

This section presented a simple econometric model with a supply equation and a demand equation—equations (2a) and (2b). The source of endogeneity bias was identified: disturbance terms turn up in formulas (3ab) for  $Q_t$  and  $P_t$ . (These “reduced form” equations are of no further interest here, although they may be helpful in other contexts.) The way to get around endogeneity bias is to estimate equations (2a) and (2b) by instrumental variables rather than OLS. This new technique will be explained in sections 2 and 3. Section 7.4 discussed endogeneity bias in a different kind of model, with a binary response variable.

### Exercise set A

1. In equation (1a), should  $a_1$  be positive or negative? What about  $a_2, a_3$ ?
2. In equation (1b), should  $b_1$  be positive or negative? What about  $b_2, b_3$ ?
3. In the butter model of this section:
  - (a) Does the law of supply and demand hold true?
  - (b) Is the supply curve concave? strictly concave?
  - (c) Is the demand curve convex? strictly convex?
 (Economists prefer log linear specifications. . . .)
4. An economist wants to use the butter model to determine how farmers will respond to price controls. Which of the following equations is the most relevant—(2a), (2b), (3a), (3b)? Explain briefly.

## 9.2 Instrumental variables

We begin with a slightly abstract linear model

$$(4) \quad Y = X\beta + \delta,$$

where  $Y$  is an observable  $n \times 1$  random vector,  $X$  is an observable  $n \times p$  random matrix, and  $\beta$  is an unobservable  $p \times 1$  parameter vector. The  $\delta_i$  are IID with mean 0 and finite variance  $\sigma^2$ ; they are unobservable random errors. This is the standard regression model, except that  $X$  is endogenous, i.e.,  $X$  and  $\delta$  are dependent. Conditional on  $X$ , the OLS estimates are biased by  $(X'X)^{-1}X'E(\delta|X)$ : see (4.9). This is *simultaneity bias*.

We can explain the bias another way. In the OLS model, we could have obtained the estimator as follows: multiply both sides of (4) by  $X'$ , drop  $X'\delta$  because it's small— $E(X'\delta) = 0$ —and solve the resulting  $p$  equations for the  $p$  unknown components of  $\beta$ . Here, however,  $E(X'\delta) \neq 0$ .

To handle simultaneity bias, economists and other social scientists would estimate (4) using *instrumental-variables regression*, also called *two-stage*

*least squares*: the acronyms are IVLS and IISLS (or 2SLS, if you prefer Arabic numerals). The method requires an  $n \times q$  matrix of *instrumental* or *exogenous* variables, with  $n > q \geq p$ . The matrix will be denoted  $Z$ . The matrices  $Z'X$  and  $Z'Z$  need to be of full rank,  $p$  and  $q$  respectively. If  $q > p$ , the system is *over-identified*. If  $q = p$ , the system is *just-identified*. If  $q < p$ , the case which is excluded by assuming  $q \geq p$ , the system is *under-identified*—parameters will not be identifiable (section 7.2). Let's make a cold list of the assumptions.

- (i)  $X$  is  $n \times p$  and  $Z$  is  $n \times q$  with  $n > q \geq p$ .
- (ii)  $Z'X$  and  $Z'Z$  have full rank,  $p$  and  $q$  respectively.
- (iii)  $Y = X\beta + \delta$ .
- (iv) The  $\delta_i$  are IID, with mean 0 and variance  $\sigma^2$ .
- (v)  $Z$  is exogenous, i.e.,  $Z \perp\!\!\!\perp \delta$ .

Assumptions (i) and (ii) are easy to check from the data. The others are substantially more mysterious.

The idea behind IVLS is to multiply both sides of (4) by  $Z'$ , getting

$$(5) \quad Z'Y = Z'X\beta + Z'\delta.$$

This is a least squares problem. The response variable is  $Z'Y$ . The design matrix is  $Z'X$  and the error term is  $Z'\delta$ . The parameter vector is still  $\beta$ .

Econometricians use GLS (example 5.1, p. 65) to estimate (5), rather than OLS. This is because  $\text{cov}(Z'\delta|Z) = \sigma^2 Z'Z \neq \sigma^2 I_{q \times q}$  (exercise 3C4). Assumptions (i)-(ii) show that  $Z'Z$  has an inverse; and the inverse has a square root (exercise B1 below). We multiply both sides of (5) by  $(Z'Z)^{-1/2}$  to get

$$(6) \quad [(Z'Z)^{-1/2} Z'Y] = [(Z'Z)^{-1/2} Z'X]\beta + \eta, \text{ where } \eta = (Z'Z)^{-1/2} Z'\delta.$$

Apart from a little wrinkle to be discussed below, equation (6) is the usual regression model. As far as the errors are concerned,

$$(7) \quad E(\eta|Z) = 0$$

because  $Z$  was assumed exogenous: see (iv)-(v). (You want to condition on  $Z$  not  $X$ , because the latter is endogeneous.) Moreover,

$$\begin{aligned} (8) \quad \text{cov}(\eta|Z) &= E[(Z'Z)^{-1/2} Z'\delta \delta' Z'(Z'Z)^{-1/2} | Z] \\ &= (Z'Z)^{-1/2} Z' E[\delta \delta' | Z] Z (Z'Z)^{-1/2} \\ &= (Z'Z)^{-1/2} Z' \sigma^2 I_{n \times n} Z (Z'Z)^{-1/2} \\ &= \sigma^2 (Z'Z)^{-1/2} (Z'Z) (Z'Z)^{-1/2} \\ &= \sigma^2 I_{q \times q}. \end{aligned}$$

The big move is in the third line:  $E[\delta\delta'|Z] = \sigma^2 I_{n \times n}$ , because  $Z$  was assumed to be exogenous, and the  $\delta_i$  were assumed to be IID with mean 0 and variance  $\sigma^2$ : see (iv)-(v). Otherwise, we're just factoring constants out of the expectation and juggling matrices.

The OLS estimate for  $\beta$  in (6) is

$$(9) \quad \tilde{\beta} = (M'M)^{-1}M'L,$$

where  $M = (Z'Z)^{-1/2}Z'X$  is the design matrix and  $L = (Z'Z)^{-1/2}Z'Y$  is the response variable. (Exercise B1 shows that all the inverses exist.)

The IVLS estimator in the original system (4) is usually given as

$$(10) \quad \hat{\beta}_{IVLS} = [X'Z(Z'Z)^{-1}Z'X]^{-1} X'Z(Z'Z)^{-1}Z'Y.$$

We will show that  $\hat{\beta}_{IVLS} = \tilde{\beta}$ , completing the derivation of the IVLS estimator. This takes a bit of algebra. For starters, because  $Z'Z$  is symmetric,

$$(11) \quad M'M = X'Z(Z'Z)^{-1/2}(Z'Z)^{-1/2}Z'X = X'Z(Z'Z)^{-1}Z'X,$$

and

$$(12) \quad M'L = X'Z(Z'Z)^{-1/2}(Z'Z)^{-1/2}Z'Y = X'Z(Z'Z)^{-1}Z'Y.$$

Substituting (11) and (12) into (9) proves that  $\hat{\beta}_{IVLS} = \tilde{\beta}$ .

Standard errors are estimated using (13–14):

$$(13) \quad \widehat{\text{cov}}(\hat{\beta}_{IVLS}|Z) = \hat{\sigma}^2 [X'Z(Z'Z)^{-1}Z'X]^{-1},$$

where

$$(14) \quad \hat{\sigma}^2 = \|Y - X\hat{\beta}_{IVLS}\|^2/(n - p).$$

Exercise C6 below provides an informal justification for definitions (13)–(14), and theorem 1 in section 8 has some rigor. It is conventional to divide by  $n - p$  in (14), but theorem 4.4 does not apply because we're not in the OLS model: see the discussion of “the little wrinkle,” below.

Equation (10) is pretty dense. For some people, it helps to check that all the multiplications make sense. For instance,  $Z$  is  $n \times q$ , so  $Z'$  is  $q \times n$ . Then  $Z'Z$  and  $(Z'Z)^{-1}$  are  $q \times q$ . Next,  $X$  is  $n \times p$ , so  $X'$  is  $p \times n$ . Thus,  $X'Z$  is  $p \times q$  and  $Z'X$  is  $q \times p$ , which makes  $X'Z(Z'Z)^{-1}Z'X$  a  $p \times p$  matrix. What about  $X'Z(Z'Z)^{-1}Z'Y$ ? Well,  $X'Z$  is  $p \times q$ ,  $(Z'Z)^{-1}$  is  $q \times q$ , and  $Z'Y$  is  $q \times 1$ . So  $X'Z(Z'Z)^{-1}Z'Y$  is  $p \times 1$ . This is pretty dense too, but there is a simple bottom line:  $\hat{\beta}_{IVLS}$  is  $p \times 1$ , like it should be.

*Identification.* The matrix equation (5) unpacks to  $q$  ordinary equations in  $p$  unknowns—the components of  $\beta$ . (i) If  $q > p$ , there usually won't be any vector  $\beta$  that satisfies (5) exactly. GLS gives a compromise solution  $\hat{\beta}_{IVLS}$ . (ii) If  $q = p$ , there is a unique solution, which is  $\hat{\beta}_{IVLS}$ : see exercise C5 below. (iii) If  $q < p$ , we don't have enough equations relative to the number of parameters that we are estimating. There will be many  $\beta$ 's satisfying (5). That is the tipoff to under-identification.

*The little wrinkle in (6).* Given  $Z$ , the design matrix  $M = (Z'Z)^{-1/2}Z'X$  is still related to the errors  $\eta = (Z'Z)^{-1/2}Z'\delta$ , because of the endogeneity of  $X$ . This leads to *small-sample bias*. However, with luck,  $M$  will be practically constant, and a little bit of correlated randomness shouldn't matter. Theorem 1 in section 8 will make these ideas more precise.

### Exercise set B

- By assumptions (i)-(ii),  $Z'X$  is  $q \times p$  of rank  $p$ , and  $Z'Z$  is  $q \times q$  of rank  $q$ . Show that:
  - $Z'Z$  is positive definite and invertible; the inverse has a square root.
  - $X'Z(Z'Z)^{-1}Z'X$  is positive definite, hence invertible. Hint. Suppose  $c$  is  $p \times 1$ . Can  $c'X'Z(Z'Z)^{-1}Z'Xc \leq 0$ ?

Note. Without assumptions (i)-(ii), equations (10) and (13) wouldn't make sense.

- Let  $U_i$  be IID random variables. Let  $\bar{U} = \frac{1}{n} \sum_{i=1}^n U_i$ . True or false, and explain:
  - $E(U_i)$  is the same for all  $i$ .
  - $\text{var}(U_i)$  is the same for all  $i$ .
  - $E(U_i) = \bar{U}$ .
  - $\text{var}(U_i) = \frac{1}{n} \sum_{i=1}^n (U_i - \bar{U})^2$ .
  - $\text{var}(U_i) = \frac{1}{n-1} \sum_{i=1}^n (U_i - \bar{U})^2$ .

### 9.3 Estimating the butter model

Our next project is to estimate the butter model using IVLS. We'll start with the supply equation (2a). The equation is often written this way:

$$(15) \quad Q_t = a_0 + a_1 P_t + a_2 W_t + a_3 H_t + \delta_t \quad \text{for } t = 1, \dots, 20.$$

The actual price and quantity in year  $t$  are substituted for the free variables  $Q$  and  $P$  that define the supply schedule. Reminder: according to the law