

# Advanced Econometrics, HT, Problem Set 5

## Question 1

We observe a random sample  $(Y_i, X_i, Z_i)$ ,  $i = 1, \dots, n$ , where  $n$  is the sample size, and  $Y_i$ ,  $X_i$  and  $Z_i$  are three scalar variables with finite second moments. For simplicity we assume that  $\mathbb{E}Y_i = \mathbb{E}X_i = \mathbb{E}Z_i = 0$ . We assume that

$$Y_i = X_i\beta + U_i$$

holds, where  $U_i$  is a mean zero error term and  $\beta$  is the parameter of interest. We also define  $\gamma_{xx} = \text{Var}(X_i)$ ,  $\gamma_{xy} = \text{Cov}(X_i, Y_i)$ ,  $\gamma_{xz} = \text{Cov}(X_i, Z_i)$  and  $\gamma_{yz} = \text{Cov}(Y_i, Z_i)$ . We assume that  $\gamma_{xx} > 0$ .

- (a) Write down consistent estimators for  $\gamma_{xx}$ ,  $\gamma_{xy}$ ,  $\gamma_{xz}$ , and  $\gamma_{yz}$ . Explain why they are consistent.
- (b) Assume that  $X_i$  is exogenous, i.e. that  $\mathbb{E}(X_i U_i) = 0$ . Use this assumption to express  $\beta$  as a function of  $\gamma_{xx}$  and  $\gamma_{xy}$ . Use this expression for  $\beta$  and your result in (a) to provide a consistent estimator for  $\beta$ . Prove that your estimator for  $\beta$  is indeed consistent.
- (c) Assume that  $X_i$  is endogenous, i.e.  $\mathbb{E}(X_i U_i) \neq 0$ , but that the instrument  $Z_i$  satisfies the exclusion restriction  $\mathbb{E}(Z_i U_i) = 0$  and relevance assumption  $\mathbb{E}(X_i Z_i) \neq 0$ . Use these two assumptions on  $Z_i$  to express  $\beta$  as a function of  $\gamma_{xz}$  and  $\gamma_{yz}$ . Use this expression for  $\beta$  and to write down a consistent estimator for  $\beta$ , which we denote by  $\widehat{\beta}_{\text{IV}}$ .
- (d) Under the assumptions in (c), show that  $\widehat{\beta}_{\text{IV}}$  is asymptotically normal, i.e. show that  $\sqrt{n}(\widehat{\beta}_{\text{IV}} - \beta) \Rightarrow \mathcal{N}(0, \Sigma_{\text{IV}})$  as  $n \rightarrow \infty$ . Provide a formula for the asymptotic variance  $\Sigma_{\text{IV}}$ .

## Question 2

Consider the model with outcome  $Y \in \mathbb{R}$ , binary treatment  $D \in \{1, 0\}$ , and binary instrument  $Z \in \{1, 0\}$ . We observe the sample  $(Y_i, D_i, Z_i)$ ,  $i = 1, \dots, n$ . For  $z, d \in \{1, 0\}$

let

$$N_{zd} = \sum_{i=1}^n \mathbb{1}\{Z_i = z \& D_i = d\}$$

be the number of observations with instrument value  $z$  and treatment status  $d$ . We have  $n = N_{11} + N_{10} + N_{01} + N_{00}$ . Similarly, let

$$\bar{Y}_{zd} = \frac{1}{N_{zd}} \sum_{i=1}^n \mathbb{1}\{Z_i = z \& D_i = d\} Y_i$$

be the average outcome of units with instrument value  $z$  and treatment status  $d$ . Let  $\hat{\gamma}_1$  and  $\hat{\delta}_1$  be the OLS estimates for  $\gamma_1$  and  $\delta_1$  in the first stage and reduced form regressions

$$\begin{aligned} D_i &= \gamma_0 + \gamma_1 Z_i + \nu_i, \\ Y_i &= \delta_0 + \delta_1 Z_i + v_i. \end{aligned}$$

Let  $\hat{\beta} = \hat{\delta}_1 / \hat{\gamma}_1$  be the 2SLS estimator for the parameter  $\beta$  in the structural equation

$$Y_i = \alpha + \beta D_i + \epsilon_i. \quad (1)$$

- (a) Provide expressions for  $\hat{\gamma}_1$ ,  $\hat{\delta}_1$ ,  $\hat{\beta}$  as functions of  $(N_{zd}, \bar{Y}_{zd} : z, d \in \{1, 0\})$ .
- (b) Show how the formula for  $\hat{\beta}$  simplifies for the case that  $N_{01} = 0$  and  $N_{10} = 0$ . This is the case where we only have compliers (i.e. we have  $D = Z$ , or  $D(z) = z$  in potential outcome notation). Interpret the formula for  $\hat{\beta}$  in that case.
- (c) Show how the formula for  $\hat{\beta}$  also simplifies for the case that  $N_{00} = 0$  and  $N_{11} = 0$ . This is the case where we only have defiers (i.e. we have  $D = 1 - Z$ , or  $D(z) = 1 - z$  in potential outcome notation). Interpret the formula for  $\hat{\beta}$  in that case.
- (d) What happens to  $\hat{\beta}$  if  $N_{10} = 0$  and  $N_{00} = 0$ ? This is the case when we only have always takers (i.e. we have  $D = 1$  for all units, or  $D(z) = 1$  in potential outcome notation).

Now, assume that we initially have three separate datasets. The first dataset (labeled by a superscript “at”) contains only always takers and therefore satisfies

$$N_{11}^{\text{at}} > 0, \quad N_{10}^{\text{at}} = 0, \quad N_{01}^{\text{at}} > 0, \quad N_{00}^{\text{at}} = 0.$$

The second dataset (labeled by a superscript “nt”) contains only never takers and therefore satisfies

$$N_{11}^{\text{nt}} = 0, \quad N_{10}^{\text{nt}} > 0, \quad N_{01}^{\text{nt}} = 0, \quad N_{00}^{\text{nt}} > 0.$$

The third dataset (labeled by a superscript “c”) contains only compliers and therefore satisfies

$$N_{11}^{\text{c}} > 0, \quad N_{10}^{\text{c}} = 0, \quad N_{01}^{\text{c}} = 0, \quad N_{00}^{\text{c}} > 0.$$

Furthermore, let  $n^{\text{at}} = N_{11}^{\text{at}} + N_{10}^{\text{at}} + N_{01}^{\text{at}} + N_{00}^{\text{at}}$  be the total number of always takers, and similarly let  $n^{\text{nt}}$  be the total number of never takes, and  $n^{\text{c}}$  be the total number of compliers. To simplify the computation we assume that

$$\frac{N_{11}^{\text{at}} + N_{10}^{\text{at}}}{n^{\text{at}}} = \frac{N_{11}^{\text{nt}} + N_{10}^{\text{nt}}}{n^{\text{nt}}} = \frac{N_{11}^{\text{c}} + N_{10}^{\text{c}}}{n^{\text{c}}}, \quad (2)$$

that is, the fraction of observations with  $Z_i = 1$  is assumed to be exactly identical for the always takers, never takes, and compliers, respectively.

Finally, let  $\bar{Y}_{zd}^{\text{at}}, \bar{Y}_{zd}^{\text{nt}}, \bar{Y}_{zd}^{\text{c}}$  be the average outcomes of units with instrument value  $z$  and treatment status  $d$  for the always takers, never takes, and compliers, respectively, and assume that

$$\bar{Y}_{0d}^{\text{at}} = \bar{Y}_{1d}^{\text{at}}, \quad \bar{Y}_{0d}^{\text{nt}} = \bar{Y}_{1d}^{\text{nt}}. \quad (3)$$

The total dataset that we actually observe combines all the observations of these three datasets such that  $N_{zd} = N_{zd}^{\text{at}} + N_{zd}^{\text{nt}} + N_{zd}^{\text{c}}$ . However, in this combined dataset we do not actually know whether a particular observation is an always taker, a never taker or a complier. The definition of the 2SLS estimator  $\hat{\beta} = \hat{\delta}_1 / \hat{\gamma}_1$  is unchanged. In particular, your result from part (a) is still applicable.

(e) Show that for the combined dataset we have

$$\hat{\beta} = \bar{Y}_{11}^{\text{c}} - \bar{Y}_{00}^{\text{c}},$$

that is,  $\hat{\beta}$  obtained from the whole dataset is identical to the 2SLS estimator that would be obtained from the complier dataset only.

**Remark:** Let

$$\hat{p}^{\text{at}} = \frac{N_{11}^{\text{at}} + N_{10}^{\text{at}}}{n^{\text{at}}}, \quad \hat{p}^{\text{nt}} = \frac{N_{11}^{\text{nt}} + N_{10}^{\text{nt}}}{n^{\text{nt}}}, \quad \hat{p}^{\text{c}} = \frac{N_{11}^{\text{c}} + N_{10}^{\text{c}}}{n^{\text{c}}},$$

be the fraction of observations with  $Z_i = 1$  for the the always takers, never takes, and compliers, respectively. In display (2) we assumed that  $\hat{p}^{\text{at}} = \hat{p}^{\text{nt}} = \hat{p}^{\text{c}}$ . In an actual dataset this will not be exactly the case, but under random assignment we expect that  $\hat{p}^{\text{at}} \approx \hat{p}^{\text{nt}} \approx \hat{p}^{\text{c}}$  in large datasets. Similarly, (3) will only hold approximately in an actual dataset under random assignment. If (2) and (3) only hold approximately, then the result in (e) will also only hold approximately in an actual dataset.