

Endogeneity

Zhentaο Shi

November 10, 2019

1 Introduction

In microeconomic analysis, exogenous variables are the factors determined outside of the economic system under consideration, and endogenous variables are those decided within the economic system. The terms “endogenous” and “exogenous” in microeconomics will be carried over into multiple-equation econometric models. While in a single-equation regression model

$$y_i = x_i' \beta + e_i \quad (1)$$

is only part of the equation system. To make it simple, in the single-equation model we say an x_{ik} is *endogenous*, or is an *endogenous variable*, if $\text{cov}(x_{ik}, e_i) \neq 0$; otherwise x_{ik} is an *exogenous variable*.

Empirical works using linear regressions are routinely challenged by questions about endogeneity. Such questions plague economic seminars and referee reports. To defend empirical strategies in quantitative economic studies, it is important to understand the source of potential endogeneity and thoroughly discuss attempts for resolving endogeneity.

Endogeneity usually implies difficulty in identifying the parameter of interest with only (y_i, x_i) . Identification is critical for the interpretation of empirical economic research. We say a parameter is *identified* if the mapping between the parameter in the model and the distribution of the observed variable is one-to-one; otherwise we say the parameter is *under-identified*, or the parameter is failed to be identified. This is an abstract definition, and let us discuss it in more family linear regression context.

Example 1 (Identification failure due to collinearity). The linear projection model implies the moment equation

$$\mathbb{E}[x_i x_i'] \beta = \mathbb{E}[x_i y_i]. \quad (2)$$

If $\mathbb{E}[x_i x_i']$ is of full rank, then $\beta = (\mathbb{E}[x_i x_i'])^{-1} \mathbb{E}[x_i y_i]$ is a function of the quantities of the population moment and it is identified. On the contrary, if some x_k 's are perfect collinear so that $\mathbb{E}[x_i x_i']$ is rank deficient, there are multiple β that satisfies the k -equation system (2). Identification fails. \square

Example 2 (Identification failure due to endogeneity). Suppose x_i is a scalar random variable,

$$\begin{pmatrix} x_i \\ e_i \end{pmatrix} \sim N \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \sigma_{xe} \\ \sigma_{xe} & 1 \end{pmatrix} \right)$$

follows a joint normal distribution, and the dependent variable y_i is generated from (1). The joint normal assumption implies that the conditional mean

$$\mathbb{E}[y_i | x_i] = \beta x_i + \mathbb{E}[e_i | x_i] = (\beta + \sigma_{xe}) x_i$$

coincides with the linear projection model, and $\beta + \sigma_{xe}$ is the linear projection coefficient. From the observable random variable (y_i, x_i) , we can only learn $\beta + \sigma_{xe}$. As we cannot learn σ_{xe} from the data due to the unobservable e_i , there is no way to recover β . This is exactly the *omitted variable bias* that we have discussed earlier in this course. The gap lies between the available data (y_i, x_i) and the identification of the model. In the special case that we assume $\sigma_{xe} = 0$, the endogeneity vanishes and β is identified.

The linear projection model is so far the most general model in this course that justifies OLS. OLS is consistent for the linear projection coefficient. By the definition of the linear projection model, $\mathbb{E}[x_i e_i] = 0$ so there is no room for endogeneity in the linear projection model. In other words, if we talk about endogeneity, we must not be working with the linear projection model, and the coefficients we pursue are not the linear projection coefficients. \square

In econometrics we are often interested in a model with economic interpretation. The common practice in empirical research assumes that the observed data are generated from a parsimonious model, and the next step is to estimate the unknown parameters in the model. Since it is often possible to name some factors not included in the regressors but they are correlated with the included regressors and in the mean time also affects y_i , endogeneity becomes a fundamental problem.

To resolve endogeneity, we seek extra variables or data structure that may guarantee the identification of the model. The most often used methods are (i) fixed effect model (ii) instrumental variables. The fixed effect model requires that multiple observations, often across time, are collected for each individual i . Moreover, the source of endogeneity is time invariant and enters the model additively in the form

$$y_{it} = x'_{it}\beta + u_{it},$$

where $u_{it} = \alpha_i + \epsilon_{it}$ is the composite error. The panel data approach extends (y_i, x_i) to $(y_{it}, x_{it})_{i=1}^T$ if data are available along the time dimension.

The instrumental variable approach extends (y_i, x_i) to (y_i, x_i, z_i) , where the extra random variable z_i is called the *instrument variable*. It is assumed that z_i is orthogonal to the error e_i . Therefore, along with the model it adds an extra variable z_i .

Before closing this section, we stress that either the panel data approach or the instrumental variable approach entails extra information beyond (y_i, x_i) . Without such extra data, there is no way to resolve the identification failure. Just as the linear project model is available for any joint distribution of (y_i, x_i) with existence of suitable moments, from a pure statistical point of view a linear IV model is an artifact depends only on the choice of (y_i, x_i, z_i) without referencing to any economics. In essence, the linear IV model seeks a linear combination $y_i - \beta x_i$ that is orthogonal to the linear space spanned by z_i .

2 Examples

As econometricians mostly work with non-experimental data, we cannot overstate the importance of the endogeneity problem. We go over a few examples.

Example 3 (Dynamic Panel Model). We know that the first-difference (FD) estimator is consistent for (static) panel data model. Nevertheless, the FD estimator encounters difficulty in a dynamic panel model

$$y_{it} = \beta_1 + \beta_2 y_{it-1} + \beta_3 x_{it} + \alpha_i + \epsilon_{it},$$

even if we assume

$$\mathbb{E} [\epsilon_{it} | \alpha_i, x_{i1}, \dots, x_{iT}, y_{it-1}, y_{it-2}, \dots, y_{i0}] = 0. \quad (3)$$

When taking difference of the above equation for periods t and $t - 1$, we have

$$(y_{it} - y_{it-1}) = \beta_2 (y_{it-1} - y_{it-2}) + \beta_3 (x_{it} - x_{it-1}) + (\epsilon_{it} - \epsilon_{it-1}).$$

Under (3), $\mathbb{E} [(x_{it} - x_{it-1}) (\epsilon_{it} - \epsilon_{it-1})] = 0$, but

$$\mathbb{E} [(y_{it-1} - y_{it-2}) (\epsilon_{it} - \epsilon_{it-1})] = -\mathbb{E} [y_{it-1} \epsilon_{it-1}] = -\mathbb{E} [\epsilon_{it-1}^2] \neq 0. \quad \square$$

Example 4 (Classical Measurement Error). Endogeneity also emerges when an explanatory variable is not directly observable but is replaced by a measurement with error. Suppose the true linear model is

$$y_i = \beta_1 + \beta_2 x_i^* + u_i, \quad (4)$$

with $\mathbb{E} [u_i | x_i^*] = 0$. We cannot observe x_i^* but we observe x_i , a measurement of x_i^* , and they are linked by

$$x_i = x_i^* + v_i$$

with $\mathbb{E} [v_i | x_i^*, u_i] = 0$. Such a formulation of the measurement error is called the *classical measurement error*. Substitute out the unobservable x_i^* in (4),

$$y_i = \beta_1 + \beta_2 (x_i - v_i) + u_i = \beta_1 + \beta_2 x_i + e_i \quad (5)$$

where $e_i = u_i - \beta_2 v_i$. The correlation

$$\mathbb{E} [x_i e_i] = \mathbb{E} [(x_i^* + v_i) (u_i - \beta_2 v_i)] = -\beta_2 \mathbb{E} [v_i^2] \neq 0.$$

OLS (5) would not deliver a consistent estimator. \square

Next, we give two examples of equation systems, one from microeconomics and the other from macroeconomics.

Example 5 (Demand-Supply System). Let p_i and q_i be a good's log-price and log-quantity on the i -th market, and they are iid across markets. We are interested in the demand curve

$$p_i = \alpha_d - \beta_d q_i + e_{di} \quad (6)$$

for some $\beta_d \geq 0$ and the supply curve

$$p_i = \alpha_s + \beta_s q_i + e_{si} \quad (7)$$

for some $\beta_s \geq 0$. We use a simple linear specification so that the coefficient β_d can be interpreted as demand elasticity and β_s as supply elasticity. Undergraduate microeconomics teaches the deterministic form but we add an error term to cope with the data. Can we learn the elasticities by regression p_i on q_i ?

The two equations can be written in a matrix form

$$\begin{pmatrix} 1 & \beta_d \\ 1 & -\beta_s \end{pmatrix} \begin{pmatrix} p_i \\ q_i \end{pmatrix} = \begin{pmatrix} \alpha_d \\ \alpha_s \end{pmatrix} + \begin{pmatrix} e_{di} \\ e_{si} \end{pmatrix}. \quad (8)$$

Microeconomic terminology calls (p_i, q_i) endogenous variables and (e_{di}, e_{si}) exogenous variables. (8) is a *structural equation* because it is motivated from economic theory so that the coefficients bear economic meaning. If we rule out the trivial case $\beta_d = \beta_s = 0$, we can solve

$$\begin{pmatrix} p_i \\ q_i \end{pmatrix} = \frac{1}{\beta_s + \beta_d} \begin{pmatrix} \beta_s & \beta_d \\ 1 & -1 \end{pmatrix} \left[\begin{pmatrix} \alpha_d \\ \alpha_s \end{pmatrix} + \begin{pmatrix} e_{di} \\ e_{si} \end{pmatrix} \right]. \quad (9)$$

This equation (9) is called the *reduced form*—the endogenous variables are expressed as explicit functions of the parameters and the exogenous variables. In particular,

$$q_i = (\alpha_d + e_{di} - \alpha_s - e_{si}) / (\beta_s + \beta_d)$$

so that the log-price is correlated with both e_{si} and e_{di} . As q_i is endogenous (in the econometric sense) in either (6) or (7), neither the demand elasticity nor the supply elasticity is identified with (p_i, q_i) . Indeed, as

$$q_i = (\beta_s \alpha_d + \beta_d \alpha_s + \beta_s e_{di} + \beta_d e_{si}) / (\beta_s + \beta_d)$$

from (9), the linear projection coefficient of p_i on q_i is

$$\frac{\text{cov}(p_i, q_i)}{\text{var}(q_i)} = \frac{\beta_s \sigma_d^2 - \beta_d \sigma_s^2 + (\beta_d - \beta_s) \sigma_{sd}}{\beta_d^2 \sigma_d^2 + \beta_s^2 \sigma_s^2 + 2\beta_d \beta_s \sigma_{sd}},$$

where $\sigma_d^2 = \text{var}(e_{di})$, $\sigma_s^2 = \text{var}(e_{si})$ and $\sigma_{sd} = \text{cov}(e_{di}, e_{si})$.

This is a classical example of the demand-supply system. The structural parameter cannot be directly identified because the observed (p_i, q_i) is the outcome of an equilibrium—the crossing of the demand curve and the supply curve. To identify the demand curve, we will need an instrument that shifts the supply curve only; and vice versa. \square

Example 6 (Keynesian-Type Macro Equations). This is a model borrowed from Hayashi (2000, p.193) but originated from Haavelmo (1943). An econometrician is interested in learning β_2 , the marginal propensity of consumption, in the Keynesian-type equation

$$C_i = \beta_1 + \beta_2 Y_i + u_i \quad (10)$$

where C_i is household consumption, Y_i is the GNP, and u_i is the unobservable error. However, Y_i and C_i are connected by an accounting equality (with no error)

$$Y_i = C_i + I_i,$$

where I_i is investment. We assume $\mathbb{E}[u_i | I_i] = 0$ as investment is determined in advance. OLS (10) will be inconsistent because in the reduced-form $Y_i = \frac{1}{1-\beta_2} (\beta_1 + u_i + I_i)$ implies $\mathbb{E}[Y_i u_i] = \mathbb{E}[u_i^2] / (1 - \beta_2) \neq 0$. \square