrict his sample only to those who began their studies four years ago when the university gave the grant.

**Problem 5.** Suppose we observe the i.i.d. random sample  $\{(Y_i, X_i)\}_{i=1}^n$ . Denote  $\bar{X}_n = n^{-1} \sum_{i=1}^n X_i$ ,  $\bar{Y}_n = n^{-1} \sum_{i=1}^n Y_i$ ,  $\mu_X = E[X_i]$  ( $\mu_X \neq 0$ ) and  $\mu_Y = E[Y_i]$ . We are interested in  $\mu_Y/\mu_X$ . Derive the asymptotic distribution of  $\sqrt{n}(\bar{Y}_n/\bar{X}_n - \mu_Y/\mu_X)$ . Hint: Write

$$\begin{aligned} \frac{\bar{Y}_n}{\bar{X}_n} &= \frac{\bar{Y}_n}{\mu_X} \cdot \left( \frac{\mu_X}{\bar{X}_n} - 1 \right) + \frac{\bar{Y}_n}{\mu_X} \\ &= -\frac{\bar{Y}_n}{\mu_X \bar{X}_n} \cdot (\bar{X}_n - \mu_X) + \frac{\bar{Y}_n}{\mu_X} \\ &= -\left( \frac{\bar{Y}_n}{\mu_X \bar{X}_n} - \frac{\mu_Y}{\mu_X^2} + \frac{\mu_Y}{\mu_X^2} \right) \cdot (\bar{X}_n - \mu_X) + \frac{\bar{Y}_n}{\mu_X} \end{aligned}$$

You may use the following result:  $W_n \rightarrow_d N(0, \sigma^2)$  and  $\theta_n \rightarrow_p 0$ , then  $\theta_n W_n \rightarrow_p 0$ .

**Solution.**

$$\begin{aligned} \frac{\bar{Y}_n}{\bar{X}_n} - \frac{\mu_Y}{\mu_X} &= \left( \frac{\bar{Y}_n}{\mu_X} - \frac{\mu_Y}{\mu_X} \right) - \left( \frac{\bar{Y}_n}{\mu_X \bar{X}_n} - \frac{\mu_Y}{\mu_X^2} + \frac{\mu_Y}{\mu_X^2} \right) \cdot (\bar{X}_n - \mu_X) \\ &= \frac{1}{n} \sum_{i=1}^n \left\{ \frac{1}{\mu_X} (Y_i - \mu_Y) - \frac{\mu_Y}{\mu_X^2} (X_i - \mu_X) \right\} - \left( \frac{\bar{Y}_n}{\mu_X \bar{X}_n} - \frac{\mu_Y}{\mu_X^2} \right) (\bar{X}_n - \mu_X). \end{aligned}$$

By Slutsky lemma,  $\frac{\bar{Y}_n}{\mu_X \bar{X}_n} - \frac{\mu_Y}{\mu_X^2} \rightarrow_p 0$ . By CLT,  $\sqrt{n}(\bar{X}_n - \mu_X) \rightarrow_d N(0, \sigma_X^2)$ . Therefore,

$$-\left( \frac{\bar{Y}_n}{\mu_X \bar{X}_n} - \frac{\mu_Y}{\mu_X^2} \right) \sqrt{n}(\bar{X}_n - \mu_X) \rightarrow_p 0.$$

By CLT,

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \left\{ \frac{1}{\mu_X} (Y_i - \mu_Y) - \frac{\mu_Y}{\mu_X^2} (X_i - \mu_X) \right\} \rightarrow_d N \left( 0, E \left[ \left( \frac{1}{\mu_X} (Y_i - \mu_Y) - \frac{\mu_Y}{\mu_X^2} (X_i - \mu_X) \right)^2 \right] \right).$$

We have

$$\sqrt{n} \left( \frac{\bar{Y}_n}{\bar{X}_n} - \frac{\mu_Y}{\mu_X} \right) \rightarrow_d N \left( 0, E \left[ \left( \frac{1}{\mu_X} (Y_i - \mu_Y) - \frac{\mu_Y}{\mu_X^2} (X_i - \mu_X) \right)^2 \right] \right).$$

**Problem 6.** Suppose that you wish to estimate the effect of class attendance on student performance.  $stndfnl$  is the standardized outcome on a final exam. A basic model is

$$stndfnl = \beta_0 + \beta_1 atndrte + \beta_2 priGPA + \beta_3 ACT + u,$$

where  $atndrte$  is percentage of classes attended,  $priGPA$  is prior college grade point average, and  $ACT$  is the achievement test score.

- (i) Let  $dist$  be the distance from the students' living quarters to the lecture hall. Do you think  $dist$  is uncorrelated with  $u$ ?
- (ii) Assuming that  $dist$  and  $u$  are uncorrelated, what other assumption must  $dist$  satisfy to be a valid IV for  $atndrte$ ?
- (iii) Suppose we add the interaction term  $priGPA \times atndrte$ :

$$stndfnl = \beta_0 + \beta_1 atndrte + \beta_2 priGPA + \beta_3 ACT + \beta_4 priGPA \times atndrte + u.$$

If  $atndrte$  is correlated with  $u$ , then, in general, so is  $priGPA \times atndrte$ . What might be a good IV for  $priGPA \times atndrte$ ? Hint: If  $E(u|priGPA, ACT, dist) = 0$ , as happens when  $priGPA$ ,  $ACT$ , and  $dist$  are all exogenous, then any function of  $priGPA$  and  $dist$  is uncorrelated with  $u$ .

**Solution.**

- (i) It seems reasonable to assume that  $dist$  and  $u$  are uncorrelated because classrooms are not usually assigned with convenience for particular students in mind.
- (ii) The variable  $dist$  must be partially correlated with  $atndrte$ . More precisely, in the reduced form

$$atndrte = \pi_0 + \pi_1 priGPA + \pi_2 ACT + \pi_3 dist + v,$$

we must have  $\pi_3 \neq 0$ . Given a sample of data we can test  $H_0 : \pi_3 = 0$  against  $H_1 : \pi_3 \neq 0$  using a  $t$  test.

- (iii) We now need instrumental variables for  $atndrte$  and the interaction term,  $priGPA \cdot atndrte$ . (Even though  $priGPA$  is exogenous,  $atndrte$  is not, and so  $priGPA \cdot atndrte$  is generally correlated with  $u$ .) Under the exogeneity assumption that  $E(u|priGPA, ACT, dist) = 0$ , any function of  $priGPA$ ,  $ACT$ , and  $dist$  is uncorrelated with  $u$ . In particular, the interaction  $priGPA \cdot dist$  is uncorrelated with  $u$ . If  $dist$  is partially correlated with  $atndrte$  then  $priGPA \cdot dist$  is partially correlated with  $priGPA \cdot atndrte$ . So, we can estimate the equation

$$stndfnl = \beta_0 + \beta_1 atndrte + \beta_2 priGPA + \beta_3 ACT + \beta_4 priGPA \times atndrte + u$$

by 2SLS using IVs  $dist$ ,  $priGPA$ ,  $ACT$ , and  $priGPA \cdot dist$ .