

Econometrics

Homework 6

Problem 1. Consider a simple regression model (with an intercept):

$$Y_i = \beta_0 + \beta_1 X_i + U_i,$$

and the IV estimator of β_1 :

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (Z_i - \bar{Z}) Y_i}{\sum_{i=1}^n (Z_i - \bar{Z}) X_i},$$

where

$$\bar{Z} = n^{-1} \sum_{i=1}^n Z_i.$$

Suppose that Z_i is a dummy variable. Show that $\hat{\beta}_1$ can be written as

$$\hat{\beta}_1 = \frac{\bar{Y}_1 - \bar{Y}_0}{\bar{X}_1 - \bar{X}_0},$$

where \bar{X}_0 and \bar{Y}_0 are the sample averages of X_i and Y_i over the part of the sample with $Z_i = 0$, and \bar{X}_1 and \bar{Y}_1 are the sample averages of X_i and Y_i over the part of the sample with $Z_i = 1$, by following the steps below. Let n_1 be the number of observations in the part of the sample with $Z_i = 1$. Let n_0 be the number of observations in the part of the sample with $Z_i = 0$. Hint: $n_1 = \sum_{i=1}^n Z_i$, $n_0 = n - \sum_{i=1}^n Z_i = \sum_{i=1}^n (1 - Z_i)$. $\bar{Y}_1 = \sum_{i=1}^n Z_i Y_i / n_1$ and $\bar{Y}_0 = \sum_{i=1}^n (1 - Z_i) Y_i / n_0$

(i) Show that $\sum_{i=1}^n (Z_i - \bar{Z}) Y_i = \sum_{i=1}^n Z_i (Y_i - \bar{Y})$.

(ii) Show that $\sum_{i=1}^n Z_i (Y_i - \bar{Y}) = n_1 (\bar{Y}_1 - \bar{Y})$.

(iii) Show that $\bar{Y} = (n_1 \bar{Y}_1 + n_0 \bar{Y}_0) / n$.

(iv) Show that $n_1 (\bar{Y}_1 - \bar{Y}) = (n_1 n_0 / n) (\bar{Y}_1 - \bar{Y}_0)$.

(v) Show how (i)-(iv) imply that $\hat{\beta}_1 = (\bar{Y}_1 - \bar{Y}_0) / (\bar{X}_1 - \bar{X}_0)$.

Solution. 1.

$$\begin{aligned} \sum_{i=1}^n (Z_i - \bar{Z}) Y_i &= \sum_{i=1}^n (Z_i - \bar{Z}) (Y_i - \bar{Y} + \bar{Y}) \\ &= \sum_{i=1}^n (Z_i - \bar{Z}) (Y_i - \bar{Y}) + \bar{Y} \sum_{i=1}^n (Z_i - \bar{Z}) \\ &= \sum_{i=1}^n Z_i (Y_i - \bar{Y}) - \bar{Z} \sum_{i=1}^n (Y_i - \bar{Y}) \\ &= \sum_{i=1}^n Z_i (Y_i - \bar{Y}), \end{aligned}$$

since $\sum_{i=1}^n (Z_i - \bar{Z}) = \sum_{i=1}^n (Y_i - \bar{Y}) = 0$.

2.

$$\sum_{i=1}^n Z_i (Y_i - \bar{Y}) = \sum_{i=1}^n Z_i Y_i - \bar{Y} \sum_{i=1}^n Z_i = n_1 \bar{Y}_1 - n_1 \bar{Y}.$$

3.

$$\begin{aligned} n_1 \bar{Y}_1 + n_0 \bar{Y}_0 &= \sum_{i=1}^n Z_i Y_i + \sum_{i=1}^n (1 - Z_i) Y_i \\ &= \sum_{i=1}^n \{Z_i Y_i + (1 - Z_i) Y_i\} \\ &= \sum_{i=1}^n Y_i. \end{aligned}$$

4.

$$\begin{aligned} n_1 (\bar{Y}_1 - \bar{Y}) &= n_1 \{ \bar{Y}_1 - (n_1 \bar{Y}_1 + n_0 \bar{Y}_0) / n \} \\ &= n_1 \{ (n_0 / n) \bar{Y}_1 - (n_0 / n) \bar{Y}_0 \} \\ &= (n_1 n_0 / n) (\bar{Y}_1 - \bar{Y}_0). \end{aligned}$$

5.

$$\begin{aligned} \hat{\beta}_1 &= \frac{\sum_{i=1}^n (Z_i - \bar{Z}) Y_i}{\sum_{i=1}^n (Z_i - \bar{Z}) X_i} \\ &= \frac{(n_1 n_0 / n) (\bar{Y}_1 - \bar{Y}_0)}{(n_1 n_0 / n) (\bar{X}_1 - \bar{X}_0)} \\ &= \frac{\bar{Y}_1 - \bar{Y}_0}{\bar{X}_1 - \bar{X}_0}. \end{aligned}$$

Problem 2. Consider the model

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + U_i, \quad (1)$$

where X_{1i} is an exogenous regressor and X_{2i} is an endogenous regressor. Assume that data are iid and conditions required for LLNs hold. For each of the following statements, indicate true or false, and explain your answer.

- (i) Let $\hat{\beta}_1$ denote the estimated coefficient on X_1 in the OLS regression of Y against a constant, X_1 , and X_2 . Since X_1 is exogenous, $\hat{\beta}_1$ consistently estimates β_1 .
- (ii) Let $\hat{\beta}_1$ denote the estimated coefficient on X_1 in the OLS regression of Y against a constant and X_1 . If $Cov(X_{1i}, X_{2i}) = 0$, then $\hat{\beta}_1$ consistently estimates β_1 .
- (iii) Consider the following IV estimator of β_2 that uses X_1 as an IV:

$$\hat{\beta}_2 = \frac{\sum_{i=1}^n (X_{1i} - \bar{X}_1) Y_i}{\sum_{i=1}^n (X_{1i} - \bar{X}_1) X_{2i}}.$$

If $Cov(X_{1i}, X_{2i}) \neq 0$ and $\beta_1 = 0$, then $\hat{\beta}_2$ consistently estimates β_2 .

Solution.

- (i) False. If X_1 and X_2 are correlated, $\hat{\beta}_1$ is inconsistent. Let \tilde{X}_{1i} denote fitted residuals in the regression of X_1 against a constant and X_2 :

$$\tilde{X}_{1i} = X_{1i} - \hat{\gamma}_0 - \hat{\gamma}_1 X_{2i},$$

where $\hat{\gamma}$'s denote the OLS estimators.

$$\begin{aligned}\hat{\beta}_1 &= \frac{\sum \tilde{X}_{1i} Y_i}{\sum \tilde{X}_{1i}^2} \\ &= \beta_1 + \frac{n^{-1} \sum \tilde{X}_{1i} U_i}{n^{-1} \sum \tilde{X}_{1i}^2}.\end{aligned}$$

Next,

$$n^{-1} \sum \tilde{X}_{1i} U_i = n^{-1} \sum X_{1i} U_i - \hat{\gamma}_0 n^{-1} \sum U_i - \hat{\gamma}_1 n^{-1} \sum X_{2i} U_i.$$

Since X_{1i} is exogenous,

$$n^{-1} \sum X_{1i} U_i \rightarrow_p 0.$$

We can also expect that

$$n^{-1} \sum U_i \rightarrow_p 0.$$

However, since X_{2i} is endogenous,

$$n^{-1} \sum X_{2i} U_i \rightarrow_p EX_{2i} U_i \neq 0.$$

Note also that

$$\hat{\gamma}_1 = \frac{n^{-1} \sum (X_{2i} - \bar{X}_2) X_{1i}}{n^{-1} \sum (X_{2i} - \bar{X}_2)^2} \rightarrow_p \frac{Cov(X_{2i}, X_{1i})}{Var(X_{2i})}.$$

Hence, if X_1 and X_2 are correlated, then $\hat{\beta}_1$ will be inconsistent.

- (ii) True. Write

$$\begin{aligned}Y_i &= \beta_0 + \beta_1 X_{1i} + V_i, \\ V_i &= \beta_2 X_{2i} + U_i.\end{aligned}$$

We have $Cov(X_{1i}, V_i) = \beta_2 Cov(X_{1i}, X_{2i}) + Cov(X_{1i}, U_i)$. Since X_1 is exogenous in the original model, $Cov(X_{1i}, U_i) = 0$. If $Cov(X_{1i}, X_{2i}) = 0$, then X_1 is uncorrelated with V in the new regression equation and, therefore, exogenous. Hence, $\hat{\beta}_1$ is a consistent estimator.

- (iii) True. Since $\beta_1 = 0$, X_1 is excluded from the structural equation. By the assumption, X_1 and U are uncorrelated. Since X_1 and X_2 are correlated, X_1 is a valid IV.

Problem 3. Suppose that the linear model

$$PS = \beta_0 + \beta_1 \text{Funds} + \beta_2 \text{Risk} + U$$

satisfies $E[U] = E[U \cdot \text{Funds}] = E[U \cdot \text{Risk}] = 0$. PS is the percentage of a person's savings invested in the stock market, Funds is the number of mutual funds that the person can choose from, and Risk is some measure of risk tolerance (larger Risk means the person has a higher tolerance for risk).

- (i) If Funds and Risk are positively correlated, does the slope coefficient in the simple regression of PS on Funds overestimate or underestimate β_1 , in large samples?
- (ii) We are unable to observe Risk directly, but we have data on the amount of life insurance a worker has, Insurance. Assume that Insurance is noisy measure of Risk, $\text{Insurance} = \text{Risk} + e$, with $E[e] = E[\text{Risk} \cdot e] = E[\text{Funds} \cdot e] = E[eU] = 0$. Will the OLS estimate of the coefficient on Funds in a regression of PS on Funds and Insurance be a consistent estimate of β_1 ?
- (iii) Suppose we also have data on how often a worker gambles, Gamble. Assume that Gamble is an independent noisy measure of Risk, $\text{Gamble} = \text{Risk} + v$, with $E[v] = E[vU] = E[ve] = E[\text{Risk} \cdot v] = E[\text{Funds} \cdot v] = 0$. Explain how we can consistently estimate β_1 using our data on PS, Funds, Insurance, and Gamble.

Solution.

- (i) The slope coefficient in the simple regression of PS on Funds converges in probability to $\beta_1 + \beta_2 \frac{\text{Cov}(\text{Funds}, \text{Risk})}{\text{Var}(\text{Funds})}$. We expect $\beta_2 \geq 0$. Therefore, the slope coefficient in the simple regression overestimates β_1 .
- (ii) No. Insurance is endogenous, due to measurement error.
- (iii) Run an IV regression, using Gamble as an instrument for Insurance.

Problem 4. Aggregate demand Q_D for a certain commodity is determined by its price P , aggregate income Y , and population, POP ,

$$Q_D = \beta_1 + \beta_2 P + \beta_3 Y + \beta_4 POP + U^D$$

and aggregate supply is given by

$$Q_S = \alpha_1 + \alpha_2 P + U^S$$

where U_D and U_S are independently distributed error terms: U_D and U_S are independent from all other variables and they are also independent from each other. Remember that the quantity and the price are determined simultaneously in the equilibrium $Q_S = Q_D = Q$. We observe only the equilibrium values Q so that the observed price must satisfy the equation (demand = supply):

$$\beta_1 + \beta_2 P + \beta_3 Y + \beta_4 POP + U^D = \alpha_1 + \alpha_2 P + U^S.$$

- (i) Show that the OLS (ordinary least squares) estimator of α_2 will be inconsistent if OLS is used to fit the supply equation.
- (ii) Show that a consistent estimator of α_2 is

$$\tilde{\alpha}_2 = \frac{\sum_{i=1}^n (Y_i - \bar{Y}) (Q_i - \bar{Q})}{\sum_{i=1}^n (Y_i - \bar{Y}) (P_i - \bar{P})}.$$

$$(\bar{Y} = n^{-1} \sum_{i=1}^n Y_i, \bar{Q} = n^{-1} \sum_{i=1}^n Q_i, \bar{P} = n^{-1} \sum_{i=1}^n P_i.)$$

Solution.

- (i) The reduced form equation (which expresses P as a function of the explanatory variables and the error terms) for P is

$$P = \frac{1}{\alpha_2 - \beta_2} (\beta_1 - \alpha_1 + \beta_3 Y + \beta_4 POP + U^D - U^S).$$

Therefore in the supply equation

$$Q_S = \alpha_1 + \alpha_2 P + U^S,$$

P is correlated with U^S . The OLS estimator is

$$\begin{aligned}\hat{\alpha}_2^{OLS} &= \frac{\sum_{i=1}^n (P_i - \bar{P}) (Q_i - \bar{Q})}{\sum_{i=1}^n (P_i - \bar{P})^2} \\ &= \alpha_2 + \frac{\sum_{i=1}^n (P_i - \bar{P}) (U_i^S - \bar{U}^S)}{\sum_{i=1}^n (P_i - \bar{P})^2} \\ &\rightarrow_p \alpha_2 + \frac{\text{Cov}(P_i, U_i^S)}{\text{Var}(P_i)}\end{aligned}$$

and

$$\begin{aligned}\text{Cov}(P_i, U_i^S) &= \text{Cov}\left(\frac{1}{\alpha_2 - \beta_2} (\beta_1 - \alpha_1 + \beta_3 Y_i + \beta_4 POP_i + U_i^D - U_i^S), U_i^S\right) \\ &= -\frac{1}{\alpha_2 - \beta_2} \text{Var}(U_i^S) \\ &\neq 0\end{aligned}$$

assuming that Y and POP are exogenous and so $\text{Cov}(U^S, Y) = \text{Cov}(U^S, POP) = 0$. We are told that U^S and U^D are distributed independently, so that $\text{Cov}(U^S, U^D) = 0$.

- (ii) The instrument variable estimator is

$$\begin{aligned}\hat{\alpha}_2^{IV} &= \frac{\sum_{i=1}^n (Y_i - \bar{Y}) (Q_i - \bar{Q})}{\sum_{i=1}^n (Y_i - \bar{Y}) (P_i - \bar{P})} \\ &= \alpha_2 + \frac{\sum_{i=1}^n (Y_i - \bar{Y}) (U^S - \bar{U}^S)}{\sum_{i=1}^n (Y_i - \bar{Y}) (P_i - \bar{P})} \\ &\rightarrow_p \alpha_2 + \frac{\text{Cov}(Y_i, U_i^S)}{\text{Cov}(P_i, Y_i)}.\end{aligned}$$

The desired result follows from the assumptions $\text{Cov}(Y_i, U_i^S) = 0$ and $\text{Cov}(P_i, Y_i) \neq 0$.

Problem 5. In an econometric model, we say that a parameter is identified if we can recover its value perfectly given the joint distribution of the observable variables. Suppose that (Y, X) is the observable variables and U is the unobservable variable.

- (i) Suppose that $Y = \beta_0 + \beta_1 X + U$ and $E[U] = E[XU] = 0$. Show that β_1 is identified. I.e., if you know the joint distribution of (Y, X) , how do you determine the value of the parameter β_1 ?
- (ii) Suppose that Y is binary and $Y = 1(\beta_0 + \beta_1 X \geq U)$ and U is a standard normal ($N(0, 1)$) random variable that is independent of X . If you know the joint distribution of (Y, X) , how do you determine the value of the parameter β_1 ? Hint: $E[Y | X] = E[1(\beta_0 + \beta_1 X \geq U) | X] = \Phi(\beta_0 + \beta_1 X)$, where Φ is the standard normal CDF.

Solution. Take

$$\begin{aligned} \text{Cov}[Y, X] &= \text{Cov}[\beta_0 + \beta_1 X + U, X] = \text{Cov}[\beta_1 X + U, X] = \\ &= \beta_1 \text{Cov}[X, X] + \text{Cov}[U, X] = \beta_1 \text{Var}[X]. \end{aligned}$$

Therefore, $\beta_1 = \text{Cov}[Y, X] / \text{Var}[X]$. This quantity can be recovered if you know the joint distribution of (Y, X) .

Similarly, $E[Y | X] = \Phi(\beta_0 + \beta_1 X)$ gives $\beta_0 + \beta_1 X = \Phi^{-1}(E[Y | X])$, where Φ^{-1} is the inverse function of Φ (Φ is strictly increasing). Then,

$$\text{Cov}[\Phi^{-1}(E[Y | X]), X] = \text{Cov}[\beta_0 + \beta_1 X, X] = \beta_1 \text{Var}[X].$$

Therefore, $\beta_1 = \text{Cov}[\Phi^{-1}(E[Y | X]), X] / \text{Var}[X]$. This quantity can be recovered if you know the joint distribution of (Y, X) .

Problem 6. Let (X_i, Y_i) , $i = 1, \dots, n$ be an i.i.d. random sample where $Y_i \geq 0$ and $X_i \geq 0$ is a discrete random variable for all i . The conditional density of Y given X belong to the family:

$$f_{Y|X}(y|x, \lambda) = \frac{\lambda \exp(-\lambda y) (\lambda y)^x}{x!},$$

$y \geq 0$, $\lambda > 0$, i.e., the conditional density of Y given X is $f_{Y|X}(\cdot | \cdot, \lambda_*)$ for some $\lambda_* > 0$. Write the likelihood function for estimating λ_* . Provide the maximum likelihood estimator for λ_* as a solution of an equation. Give the asymptotic distribution for the maximum likelihood estimator, i.e. find the asymptotic variance V_{ML} of

$$\sqrt{n}(\hat{\lambda}_{ML} - \lambda_*) \rightarrow_d N(0, V_{ML}).$$

Suggest a consistent estimator of V_{ML} .

Solution. The log likelihood function is

$$\ell(\lambda) = \sum_{i=1}^N \log \left(\frac{\lambda \exp(-\lambda Y_i) (\lambda Y_i)^{X_i}}{X_i!} \right) = \sum_{i=1}^N \{ \log(\lambda) - \lambda Y_i + X_i \log(\lambda Y_i) - \log(X_i!) \}.$$

Take derivative with respect to λ :

$$\frac{d\ell(\lambda)}{d\lambda} = \sum_{i=1}^N \left\{ \frac{1}{\lambda} - Y_i + \frac{X_i}{\lambda} \right\}.$$

The maximum likelihood estimator is the solution to the first order condition

$$\sum_{i=1}^N \left\{ \frac{1}{\lambda} - Y_i + \frac{X_i}{\lambda} \right\} = 0.$$

Solving this gives

$$\hat{\lambda}_{ML} = \frac{N + \sum_{i=1}^N X_i}{\sum_{i=1}^N Y_i}.$$

To derive the asymptotic distribution of $\hat{\lambda}^{MLE}$, take

$$\frac{d^2 \ell(\lambda)}{d\lambda^2} = -\frac{\sum_{i=1}^N 1 + X_i}{\lambda^2}.$$

The asymptotic variance V_{ML} is

$$V_{ML} = -N \left(E \left[\frac{d^2 \ell(\lambda)}{d\lambda^2} \Big|_{\lambda=\lambda_*} \right] \right)^{-1} = -N \left(E \left[-\frac{\sum_{i=1}^N 1 + X_i}{\lambda_*^2} \right] \right)^{-1} = \frac{\lambda_*^2}{1 + E[X_i]}.$$

A consistent estimator is just

$$\frac{\hat{\lambda}_{ML}^2}{1 + \frac{1}{N} \sum_{i=1}^N X_i}.$$

Since the sample is i.i.d. $\frac{1}{N} \sum_{i=1}^N X_i \rightarrow_p E[X_i]$. It follows from this result, consistency of $\hat{\lambda}_{ML}$, and Slutsky's lemma that

$$\frac{\hat{\lambda}_{ML}^2}{1 + \frac{1}{N} \sum_{i=1}^N X_i} \rightarrow_p \frac{\lambda_*^2}{1 + E[X_i]}.$$

Problem 7. Define a density function

$$f(x | \theta) = \begin{cases} \left(1 + \frac{1-2\theta}{\theta-1}\right) x^{\frac{1-2\theta}{\theta-1}} & x \in (0, 1) \\ 0 & x \notin (0, 1), \end{cases}$$

where $0 < \theta < 1$ is a parameter. X_1, \dots, X_n is an independent and identically distributed sample with true density $f(\cdot | \theta_*)$ for some θ_* .

(i) Show that $f(\cdot | \theta)$ is a probability density function, for all $0 < \theta < 1$.

(ii) Show that $\theta_* = \int_0^1 x f(x | \theta_*) dx$. I.e., in this parametrization, θ_* is also the population mean.

Solution.

(i) Compute

$$\int_0^1 f(x | \theta) dx = \left(1 + \frac{1-2\theta}{\theta-1}\right) \int_0^1 x^{\frac{1-2\theta}{\theta-1}} dx = \left(1 + \frac{1-2\theta}{\theta-1}\right) \frac{1}{1 + \frac{1-2\theta}{\theta-1}} x^{1 + \frac{1-2\theta}{\theta-1}} \Big|_0^1 = 1.$$

Therefore, $f(x | \theta) \geq 0$ and $\int_0^1 f(x | \theta) dx = 1$.

(ii) Compute

$$\int_0^1 x f(x | \theta_*) dx = \left(1 + \frac{1 - 2\theta_*}{\theta_* - 1}\right) \int_0^1 x \cdot x^{\frac{1-2\theta_*}{\theta_*-1}} dx = \left(1 + \frac{1 - 2\theta_*}{\theta_* - 1}\right) \frac{1}{1 - \frac{\theta_*}{\theta_*-1}} x^{1 - \frac{\theta_*}{\theta_*-1}} \Big|_0^1 = \theta_*.$$

The method of moment estimator: $n^{-1} \sum_{i=1}^n X_i$.

- The log-maximum likelihood function is

$$\log L(\theta; X_1, \dots, X_n) = n \log \left(\frac{\theta}{1 - \theta} \right) + \frac{1 - 2\theta}{\theta - 1} \sum_{i=1}^n \log(X_i).$$

Differentiating with respect to θ :

$$\frac{\partial \log L}{\partial \theta} = \frac{n}{\theta(1 - \theta)} + \frac{1}{(1 - \theta)^2} \sum_{i=1}^n \log(X_i).$$

Solving the first order condition, the maximum likelihood estimator is

$$\hat{\theta} = \frac{n}{n - \sum_{i=1}^n \log(X_i)},$$

which is different from the method of moments estimator.