

# IIA-3 Econometrics: Supervision 3

Emre Usenmez

Christmas Break 2024

Very grateful to Dr Oleg Kitov and Dr Clive Lawson for the very informative stylized answers to previous iterations of the supervision questions.

## FACULTY QUESTIONS

### QUESTION 1

Consider the following bivariate linear regression

$$y = \alpha + T\beta + u$$

where  $T$  is a binary treatment regressor,  $\alpha$  and  $\beta$  are unknown parameters, and  $u$  is an error term.

(a) Describe in two sentences an empirical, real-life example where such an equation might arise.

**Answer:** We can think of  $T$  as "graduated from university" and  $y$  as "annual earning after 10 years of graduation."

---

(b) Why might  $u$  be heteroskedastic in your example.

**Answer:** The variance of earnings will likely to be smaller across people who did not graduate from a university compared to those who did it. This may be because those who did not go to university are less likely to be in the professions such as lawyers or doctors, and more likely to be in lower-paying jobs, or unemployed, or out of labor force.

---

(c) Why might  $T$  be endogenous in your example?

**Answer:** Broadly, variables that are correlated with the error term are called *endogeneous variables*, and those that are uncorrelated with the error term are called *exogeneous variables*.<sup>1</sup>

Thus the question is asking us to consider some of the reasons as to why  $T$  might be correlated with the error term. There are certainly nonnegligible number of high earners who either never went to a university or dropped out. There may be omitted variable or even simultaneity is possible.

Let's consider what the implications of of endogeneity are for the OLS estimator of  $\beta$ .

Variable  $T$  would be endogenous if  $\mathbb{E}(u|T) \neq 0$ . Endogeneity would imply that  $Cov(T, u) \neq 0$ .

We can first look at whether it is biased. For that, we need to use the law of iterated expectations whereby

$$\mathbb{E}(\hat{\beta}^{OLS}) = \mathbb{E}[\mathbb{E}(\hat{\beta}^{OLS} | T_1, \dots, T_n)]$$

The OLS estimator of  $\beta$  would be:

$$\begin{aligned} \mathbb{E}(\hat{\beta}^{OLS} | T_1, \dots, T_n) &= \mathbb{E}\left(\frac{\widehat{Cov}(T_i, Y_i)}{\widehat{Var}(T_i)} \mid T_1, \dots, T_n\right) = \mathbb{E}\left(\frac{\hat{\sigma}_{TY}}{\hat{\sigma}_{TT}} \mid T_1, \dots, T_n\right) \\ &= \mathbb{E}\left(\frac{\sum_{i=1}^n (T_i - \bar{T})(Y_i - \bar{Y})}{\sum_{i=1}^n (T_i - \bar{T})^2} \mid T_1, \dots, T_n\right) \\ &= \mathbb{E}\left(\frac{\sum_{i=1}^n (T_i - \bar{T})((\alpha + \beta T_i + u_i) - (\alpha + \beta \bar{T} + \bar{u}))}{\sum_{i=1}^n (T_i - \bar{T})^2} \mid T_1, \dots, T_n\right) \\ &= \mathbb{E}\left(\frac{\sum_{i=1}^n (T_i - \bar{T})(\beta(T_i - \bar{T}) + u_i - \bar{u})}{\sum_{i=1}^n (T_i - \bar{T})^2} \mid T_1, \dots, T_n\right) \\ &= \mathbb{E}\left(\frac{\sum_{i=1}^n \beta(T_i - \bar{T})^2 + \sum_{i=1}^n (T_i - \bar{T})(u_i - \bar{u})}{\sum_{i=1}^n (T_i - \bar{T})^2} \mid T_1, \dots, T_n\right) \\ &= \mathbb{E}\left(\beta + \frac{\sum_{i=1}^n (T_i - \bar{T})(u_i - \bar{u})}{\sum_{i=1}^n (T_i - \bar{T})^2} \mid T_1, \dots, T_n\right) \end{aligned}$$

---

<sup>1</sup>See Chapter 12: Instrumental Variables Regression p.428 in Stock J H, and Watson M W (2020) Introduction to Econometrics, 4<sup>th</sup> Global Ed, Pearson; and Section 8.5: Instrumental Variables in Dougherty C (2016) Introduction to Econometrics 5<sup>th</sup> ed, OUP

$$\begin{aligned}
&= \mathbb{E} \left( \beta + \frac{\sum_{i=1}^n (T_i - \bar{T}) u_i - \bar{u} \sum_{i=1}^n (T_i - \bar{T})}{\sum_{i=1}^n (T_i - \bar{T})^2} \middle| T_1, \dots, T_n \right) \\
&= \mathbb{E} \left( \beta + \frac{\sum_{i=1}^n (T_i - \bar{T}) u_i - \bar{u} \left( \sum_{i=1}^n T_i - n\bar{T} \right)}{\sum_{i=1}^n (T_i - \bar{T})^2} \middle| T_1, \dots, T_n \right) \\
&= \mathbb{E} \left( \beta + \frac{\sum_{i=1}^n (T_i - \bar{T}) u_i - \bar{u} (n\bar{T} - n\bar{T})}{\sum_{i=1}^n (T_i - \bar{T})^2} \middle| T_1, \dots, T_n \right) \\
&= \mathbb{E} \left( \beta + \frac{\sum_{i=1}^n (T_i - \bar{T}) u_i}{\sum_{i=1}^n (T_i - \bar{T})^2} \middle| T_1, \dots, T_n \right) \\
&= \beta + \frac{\sum_{i=1}^n (T_i - \bar{T}) \mathbb{E}(u_i \mid T_1, \dots, T_n)}{\mathbb{E} \left( \sum_{i=1}^n (T_i - \bar{T})^2 \mid T_1, \dots, T_n \right)}
\end{aligned}$$

Notice that since  $\mathbb{E}(u|T) \neq 0$ , the numerator of this last expression is also nonzero. That is,  $\sum_{i=1}^n (T_i - \bar{T}) \mathbb{E}(u_i \mid T_1, \dots, T_n) \neq 0$ . Therefore the expectation of this expectation is also not equal to  $\beta$ :

$$\mathbb{E}(\hat{\beta}^{OLS}) = \mathbb{E} \left[ \mathbb{E}(\hat{\beta}^{OLS} \mid T_1, \dots, T_n) \right] = \mathbb{E} \left[ \mathbb{E} \left( \beta + \frac{\sum_{i=1}^n (T_i - \bar{T}) u_i}{\sum_{i=1}^n (T_i - \bar{T})^2} \middle| T_1, \dots, T_n \right) \right] \neq \beta$$

which means the OLS estimator is *not* unbiased.

We can also check for consistency by examining the probability limit of this expression as  $n$  tends towards infinity. For that, we can rewrite the OLS estimator as:

$$\hat{\beta}^{OLS} = \beta + \frac{\frac{1}{n} \sum_{i=1}^n (T_i - \bar{T}) u_i}{\frac{1}{n} \sum_{i=1}^n (T_i - \bar{T})^2}$$

Using the law of large numbers, we can see that as  $n \rightarrow \infty$

$$\frac{1}{n} \sum_{i=1}^n (T_i - \bar{T}) u_i \xrightarrow{p} \mathbb{E}[(T_i - \bar{T}) u_i] = \text{Cov}(T_i, u_i) \neq 0$$

and

$$\frac{1}{n} \sum_{i=1}^n (T_i - \bar{T})^2 \xrightarrow{p} \mathbb{E}[(T_i - \bar{T})^2] = \text{Var}(T_i) = \sigma_T^2 < \infty$$

Note that  $\text{Var}(T_i) = \sigma_T^2 < \infty$  is an additional assumption.

Since  $\text{Cov}(T_i, u_i) \neq 0$ , the OLS estimator as  $n \rightarrow \infty$  (using Slutsky's theorem):

$$\hat{\beta}^{OLS} \xrightarrow{p} \beta + \frac{\text{Cov}(T_i, u_i)}{\text{Var}(T_i)} \neq \beta$$

which means the OLS estimator is not only biased but also inconsistent for  $\beta$ .

(d) Suppose a single instrument  $z$  is available. Show that the population coefficient  $\beta$  satisfies

$$\beta = \frac{\text{Cov}(z, y)}{\text{Cov}(z, T)}$$

where  $\text{Cov}(z, y)$  and  $\text{Cov}(z, T)$  denote, respectively, the population covariance between  $z$  and  $y$ , and  $z$  and  $T$ . How can you use this information to obtain a consistent estimate of  $\beta$ ?

**Answer:** Instrument  $z$  needs to satisfy the following conditions:

- *Instrument relevance:*  $z$  must have non-trivial explanatory power for  $T$ , namely  $\text{Cov}(z, T) \neq 0$ .
- *Instrument exogeneity:*  $z$  must affect  $Y$  only through its influence on  $T$  and not in any other way. That is,  $z$  must be exogenous with respect to  $u$  in regression  $y = \alpha + \beta T + u$ . Formally,  $\mathbb{E}(u|z) = 0$ . This is why it is said " $z$  is exogenous in  $y = \alpha + \beta T + u$ ". Exogeneity of instrument  $z$  implies that  $\text{Cov}(z, u) = 0$ .

In the context of omitted variables, instrument exogeneity means that  $z$  should be uncorrelated with the omitted variables, i.e.  $\text{Cov}(z, u) = 0$ , and  $z$  should be related, positively or negatively, to the endogenous explanatory variable  $T$ , i.e.  $\text{Cov}(z, T) \neq 0$ .<sup>2</sup>

The underlying reasoning is that if an instrument is relevant, then variation in that instrument  $z$  is related to variation in  $T$ , and if it is also exogenous, then that part of the variation of  $T$  captured by  $z$  is exogenous. Therefore, an instrument that is relevant and exogenous can capture movements in  $T$  that are exogenous. This exogenous variation can in turn be used to estimate the population coefficient  $\beta$ .<sup>3</sup>

These conditions serve to *identify* the parameter  $\beta$ . In this context, *identification of a parameter* means that we can write  $\beta$  in terms of population moments that can be estimated using a sample of data.

To write  $\beta$  in terms of population covariances we use  $y = \alpha + \beta T + u$ :

$$\text{Cov}(z, y) = \text{Cov}(z, \alpha + \beta T + u) = \beta \text{Cov}(z, T) + \text{Cov}(z, u)$$

<sup>2</sup>see Section 15-1: Omitted Variables in a Simple Regression Model in Wooldridge J M (2021) Introductory Econometrics: A Modern Approach, 7<sup>th</sup> ed, Cengage

<sup>3</sup>see Section 12.1: The IV Estimator with a Single Regressor and a Single Instrument in Stock and Watson (2020, 4<sup>th</sup> ed.).

Since instrument exogeneity condition assumes that  $Cov(z, u) = 0$  then  $Cov(z, y) = \beta Cov(z, T)$ . Rearranging this gives:

$$\beta = \frac{Cov(z, y)}{Cov(z, T)}$$

as desired. Notice that this only holds if instrument relevance also holds, since this expression would fail if  $Cov(z, T) = 0$ . What this expression is telling us is that  $\beta$  is identified by the ratio of population covariance between  $z$  and  $y$  to population covariance between  $z$  and  $T$ .

Given a random sample, we estimate the population quantities by the sample analogs:

$$\hat{\beta}^{IV} = \frac{\frac{1}{n} \sum_{i=1}^n (z_i - \bar{z})(y_i - \bar{y})}{\frac{1}{n} \sum_{i=1}^n (z_i - \bar{z})(T_i - \bar{T})} = \frac{\sum_{i=1}^n (z_i - \bar{z})(y_i - \bar{y})}{\sum_{i=1}^n (z_i - \bar{z})(T_i - \bar{T})}.$$

With a sample data on  $T$ ,  $y$ , and  $z$  we can obtain the IV estimator above. The IV estimator for the intercept  $\alpha$  is  $\alpha = \bar{y} - \hat{\beta}^{IV} \bar{T}$ . Also notice that when  $z = T$ , we get the OLS estimator of  $\beta$ . That is, when  $T$  is exogeneous, it can be used as its own IV, and the IV estimator is then identical to the OLS estimator.

A similar set of steps we used in part (c) will show that IV estimator is consistent for  $\beta$ . That is,  $\text{plim}_{n \rightarrow \infty}(\hat{\beta}) = \beta$ .

Note that, an important feature of IV estimator is that when  $T$  and  $u$  are in fact correlated, and thus instrumental variables estimation is actually needed, it is essentially never unbiased. This means, in small samples, the IV estimator can have a substantial bias, which is one reason why large samples are preferred.

---

(e) Can you give an example of an instrument in your example? Argue why it might be a sensible IV.

**Answer:** Distance from nearest college can be an example of an instrument, where  $z = 1$  if individual lived near college and 0 otherwise. This may be violated for a number of reasons, though; for e.g. if wealthy parents choose to live near college. This would mean that  $z$  is correlated with unobserved factors that also affect wage, our  $y$ . For any example, exogeneity and relevance conditions need to be checked.

## QUESTION 2

Consider the following wage equation that explicitly recognizes that ability affects  $\log(wage)$

$$\log(wage) = \alpha + \beta_1 educ + \beta_2 ability + u$$

The above model shows explicitly that we would like to hold ability fixed when measuring the returns on education. Assuming that the primary interest is in obtaining a reliable estimate of the slope parameters  $\beta_1$ , and that there is no direct measurement for ability, explain how you would do this using a method based upon a proxy variable and an IV estimator. In doing so evaluate the following statement:

*“whilst IQ is a good candidate as a proxy for variable for ability, it is not a good instrumental variable for educ.”*