

# Instrumental Variables Estimation and Two Stage Least Squares

In this chapter, we further study the problem of **endogenous explanatory variables** in multiple regression models. In Chapter 3, we derived the bias in the OLS estimators when an important variable is omitted; in Chapter 5, we showed that OLS is generally inconsistent under **omitted variables**. Chapter 9 demonstrated that omitted variables bias can be eliminated (or at least mitigated) when a suitable proxy variable is given for an unobserved explanatory variable. Unfortunately, suitable proxy variables are not always available.

In the previous two chapters, we explained how fixed effects estimation or first differencing can be used with panel data to estimate the effects of time-varying independent variables in the presence of *time-constant* omitted variables. Although such methods are very useful, we do not always have access to panel data. Even if we can obtain panel data, it does us little good if we are interested in the effect of a variable that does not change over time: first differencing or fixed effects estimation eliminates time-constant explanatory variables. In addition, the panel data methods that we have studied so far do not solve the problem of time-varying omitted variables that are correlated with the explanatory variables.

In this chapter, we take a different approach to the endogeneity problem. You will see how the method of instrumental variables (IV) can be used to solve the problem of endogeneity of one or more explanatory variables. The method of two stage least squares (2SLS or TSLS) is second in popularity only to ordinary least squares for estimating linear equations in applied econometrics.

We begin by showing how IV methods can be used to obtain consistent estimators in the presence of omitted variables. IV can also be used to solve the **errors-in-variables** problem, at least under

certain assumptions. Chapter 16 will demonstrate how to estimate simultaneous equations models using IV methods.

Our treatment of instrumental variables estimation closely follows our development of ordinary least squares in Part 1, where we assumed that we had a random sample from an underlying population. This is a desirable starting point because, in addition to simplifying the notation, it emphasizes that the important assumptions for IV estimation are stated in terms of the underlying population (just as with OLS). As we showed in Part 2, OLS can be applied to time series data, and the same is true of instrumental variables methods. Section 15-7 discusses some special issues that arise when IV methods are applied to time series data. In Section 15-8, we cover applications to pooled cross sections and panel data.

## 15-1 Motivation: Omitted Variables in a Simple Regression Model

When faced with the prospect of omitted variables bias (or unobserved heterogeneity), we have so far discussed three options: (1) we can ignore the problem and suffer the consequences of biased and inconsistent estimators; (2) we can try to find and use a suitable proxy variable for the unobserved variable; or (3) we can assume that the omitted variable does not change over time and use the fixed effects or first-differencing methods from Chapters 13 and 14. The first response can be satisfactory if the estimates are coupled with the direction of the biases for the key parameters. For example, if we can say that the estimator of a positive parameter, say, the effect of job training on subsequent wages, is biased toward zero and we have found a statistically significant positive estimate, we have still learned something: job training has a positive effect on wages, and it is likely that we have underestimated the effect. Unfortunately, the opposite case, where our estimates may be too large in magnitude, often occurs, which makes it very difficult for us to draw any useful conclusions.

The proxy variable solution discussed in Section 9-2 can also produce satisfying results, but it is not always possible to find a good proxy. This approach attempts to solve the omitted variable problem by replacing the unobservable with one or more proxy variables.

Another approach leaves the unobserved variable in the error term, but rather than estimating the model by OLS, it uses an estimation method that recognizes the presence of the omitted variable. This is what the method of instrumental variables does.

For illustration, consider the problem of unobserved ability in a wage equation for working adults. A simple model is

$$\log(\text{wage}) = \beta_0 + \beta_1 \text{educ} + \beta_2 \text{abil} + e,$$

where  $e$  is the error term. In Chapter 9, we showed how, under certain assumptions, a proxy variable such as  $IQ$  can be substituted for ability, and then a consistent estimator of  $\beta_1$  is available from the regression of

$$\log(\text{wage}) \text{ on } \text{educ}, IQ.$$

Suppose, however, that a proxy variable is not available (or does not have the properties needed to produce a consistent estimator of  $\beta_1$ ). Then, we put  $\text{abil}$  into the error term, and we are left with the simple regression model

$$\log(\text{wage}) = \beta_0 + \beta_1 \text{educ} + u, \quad [15.1]$$

where  $u$  contains  $\text{abil}$ . Of course, if equation (15.1) is estimated by OLS, a biased and inconsistent estimator of  $\beta_1$  results if  $\text{educ}$  and  $\text{abil}$  are correlated.

It turns out that we can still use equation (15.1) as the basis for estimation, provided we can find an instrumental variable for *educ*. To describe this approach, the simple regression model is written as

$$y = \beta_0 + \beta_1 x + u, \quad [15.2]$$

where we think that  $x$  and  $u$  are correlated (have nonzero covariance):

$$\text{Cov}(x, u) \neq 0. \quad [15.3]$$

The method of instrumental variables works whether or not  $x$  and  $u$  are correlated, but, for reasons we will see later, OLS should be used if  $x$  is uncorrelated with  $u$ .

In order to obtain consistent estimators of  $\beta_0$  and  $\beta_1$  when  $x$  and  $u$  are correlated, we need some additional information. The information comes by way of a new variable that satisfies certain properties. Suppose that we have an observable variable  $z$  that satisfies these two assumptions: (1)  $z$  is uncorrelated with  $u$ , that is,

$$\text{Cov}(z, u) = 0; \quad [15.4]$$

(2)  $z$  is correlated with  $x$ , that is,

$$\text{Cov}(z, x) \neq 0. \quad [15.5]$$

Then, we call  $z$  an **instrumental variable** for  $x$ , or sometimes simply an **instrument** for  $x$ .

The requirement that the instrument  $z$  satisfies (15.4) is summarized by saying “ $z$  is exogenous in equation (15.2),” and so we often refer to (15.4) as **instrument exogeneity**. In the context of omitted variables, instrument exogeneity means that  $z$  should have no partial effect on  $y$  (after  $x$  and omitted variables have been controlled for), and  $z$  should be uncorrelated with the omitted variables. Equation (15.5) means that  $z$  must be related, either positively or negatively, to the endogenous explanatory variable  $x$ . This condition is sometimes referred to as **instrument relevance** (as in “ $z$  is relevant for explaining variation in  $x$ ”).

There is a very important difference between the two requirements for an instrumental variable. Because (15.4) involves the covariance between  $z$  and the unobserved error  $u$ , we cannot generally hope to test this assumption: in most cases, we must maintain  $\text{Cov}(z, u) = 0$  by appealing to economic behavior or introspection. Sometimes, we might have an observable proxy variable for some factor contained in  $u$ , in which case we can check to see if  $z$  and the proxy variable are roughly uncorrelated. Of course, if we have a good proxy for an important element of  $u$ , we might just add the proxy as an explanatory variable and estimate the expanded equation by ordinary least squares. See Section 9-2.

Some readers may be wondering why we do not attempt to check (15.4) by using the following procedure. Given a sample of size  $n$ , obtain the OLS residuals,  $\hat{u}_i$ , from the regression  $y_i$  on  $x_i$ . Then, devise a test based on the sample correlation between  $z_i$  and  $\hat{u}_i$  as a check on whether  $z_i$  and the unobserved errors  $u_i$  are correlated. A moment’s thought reveals the logical problem with this procedure. The entire reason for moving beyond OLS is that we think the OLS estimators of  $\beta_0$  and  $\beta_1$  are inconsistent due to correlation between  $x$  and  $u$ . Therefore, in computing the OLS residuals  $\hat{u}_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i$ , we are not getting useful estimates of the  $u_i$ . Therefore, we can learn nothing by studying the correlation between  $z_i$  and  $\hat{u}_i$ . A related suggestion is to use the OLS regression  $y_i$  on  $x_i$ ,  $z_i$  and to conclude  $z_i$  satisfies the exogeneity requirement if its coefficient is statistically insignificant. Again, this procedure does not work, regardless of the outcome of the test, because  $x$  is allowed to be endogenous. The bottom line is that, in the current setting, we have no way of testing (15.4) unless we use external information.

By contrast, the condition that  $z$  is correlated with  $x$  (in the population) can be tested, given a random sample from the population. The easiest way to do this is to estimate a simple regression between  $x$  and  $z$ . In the population, we have

$$x = \pi_0 + \pi_1 z + v. \quad [15.6]$$

Then, because  $\pi_1 = \text{Cov}(z, x)/\text{Var}(z)$ , assumption (15.5) holds if, and only if,  $\pi_1 \neq 0$ . Thus, we should be able to *reject* the null hypothesis

$$H_0: \pi_1 = 0 \quad [15.7]$$

against the two-sided alternative  $H_0: \pi_1 \neq 0$ , at a sufficiently small significance level. If this is the case, then we can be fairly confident that (15.5) holds.

For the  $\log(\text{wage})$  equation in (15.1), an instrumental variable  $z$  for  $\text{educ}$  must be (1) uncorrelated with ability (and any other unobserved factors affecting wage) and (2) correlated with education. Something such as the last digit of an individual's Social Security Number almost certainly satisfies the first requirement: it is uncorrelated with ability because it is determined randomly. However, it is precisely because of the randomness of the last digit of the SSN that it is not correlated with education, either; therefore it makes a poor instrumental variable for  $\text{educ}$  because it violates the instrument relevance requirement in equation (15.5).

What we have called a *proxy variable* for the omitted variable makes a poor IV for the opposite reason. For example, in the  $\log(\text{wage})$  example with omitted ability, a proxy variable for  $\text{abil}$  should be as highly correlated as possible with  $\text{abil}$ . An instrumental variable must be *uncorrelated* with  $\text{abil}$ . Therefore, while  $\text{IQ}$  is a good candidate as a proxy variable for  $\text{abil}$ , it is not a good instrumental variable for  $\text{educ}$  because it violates the instrument exogeneity requirement in equation (15.4).

Whether other possible instrumental variable candidates satisfy the exogeneity requirement in (15.4) is less clear-cut. In wage equations, labor economists have used family background variables as IVs for education. For example, mother's education ( $\text{motheduc}$ ) is positively correlated with child's education, as can be seen by collecting a sample of data on working people and running a simple regression of  $\text{educ}$  on  $\text{motheduc}$ . Therefore,  $\text{motheduc}$  satisfies equation (15.5). The problem is that mother's education might also be correlated with child's ability (through mother's ability and perhaps quality of nurturing at an early age), in which case (15.4) fails.

Another IV choice for  $\text{educ}$  in (15.1) is number of siblings while growing up ( $\text{sibs}$ ). Typically, having more siblings is associated with lower average levels of education. Thus, if number of siblings is uncorrelated with ability, it can act as an instrumental variable for  $\text{educ}$ .

As a second example, consider the problem of estimating the causal effect of skipping classes on final exam score. In a simple regression framework, we have

$$\text{score} = \beta_0 + \beta_1 \text{skipped} + u, \quad [15.8]$$

where  $\text{score}$  is the final exam score and  $\text{skipped}$  is the total number of lectures missed during the semester. We certainly might be worried that  $\text{skipped}$  is correlated with other factors in  $u$ : more able, highly motivated students might miss fewer classes. Thus, a simple regression of  $\text{score}$  on  $\text{skipped}$  may not give us a good estimate of the causal effect of missing classes.

What might be a good IV for  $\text{skipped}$ ? We need something that has no direct effect on  $\text{score}$  and is not correlated with student ability and motivation. At the same time, the IV must be correlated with  $\text{skipped}$ . One option is to use distance between living quarters and classrooms. Especially at large universities, some living quarters will be further from a student's classrooms, and this may essentially be a random occurrence. Some students live off campus while others commute long distances. Living further away from classrooms may increase the likelihood of missing lectures due to bad weather, oversleeping, and so on. Thus,  $\text{skipped}$  may be positively correlated with  $\text{distance}$ ; this can be checked by regressing  $\text{skipped}$  on  $\text{distance}$  and doing a  $t$  test, as described earlier.

Is  $\text{distance}$  uncorrelated with  $u$ ? In the simple regression model (15.8), some factors in  $u$  may be correlated with  $\text{distance}$ . For example, students from low-income families may live off campus; if income affects student performance, this could cause  $\text{distance}$  to be correlated with  $u$ . Section 15-2 shows how to use IV in the context of multiple regression, so that other factors affecting  $\text{score}$  can be included directly in the model. Then,  $\text{distance}$  might be a good IV for  $\text{skipped}$ . An IV approach may not be necessary at all if a good proxy exists for student ability, such as cumulative GPA prior to the semester.

There is a final point worth emphasizing before we turn to the mechanics of IV estimation: namely, in using the simple regression in equation (15.6) to test (15.7), it is important to take note of the sign (and even magnitude) of  $\hat{\pi}_1$  and not just its statistical significance. Arguments for why a variable  $z$  makes a good IV candidate for an endogenous explanatory variable  $x$  should include a discussion about the nature of the relationship between  $x$  and  $z$ . For example, due to genetics and background influences

it makes sense that child's education ( $x$ ) and mother's education ( $z$ ) are positively correlated. If in your sample of data you find that they are actually negatively correlated—that is,  $\hat{\pi}_1 < 0$ —then your use of mother's education as an IV for child's education is likely to be unconvincing. [And this has nothing to do with whether condition (15.4) is likely to hold.] In the example of measuring whether skipping classes has an effect on test performance, one should find a positive, statistically significant relationship between *skipped* and *distance* in order to justify using *distance* as an IV for *skipped*: a negative relationship would be difficult to justify [and would suggest that there are important omitted variables driving a negative correlation—variables that might themselves have to be included in the model (15.8)].

We now demonstrate that the availability of an instrumental variable can be used to estimate consistently the parameters in equation (15.2). In particular, we show that assumptions (15.4) and (15.5) serve to *identify* the parameter  $\beta_1$ . **Identification** of a parameter in this context means that we can write  $\beta_1$  in terms of population moments that can be estimated using a sample of data. To write  $\beta_1$  in terms of population covariances, we use equation (15.2): the covariance between  $z$  and  $y$  is

$$\text{Cov}(z, y) = \beta_1 \text{Cov}(z, x) + \text{Cov}(z, u).$$

Now, under assumption (15.4),  $\text{Cov}(z, u) = 0$ , and under assumption (15.5),  $\text{Cov}(z, x) \neq 0$ . Thus, we can solve for  $\beta_1$  as

$$\beta_1 = \frac{\text{Cov}(z, y)}{\text{Cov}(z, x)}. \quad [15.9]$$

[Notice how this simple algebra fails if  $z$  and  $x$  are uncorrelated, that is, if  $\text{Cov}(z, x) = 0$ .] Equation (15.9) shows that  $\beta_1$  is the population covariance between  $z$  and  $y$  divided by the population covariance between  $z$  and  $x$ , which shows that  $\beta_1$  is identified. Given a random sample, we estimate the population quantities by the sample analogs. After canceling the sample sizes in the numerator and denominator, we get the **instrumental variables (IV) estimator** of  $\beta_1$ :

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (z_i - \bar{z})(y_i - \bar{y})}{\sum_{i=1}^n (z_i - \bar{z})(x_i - \bar{x})}. \quad [15.10]$$

Given a sample of data on  $x$ ,  $y$ , and  $z$ , it is simple to obtain the IV estimator in (15.10). The IV estimator of  $\beta_0$  is simply  $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$ , which looks just like the OLS intercept estimator except that the slope estimator,  $\hat{\beta}_1$ , is now the IV estimator.

It is no accident that when  $z = x$  we obtain the OLS estimator of  $\beta_1$ . In other words, when  $x$  is exogenous, it can be used as its own IV, and the IV estimator is then identical to the OLS estimator.

A simple application of the law of large numbers shows that the IV estimator is consistent for  $\beta_1$ :  $\text{plim}(\hat{\beta}_1) = \beta_1$ , provided assumptions (15.4) and (15.5) are satisfied. If either assumption fails, the IV estimators are not consistent (more on this later). One feature of the IV estimator is that, when  $x$  and  $u$  are in fact correlated—so that instrumental variables estimation is actually needed—it is essentially never unbiased. This means that, in small samples, the IV estimator can have a substantial bias, which is one reason why large samples are preferred.

When discussing the application of instrumental variables it is important to be careful with language. Like OLS, IV is an *estimation* method. It makes little sense to refer to “an instrumental variables model”—just as the phrase “OLS model” makes little sense. As we know, a model is an equation such as (15.8), which is a special case of the generic model in equation (15.2). When we have a model such as (15.2), we can choose to estimate the parameters of that model in many different ways. Prior to this chapter we focused primarily on OLS, but, for example, we also know from Chapter 8 that one can use weighted least squares as an alternative estimation method (and there are



unlimited possibilities for the weights). If we have an instrumental variable candidate  $z$  for  $x$ , then we can instead apply instrumental variables estimation. It is certainly true that the estimation method we apply is motivated by the model and assumptions we make about that model. But the estimators are well defined and exist apart from any underlying model or assumptions: remember, an estimator is simply a rule for combining data. The bottom line is that while we probably know what a researcher means when using a phrase such as “I estimated an IV model,” such language betrays a lack of understanding about the difference between a model and an estimation method.

### 15-1a Statistical Inference with the IV Estimator

Given the similar structure of the IV and OLS estimators, it is not surprising that the IV estimator has an approximate normal distribution in large sample sizes. To perform inference on  $\beta_1$ , we need a standard error that can be used to compute  $t$  statistics and confidence intervals. The usual approach is to impose a homoskedasticity assumption, just as in the case of OLS. Now, the homoskedasticity assumption is stated conditional on the instrumental variable,  $z$ , not the endogenous explanatory variable,  $x$ . Along with the previous assumptions on  $u$ ,  $x$ , and  $z$ , we add

$$E(u^2|z) = \sigma^2 = \text{Var}(u). \quad [15.11]$$

It can be shown that, under (15.4), (15.5), and (15.11), the asymptotic variance of  $\hat{\beta}_1$  is

$$\frac{\sigma^2}{n\sigma_x^2\rho_{x,z}^2}, \quad [15.12]$$

where  $\sigma_x^2$  is the population variance of  $x$ ,  $\sigma^2$  is the population variance of  $u$ , and  $\rho_{x,z}^2$  is the square of the population correlation between  $x$  and  $z$ . This tells us how highly correlated  $x$  and  $z$  are in the population. As with the OLS estimator, the asymptotic variance of the IV estimator decreases to zero at the rate of  $1/n$ , where  $n$  is the sample size.

Equation (15.12) is interesting for two reasons. First, it provides a way to obtain a standard error for the IV estimator. All quantities in (15.12) can be consistently estimated given a random sample. To estimate  $\sigma_x^2$ , we simply compute the sample variance of  $x_i$ ; to estimate  $\rho_{x,z}^2$ , we can run the regression of  $x_i$  on  $z_i$  to obtain the  $R$ -squared, say,  $R_{x,z}^2$ . Finally, to estimate  $\sigma^2$ , we can use the IV residuals,

$$\hat{u}_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i, \quad i = 1, 2, \dots, n,$$

where  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are the IV estimates. A consistent estimator of  $\sigma^2$  looks just like the estimator of  $\sigma^2$  from a simple OLS regression:

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n \hat{u}_i^2,$$

where it is standard to use the degrees of freedom correction (even though this has little effect as the sample size grows).

The (asymptotic) standard error of  $\hat{\beta}_1$  is the square root of the estimated asymptotic variance, the latter of which is given by

$$\frac{\hat{\sigma}^2}{\text{SST}_x \cdot R_{x,z}^2}, \quad [15.13]$$

where  $\text{SST}_x$  is the total sum of squares of the  $x_i$ . [Recall that the sample variance of  $x_i$  is  $\text{SST}_x/n$ , and so the sample sizes cancel to give us (15.13).] The resulting standard error can be used to construct either  $t$  statistics for hypotheses involving  $\beta_1$  or confidence intervals for  $\beta_1$ .  $\hat{\beta}_0$  also has a standard error that we do not present here. Any modern econometrics package computes the standard error after any IV estimation; there is rarely any reason to perform the calculations by hand.

A second reason (15.12) is interesting is that it allows us to compare the asymptotic variances of the IV and the OLS estimators (when  $x$  and  $u$  are uncorrelated). Under the Gauss-Markov assumptions, the variance of the OLS estimator is  $\sigma^2/\text{SST}_x$ , while the comparable formula for the IV estimator is  $\sigma^2/(\text{SST}_x \cdot R_{x,z}^2)$ ; they differ only in that  $R_{x,z}^2$  appears in the denominator of the IV variance. Because an  $R$ -squared is always less than one, the IV variance is always larger than the OLS variance (when OLS is valid). If  $R_{x,z}^2$  is small, then the IV variance can be much larger than the OLS variance. Remember,  $R_{x,z}^2$  measures the strength of the linear relationship between  $x$  and  $z$  in the sample. If  $x$  and  $z$  are only slightly correlated,  $R_{x,z}^2$  can be small, and this can translate into a very large sampling variance for the IV estimator. The more highly correlated  $z$  is with  $x$ , the closer  $R_{x,z}^2$  is to one, and the smaller is the variance of the IV estimator. In the case that  $z = x$ ,  $R_{x,z}^2 = 1$ , and we get the OLS variance, as expected.

The previous discussion highlights an important cost of performing IV estimation when  $x$  and  $u$  are uncorrelated: the asymptotic variance of the IV estimator is always larger, and sometimes much larger, than the asymptotic variance of the OLS estimator.

### EXAMPLE 15.1 Estimating the Return to Education for Married Women

We use the data on married working women in MROZ to estimate the return to education in the simple regression model

$$\log(\text{wage}) = \beta_0 + \beta_1 \text{educ} + u. \quad [15.14]$$

For comparison, we first obtain the OLS estimates:

$$\begin{aligned} \widehat{\log(\text{wage})} &= -.185 + .109 \text{educ} \\ &(.185) \quad (.014) \\ n &= 428, R^2 = .118. \end{aligned} \quad [15.15]$$

The estimate for  $\beta_1$  implies an almost 11% return for another year of education.

Next, we use father's education (*fatheduc*) as an instrumental variable for *educ*. We have to maintain that *fatheduc* is uncorrelated with  $u$ . The second requirement is that *educ* and *fatheduc* are correlated. We can check this very easily using a simple regression of *educ* on *fatheduc* (using only the working women in the sample):

$$\begin{aligned} \widehat{\text{educ}} &= 10.24 + .269 \text{fatheduc} \\ &(.28) \quad (.029) \\ n &= 428, R^2 = .173. \end{aligned} \quad [15.16]$$

The  $t$  statistic on *fatheduc* is 9.28, which indicates that *educ* and *fatheduc* have a statistically significant positive correlation. (In fact, *fatheduc* explains about 17% of the variation in *educ* in the sample.) Using *fatheduc* as an IV for *educ* gives

$$\begin{aligned} \widehat{\log(\text{wage})} &= .441 + .059 \text{educ} \\ &(.446) \quad (.035) \\ n &= 428, R^2 = .093. \end{aligned} \quad [15.17]$$

The IV estimate of the return to education is 5.9%, which is barely more than one-half of the OLS estimate. This suggests that the OLS estimate is too high and is consistent with omitted ability bias. But we should remember that these are estimates from just one sample: we can never know whether .109 is above the true return to education, or whether .059 is closer to the true return to education. Further, the standard error of the IV estimate is two and one-half times as large as the OLS standard error (this is expected, for the reasons we gave earlier). The 95% confidence interval for  $\beta_1$  using OLS is much tighter than that using the IV; in fact, the IV confidence interval actually contains the OLS estimate. Therefore, although the differences between (15.15) and (15.17) are practically large, we cannot say whether the difference is statistically significant. We will show how to test this in Section 15-5.

In the previous example, the estimated return to education using IV was less than that using OLS, which corresponds to our expectations. But this need not have been the case, as the following example demonstrates.

### EXAMPLE 15.2 Estimating the Return to Education for Men

We now use WAGE2 to estimate the return to education for men. We use the variable *sibs* (number of siblings) as an instrument for *educ*. These are negatively correlated, as we can verify from a simple regression:

$$\begin{aligned}\widehat{educ} &= 14.14 - .228 \text{ sibs} \\ &(.11) \quad (.030) \\ n &= 935, R^2 = .057.\end{aligned}$$

This equation implies that every sibling is associated with, on average, about .23 less of a year of education. If we assume that *sibs* is uncorrelated with the error term in (15.14), then the IV estimator is consistent. Estimating equation (15.14) using *sibs* as an IV for *educ* gives

$$\begin{aligned}\widehat{\log(wage)} &= 5.13 + .122 \text{ educ} \\ &(.36) \quad (.026) \\ n &= 935.\end{aligned}$$

(The *R*-squared is computed to be negative, so we do not report it. A discussion of *R*-squared in the context of IV estimation follows.) For comparison, the OLS estimate of  $\beta_1$  is .059 with a standard error of .006. Unlike in the previous example, the IV estimate is now much higher than the OLS estimate. Although we do not know whether the difference is statistically significant, this does not mesh with the omitted ability bias from OLS. It could be that *sibs* is also correlated with ability: more siblings means, on average, less parental attention, which could result in lower ability. Another interpretation is that the OLS estimator is biased toward zero because of measurement error in *educ*. This is not entirely convincing because, as we discussed in Section 9-3, *educ* is unlikely to satisfy the classical errors-in-variables model.

In the previous examples, the endogenous explanatory variable (*educ*) and the instrumental variables (*fatheduc*, *sibs*) have quantitative meaning. But nothing prevents the explanatory variable or IV from being binary variables. Angrist and Krueger (1991), in their simplest analysis, came up with a clever binary instrumental variable for *educ*, using census data on men in the United States. Let *firstqtr* be equal to one if the man was born in the first quarter of the year, and zero otherwise. It seems that the error term in (15.14)—and, in particular, ability—should be unrelated to quarter of birth. But *firstqtr* also needs to be correlated with *educ*. It turns out that years of education *do* differ systematically in the population based on quarter of birth. Angrist and Krueger argued persuasively that this is due to compulsory school attendance laws in effect in all states. Briefly, students born early in the year typically begin school at an older age. Therefore, they reach the compulsory schooling age (16 in most states) with somewhat less education than students who begin school at a younger age. For students who finish high school, Angrist and Krueger verified that there is no relationship between years of education and quarter of birth.

Because years of education varies only slightly across quarter of birth—which means  $R^2_{x,z}$  in (15.13) is very small—Angrist and Krueger needed a very large sample size to get a reasonably precise IV estimate. Using 247,199 men born between 1920 and 1929, the OLS estimate of the return to education was .0801 (standard error .0004), and the IV estimate was .0715 (.0219); these are reported in Table III of Angrist and Krueger's paper. Note how large the *t* statistic is for the OLS estimate (about 200), whereas the *t* statistic for the IV estimate is only 3.26. Thus, the IV estimate is statistically different from zero, but its confidence interval is much wider than that based on the OLS estimate.

An interesting finding by Angrist and Krueger is that the IV estimate does not differ much from the OLS estimate. In fact, using men born in the next decade, the IV estimate is somewhat higher



than the OLS estimate. One could interpret this as showing that there is no omitted ability bias when wage equations are estimated by OLS. However, the Angrist and Krueger paper has been criticized on econometric grounds. As discussed by Bound, Jaeger, and Baker (1995), it is not obvious that season of birth is unrelated to unobserved factors that affect wage. As we will explain in the next subsection, even a small amount of correlation between  $z$  and  $u$  can cause serious problems for the IV estimator.

For policy analysis, the endogenous explanatory variable is often a binary variable. For example, Angrist (1990) studied the effect that being a veteran of the Vietnam War had on lifetime earnings. A simple model is

$$\log(\text{earnings}) = \beta_0 + \beta_1 \text{veteran} + u, \quad [15.18]$$

where *veteran* is a binary variable. The problem with estimating this equation by OLS is that there may be a *self-selection* problem, as we mentioned in Chapter 7: perhaps people who get the most out of the military choose to join, or the decision to join is correlated with other characteristics that affect earnings. These will cause *veteran* and  $u$  to be correlated.

### GOING FURTHER 15.1

If some men who were assigned low draft lottery numbers obtained additional schooling to reduce the probability of being drafted, is lottery number a good instrument for *veteran* in (15.18)?

Angrist pointed out that the Vietnam draft lottery provided a **natural experiment** (see also Chapter 13) that created an instrumental variable for *veteran*. Young men were given lottery numbers that determined whether they would be called to serve in Vietnam. Because the numbers given were (eventually) randomly assigned, it seems plausible that draft lottery number is uncorrelated with the error term  $u$ .

But those with a low enough number had to serve in Vietnam, so that the probability of being a veteran is correlated with lottery number. If both of these assertions are true, draft lottery number is a good IV candidate for *veteran*.

It is also possible to have a binary endogenous explanatory variable and a binary instrumental variable. See Problem 1 for an example.

## 15-1b Properties of IV with a Poor Instrumental Variable

We have already seen that, though IV is consistent when  $z$  and  $u$  are uncorrelated and  $z$  and  $x$  have any positive or negative correlation, IV estimates can have large standard errors, especially if  $z$  and  $x$  are only weakly correlated. Weak correlation between  $z$  and  $x$  can have even more serious consequences: the IV estimator can have a large asymptotic bias even if  $z$  and  $u$  are only moderately correlated.

We can see this by studying the probability limit of the IV estimator when  $z$  and  $u$  are possibly correlated. Letting  $\hat{\beta}_{1,IV}$  denote the IV estimator, we can write

$$\text{plim } \hat{\beta}_{1,IV} = \beta_1 + \frac{\text{Corr}(z,u)}{\text{Corr}(z,x)} \cdot \frac{\sigma_u}{\sigma_x}, \quad [15.19]$$

where  $\sigma_u$  and  $\sigma_x$  are the standard deviations of  $u$  and  $x$  in the population, respectively. The interesting part of this equation involves the correlation terms. It shows that, even if  $\text{Corr}(z,u)$  is small, the inconsistency in the IV estimator can be very large if  $\text{Corr}(z,x)$  is also small. Thus, even if we focus only on consistency, it is not necessarily better to use IV than OLS if the correlation between  $z$  and  $u$  is smaller than that between  $x$  and  $u$ . Using the fact that  $\text{Corr}(x,u) = \text{Cov}(x,u)/(\sigma_x \sigma_u)$  along with equation (5.3), we can write the plim of the OLS estimator—call it  $\hat{\beta}_{1,OLS}$ —as

$$\text{plim } \hat{\beta}_{1,OLS} = \beta_1 + \text{Corr}(x,u) \cdot \frac{\sigma_u}{\sigma_x}. \quad [15.20]$$

Comparing these formulas shows that it is possible for the directions of the asymptotic biases to be different for IV and OLS. For example, suppose  $\text{Corr}(x,u) > 0$ ,  $\text{Corr}(z,x) > 0$ , and  $\text{Corr}(z,u) < 0$ .

Then the IV estimator has a downward bias, whereas the OLS estimator has an upward bias (asymptotically). In practice, this situation is probably rare. More problematic is when the direction of the bias is the same and the correlation between  $z$  and  $x$  is small. For concreteness, suppose  $x$  and  $z$  are both positively correlated with  $u$  and  $\text{Corr}(z, x) > 0$ . Then the asymptotic bias in the IV estimator is less than that for OLS only if  $\text{Corr}(z, u)/\text{Corr}(z, x) < \text{Corr}(x, u)$ . If  $\text{Corr}(z, x)$  is small, then a seemingly small correlation between  $z$  and  $u$  can be magnified and make IV worse than OLS, even if we restrict attention to bias. For example, if  $\text{Corr}(z, x) = .2$ ,  $\text{Corr}(z, u)$  must be less than one-fifth of  $\text{Corr}(x, u)$  before IV has less asymptotic bias than OLS. In many applications, the correlation between the instrument and  $x$  is less than .2. Unfortunately, because we rarely have an idea about the relative magnitudes of  $\text{Corr}(z, u)$  and  $\text{Corr}(x, u)$ , we can never know for sure which estimator has the largest asymptotic bias [unless, of course, we assume  $\text{Corr}(z, u) = 0$ ].

In the Angrist and Krueger (1991) example mentioned earlier, where  $x$  is years of schooling and  $z$  is a binary variable indicating quarter of birth, the correlation between  $z$  and  $x$  is very small. Bound, Jaeger, and Baker (1995) discussed reasons why quarter of birth and  $u$  might be somewhat correlated. From equation (15.19), we see that this can lead to a substantial bias in the IV estimator.

When  $z$  and  $x$  are not correlated at all, things are especially bad, whether or not  $z$  is uncorrelated with  $u$ . The following example illustrates why we should always check to see if the endogenous explanatory variable is correlated with the IV candidate.

### EXAMPLE 15.3 Estimating the Effect of Smoking on Birth Weight

In Chapter 6, we estimated the effect of cigarette smoking on child birth weight. Without other explanatory variables, the model is

$$\log(bwght) = \beta_0 + \beta_1 \text{packs} + u, \quad [15.21]$$

where *packs* is the number of packs smoked by the mother per day. We might worry that *packs* is correlated with other health factors or the availability of good prenatal care, so that *packs* and  $u$  might be correlated. A possible instrumental variable for *packs* is the average price of cigarettes in the state of residence, *cigprice*. We will assume that *cigprice* and  $u$  are uncorrelated (even though state support for health care could be correlated with cigarette taxes).

If cigarettes are a typical consumption good, basic economic theory suggests that *packs* and *cigprice* are negatively correlated, so that *cigprice* can be used as an IV for *packs*. To check this, we regress *packs* on *cigprice*, using the data in BWGHT:

$$\begin{aligned} \widehat{\text{packs}} &= .067 + .0003 \text{cigprice} \\ &(.103) \quad (.0008) \\ n &= 1,388, R^2 = .0000, \bar{R}^2 = -.0006. \end{aligned}$$

This indicates no relationship between smoking during pregnancy and cigarette prices, which is perhaps not too surprising given the addictive nature of cigarette smoking.

Because *packs* and *cigprice* are not correlated, we should not use *cigprice* as an IV for *packs* in (15.21). But what happens if we do? The IV results would be

$$\begin{aligned} \widehat{\log(bwght)} &= 4.45 + 2.99 \text{packs} \\ &(.91) \quad (8.70) \\ n &= 1,388 \end{aligned}$$

(the reported  $R$ -squared is negative). The coefficient on *packs* is huge and of an unexpected sign. The standard error is also very large, so *packs* is not significant. But the estimates are meaningless because *cigprice* fails the one requirement of an IV that we can always test: assumption (15.5).

The previous example shows that IV estimation can produce strange results when the instrument relevance condition,  $\text{Corr}(z, x) \neq 0$ , fails. Of practically greater interest is the so-called problem of **weak instruments**, which is loosely defined as the problem of “low” (but not zero) correlation between  $z$  and  $x$ . In a particular application, it is difficult to define how low is too low, but recent theoretical research, supplemented by simulation studies, has shed considerable light on the issue. Staiger and Stock (1997) formalized the problem of weak instruments by modeling the correlation between  $z$  and  $x$  as a function of the sample size; in particular, the correlation is assumed to shrink to zero at the rate  $1/\sqrt{n}$ . Not surprisingly, the asymptotic distribution of the instrumental variables estimator is different compared with the usual asymptotics, where the correlation is assumed to be fixed and nonzero. One of the implications of the Stock–Staiger work is that the usual statistical inference, based on  $t$  statistics and the standard normal distribution, can be seriously misleading. We discuss this further in Section 15-3.

### 15-1c Computing $R$ -Squared after IV Estimation

Most regression packages compute an  $R$ -squared after IV estimation, using the standard formula:  $R^2 = 1 - \text{SSR}/\text{SST}$ , where SSR is the sum of squared *IV residuals* and SST is the total sum of squares of  $y$ . Unlike in the case of OLS, the  $R$ -squared from IV estimation can be negative because SSR for IV can actually be larger than SST. Although it does not really hurt to report the  $R$ -squared for IV estimation, it is not very useful, either. When  $x$  and  $u$  are correlated, we cannot decompose the variance of  $y$  into  $\beta_1^2 \text{Var}(x) + \text{Var}(u)$ , and so the  $R$ -squared has no natural interpretation. In addition, as we will discuss in Section 15-3, these  $R$ -squareds *cannot* be used in the usual way to compute  $F$  tests of joint restrictions.

If our goal was to produce the largest  $R$ -squared, we would always use OLS. IV methods are intended to provide better estimates of the *ceteris paribus* effect of  $x$  on  $y$  when  $x$  and  $u$  are correlated; goodness-of-fit is not a factor. A high  $R$ -squared resulting from OLS is of little comfort if we cannot consistently estimate  $\beta_1$ .

## 15-2 IV Estimation of the Multiple Regression Model

The IV estimator for the simple regression model is easily extended to the multiple regression case. We begin with the case where only one of the explanatory variables is correlated with the error. In fact, consider a standard linear model with two explanatory variables:

$$y_1 = \beta_0 + \beta_1 y_2 + \beta_2 z_1 + u_1. \quad [15.22]$$

We call this a **structural equation** to emphasize that we are interested in the  $\beta_j$ , which simply means that the equation is supposed to measure a causal relationship. We use a new notation here to distinguish endogenous from **exogenous variables**. The dependent variable  $y_1$  is clearly endogenous, as it is correlated with  $u_1$ . The variables  $y_2$  and  $z_1$  are the explanatory variables, and  $u_1$  is the error. As usual, we assume that the expected value of  $u_1$  is zero:  $E(u_1) = 0$ . We use  $z_1$  to indicate that this variable is exogenous in (15.22) ( $z_1$  is uncorrelated with  $u_1$ ). We use  $y_2$  to indicate that this variable is suspected of being correlated with  $u_1$ . We do not specify why  $y_2$  and  $u_1$  are correlated, but for now it is best to think of  $u_1$  as containing an omitted variable correlated with  $y_2$ . The notation in equation (15.22) originates in simultaneous equations models (which we cover in Chapter 16), but we use it more generally to easily distinguish exogenous from endogenous explanatory variables in a multiple regression model.

An example of (15.22) is

$$\log(\text{wage}) = \beta_0 + \beta_1 \text{educ} + \beta_2 \text{exper} + u_1, \quad [15.23]$$

where  $y_1 = \log(\text{wage})$ ,  $y_2 = \text{educ}$ , and  $z_1 = \text{exper}$ . In other words, we assume that *exper* is exogenous in (15.23), but we allow that *educ*—for the usual reasons—is correlated with  $u_1$ .

We know that if (15.22) is estimated by OLS, *all* of the estimators will be biased and inconsistent. Thus, we follow the strategy suggested in the previous section and seek an instrumental variable for  $y_2$ . Because  $z_1$  is assumed to be uncorrelated with  $u_1$ , can we use  $z_1$  as an instrument for  $y_2$ , assuming  $y_2$  and  $z_1$  are correlated? The answer is no. Because  $z_1$  itself appears as an explanatory variable in (15.22), it cannot serve as an instrumental variable for  $y_2$ . We need another exogenous variable—call it  $z_2$ —that does *not* appear in (15.22). Therefore, key assumptions are that  $z_1$  and  $z_2$  are uncorrelated with  $u_1$ ; we also assume that  $u_1$  has zero expected value, which is without loss of generality when the equation contains an intercept:

$$E(u_1) = 0, \text{Cov}(z_1, u_1) = 0, \text{and } \text{Cov}(z_2, u_1) = 0. \quad [15.24]$$

Given the zero mean assumption, the latter two assumptions are equivalent to  $E(z_1 u_1) = E(z_2 u_1) = 0$ , and so the method of moments approach suggests obtaining estimators  $\hat{\beta}_0$ ,  $\hat{\beta}_1$ , and  $\hat{\beta}_2$  by solving the sample counterparts of (15.24):

$$\begin{aligned} \sum_{i=1}^n (y_{i1} - \hat{\beta}_0 - \hat{\beta}_1 y_{i2} - \hat{\beta}_2 z_{i1}) &= 0 \\ \sum_{i=1}^n z_{i1} (y_{i1} - \hat{\beta}_0 - \hat{\beta}_1 y_{i2} - \hat{\beta}_2 z_{i1}) &= 0 \\ \sum_{i=1}^n z_{i2} (y_{i1} - \hat{\beta}_0 - \hat{\beta}_1 y_{i2} - \hat{\beta}_2 z_{i1}) &= 0. \end{aligned} \quad [15.25]$$

This is a set of three linear equations in the three unknowns  $\hat{\beta}_0$ ,  $\hat{\beta}_1$ , and  $\hat{\beta}_2$ , and it is easily solved given the data on  $y_1$ ,  $y_2$ ,  $z_1$ , and  $z_2$ . The estimators are called *instrumental variables estimators*. If we think  $y_2$  is exogenous and we choose  $z_2 = y_2$ , equations (15.25) are exactly the first order conditions for the OLS estimators; see equations (3.13).

We still need the instrumental variable  $z_2$  to be correlated with  $y_2$ , but the sense in which these two variables must be correlated is complicated by the presence of  $z_1$  in equation (15.22). We now need to state the assumption in terms of *partial* correlation. The easiest way to state the condition is to write the endogenous explanatory variable as a linear function of the exogenous variables and an error term:

$$y_2 = \pi_0 + \pi_1 z_1 + \pi_2 z_2 + v_2, \quad [15.26]$$

where, by construction,  $E(v_2) = 0$ ,  $\text{Cov}(z_1, v_2) = 0$ , and  $\text{Cov}(z_2, v_2) = 0$ , and the  $\pi_j$  are unknown parameters. The key identification condition [along with (15.24)] is that

$$\pi_2 \neq 0. \quad [15.27]$$

In other words, after partialling out  $z_1$ ,  $y_2$  and  $z_2$  are still correlated. This correlation can be positive or negative, but it cannot be zero. Testing (15.27) is easy: we estimate (15.26) by OLS and use a  $t$  test (possibly making it robust to heteroskedasticity). We should always test this assumption. Unfortunately, we cannot test that  $z_1$  and  $z_2$  are uncorrelated with  $u_1$ ; hopefully, we can make the case based on economic reasoning or introspection.

### GOING FURTHER 15.2

Suppose we wish to estimate the effect of marijuana usage on college grade point average. For the population of college seniors at a university, let *daysused* denote the number of days in the past month on which a student smoked marijuana and consider the structural equation

$$\text{colGPA} = \beta_0 + \beta_1 \text{daysused} + \beta_2 \text{SAT} + u.$$

(i) Let *perchs* denote the percentage of a student's high school graduating class that reported regular use of marijuana. If this is an IV candidate for *daysused*, write the reduced form for *daysused*. Do you think (15.27) is likely to be true?

(ii) Do you think *perchs* is truly exogenous in the structural equation? What problems might there be?

Equation (15.26) is an example of a **reduced form equation**, which means that we have written an endogenous variable in terms of exogenous variables. This name comes from simultaneous equations models—which we study in Chapter 16—but it is a useful concept whenever we have an endogenous explanatory variable. The name helps distinguish it from the structural equation (15.22).

Adding more **exogenous explanatory variables** to the model is straightforward. Write the structural model as

$$y_1 = \beta_0 + \beta_1 y_2 + \beta_2 z_1 + \cdots + \beta_k z_{k-1} + u_1, \quad [15.28]$$

where  $y_2$  is thought to be correlated with  $u_1$ . Let  $z_k$  be a variable not in (15.28) that is also exogenous. Therefore, we assume that

$$E(u_1) = 0, \text{Cov}(z_j, u_1) = 0, \quad j = 1, \dots, k. \quad [15.29]$$

Under (15.29),  $z_1, \dots, z_{k-1}$  are the exogenous variables appearing in (15.28). In effect, these act as their own instrumental variables in estimating the  $\beta_j$  in (15.28). The special case of  $k = 2$  is given in the equations in (15.25); along with  $z_2$ ,  $z_1$  appears in the set of moment conditions used to obtain the IV estimates. More generally,  $z_1, \dots, z_{k-1}$  are used in the moment conditions along with the instrumental variable for  $y_2$ ,  $z_k$ .

The reduced form for  $y_2$  is

$$y_2 = \pi_0 + \pi_1 z_1 + \cdots + \pi_{k-1} z_{k-1} + \pi_k z_k + v_2, \quad [15.30]$$

and we need some partial correlation between  $z_k$  and  $y_2$ :

$$\pi_k \neq 0. \quad [15.31]$$

Under (15.29) and (15.31),  $z_k$  is a valid IV for  $y_2$ . [We do not care about the remaining  $\pi_j$  in (15.30); some or all of them could be zero.] A minor additional assumption is that there are no perfect linear relationships among the exogenous variables; this is analogous to the assumption of no perfect collinearity in the context of OLS.

For standard statistical inference, we need to assume homoskedasticity of  $u_1$ . We give a careful statement of these assumptions in a more general setting in Section 15-3.

#### EXAMPLE 15.4 Using College Proximity as an IV for Education

Card (1995) used wage and education data for a sample of men in 1976 to estimate the return to education. He used a dummy variable for whether someone grew up near a four-year college (*nearc4*) as an instrumental variable for education. In a log(*wage*) equation, he included other standard controls: experience, a black dummy variable, dummy variables for living in an SMSA and living in the South, and a full set of regional dummy variables and an SMSA dummy for where the man was living in 1966. In order for *nearc4* to be a valid instrument, it must be uncorrelated with the error term in the wage equation—we assume this—and it must be partially correlated with *educ*. To check the latter requirement, we regress *educ* on *nearc4* and all of the exogenous variables appearing in the equation. (That is, we estimate the reduced form for *educ*.) Using the data in CARD, we obtain, in condensed form,

$$\begin{aligned} educ &= 16.64 + .320 \text{ nearc4} - .413 \text{ exper} + \cdots \\ &\quad (.24) \quad (.088) \quad (.034) \\ n &= 3,010, R^2 = .477. \end{aligned}$$

We are interested in the coefficient and  $t$  statistic on *nearc4*. The coefficient implies that in 1976, other things being fixed (experience, race, region, and so on), people who lived near a college in 1966 had, on average, about one-third of a year more education than those who did not grow up near a college. The  $t$  statistic on *nearc4* is 3.64, which gives a  $p$ -value that is zero in the first three decimals.

Therefore, if *nearc4* is uncorrelated with unobserved factors in the error term, we can use *nearc4* as an IV for *educ*.

The OLS and IV estimates are given in Table 15.1. Like the OLS standard errors, the reported IV standard errors employ a degrees-of-freedom adjustment in estimating the error variance. In some statistical packages the degrees-of-freedom adjustment is the default; in others it is not.

Interestingly, the IV estimate of the return to education is almost twice as large as the OLS estimate, but the standard error of the IV estimate is over 18 times larger than the OLS standard error. The 95% confidence interval for the IV estimate is between .024 and .239, which is a very wide range. The presence of larger confidence intervals is a price we must pay to get a consistent estimator of the return to education when we think *educ* is endogenous.

**TABLE 15.1** Dependent Variable:  $\log(\text{wage})$

Explanatory Variables	OLS	IV
<i>educ</i>	.075 (.003)	.132 (.055)
<i>exper</i>	.085 (.007)	.108 (.024)
<i>exper</i> <sup>2</sup>	−.0023 (.0003)	−.0023 (.0003)
<i>black</i>	−.199 (.018)	−.147 (.054)
<i>smsa</i>	.136 (.020)	.112 (.032)
<i>south</i>	−.148 (.026)	−.145 (.027)
Observations	3,010	3,010
<i>R</i> -squared	.300	.238
Other controls: <i>smsa66</i> , <i>reg662</i> , ..., <i>reg669</i>		

As discussed earlier, we should not make anything of the smaller *R*-squared in the IV estimation: by definition, the OLS *R*-squared will always be larger because OLS minimizes the sum of squared residuals.

It is worth noting, especially for studying the effects of policy interventions, that a reduced form equation exists for  $y_1$ , too. In the context of equation (15.28) with  $z_k$  an IV for  $y_2$ , the reduced form for  $y_1$  always has the form

$$y_1 = \gamma_0 + \gamma_1 z_1 + \cdots + \gamma_k z_k + e_1, \quad [15.32]$$

where  $\gamma_j = \beta_j + \beta_1 \pi_j$  for  $j < k$ ,  $\gamma_k = \beta_1 \pi_k$ , and  $e_1 = u_1 + \beta_1 v_2$ —as can be verified by plugging (15.30) into (15.28) and rearranging. Because the  $z_j$  are exogenous in (15.32), the  $\gamma_j$  can be consistently estimated by OLS. In other words, we regress  $y_1$  on all of the exogenous variables, including  $z_k$ , the IV for  $y_2$ . Only if we want to estimate  $\beta_1$  in (15.28) do we need to apply IV.

When  $y_2$  is a zero-one variable denoting participation and  $z_k$  is a zero-one variable representing *eligibility* for program participation—which is, hopefully, either randomized across individuals or, at most, a function of the other exogenous variables  $z_1, \dots, z_{k-1}$  (such as income)—the coefficient  $\gamma_k$  has an interesting interpretation. Rather than an estimate of the effect of the program itself, it is an



estimate of the effect of *offering* the program. Unlike  $\beta_1$  in (15.28)—which measures the effect of the program itself— $\gamma_k$  accounts for the possibility that some units made eligible will choose not to participate. In the program evaluation literature,  $\gamma_k$  is an example of an *intention-to-treat* parameter: it measures the effect of being made *eligible* and not the effect of actual participation. The intention-to-treat coefficient,  $\gamma_k = \beta_1 \pi_k$ , depends on the effect of participating,  $\beta_1$ , and the change (typically, increase) in the probability of participating due to being eligible,  $\pi_k$ . [When  $y_2$  is binary, equation (15.30) is a linear probability model, and therefore  $\pi_k$  measures the ceteris paribus change in probability that  $y_2 = 1$  as  $z_k$  switches from zero to one.]

## 15-3 Two Stage Least Squares

In the previous section, we assumed that we had a single endogenous explanatory variable ( $y_2$ ), along with one instrumental variable for  $y_2$ . It often happens that we have more than one exogenous variable that is excluded from the structural model and might be correlated with  $y_2$ , which means they are valid IVs for  $y_2$ . In this section, we discuss how to use multiple instrumental variables.

### 15-3a A Single Endogenous Explanatory Variable

Consider again the structural model (15.22), which has one endogenous and one exogenous explanatory variable. Suppose now that we have *two* exogenous variables excluded from (15.22):  $z_2$  and  $z_3$ . Our assumptions that  $z_2$  and  $z_3$  do not appear in (15.22) and are uncorrelated with the error  $u_1$  are known as **exclusion restrictions**.

If  $z_2$  and  $z_3$  are both correlated with  $y_2$ , we could just use each as an IV, as in the previous section. But then we would have two IV estimators, and neither of these would, in general, be efficient. Since each of  $z_1$ ,  $z_2$ , and  $z_3$  is uncorrelated with  $u_1$ , any linear combination is also uncorrelated with  $u_1$ , and therefore any linear combination of the exogenous variables is a valid IV. To find the best IV, we choose the linear combination that is most highly correlated with  $y_2$ . This turns out to be given by the reduced form equation for  $y_2$ . Write

$$y_2 = \pi_0 + \pi_1 z_1 + \pi_2 z_2 + \pi_3 z_3 + v_2, \quad [15.33]$$

where

$$E(v_2) = 0, \text{ Cov}(z_1, v_2) = 0, \text{ Cov}(z_2, v_2) = 0, \text{ and } \text{Cov}(z_3, v_2) = 0.$$

Then, the best IV for  $y_2$  (under the assumptions given in the chapter appendix) is the linear combination of the  $z_j$  in (15.33), which we call  $y_2^*$ :

$$y_2^* = \pi_0 + \pi_1 z_1 + \pi_2 z_2 + \pi_3 z_3. \quad [15.34]$$

For this IV not to be perfectly correlated with  $z_1$  we need at least one of  $\pi_2$  or  $\pi_3$  to be different from zero:

$$\pi_2 \neq 0 \text{ or } \pi_3 \neq 0. \quad [15.35]$$

This is the key identification assumption, once we assume the  $z_j$  are all exogenous. (The value of  $\pi_1$  is irrelevant.) The structural equation (15.22) is not identified if  $\pi_2 = 0$  and  $\pi_3 = 0$ . We can test  $H_0: \pi_2 = 0$  and  $\pi_3 = 0$  against (15.35) using an  $F$  statistic.

A useful way to think of (15.33) is that it breaks  $y_2$  into two pieces. The first is  $y_2^*$ ; this is the part of  $y_2$  that is uncorrelated with the error term,  $u_1$ . The second piece is  $v_2$ , and this part is possibly correlated with  $u_1$ —which is why  $y_2$  is possibly endogenous.

Given data on the  $z_j$ , we can compute  $y_2^*$  for each observation, provided we know the population parameters  $\pi_j$ . This is never true in practice. Nevertheless, as we saw in the previous section, we can

always estimate the reduced form by OLS. Thus, using the sample, we regress  $y_2$  on  $z_1$ ,  $z_2$ , and  $z_3$  and obtain the fitted values:

$$\hat{y}_2 = \hat{\pi}_0 + \hat{\pi}_1 z_1 + \hat{\pi}_2 z_2 + \hat{\pi}_3 z_3 \quad [15.36]$$

(that is, we have  $\hat{y}_{i2}$  for each  $i$ ). At this point, we should verify that  $z_2$  and  $z_3$  are jointly significant in (15.33) at a reasonably small significance level (no larger than 5%). If  $z_2$  and  $z_3$  are not jointly significant in (15.33), then we are wasting our time with IV estimation.

Once we have  $\hat{y}_2$ , we can use it as the IV for  $y_2$ . The three equations for estimating  $\beta_0$ ,  $\beta_1$ , and  $\beta_2$  are the first two equations of (15.25), with the third replaced by

$$\sum_{i=1}^n \hat{y}_{i2} (y_{i1} - \hat{\beta}_0 - \hat{\beta}_1 y_{i2} - \hat{\beta}_2 z_{i1}) = 0. \quad [15.37]$$

Solving the three equations in three unknowns gives us the IV estimators.

With multiple instruments, the IV estimator using  $\hat{y}_2$  as the instrument is also called the **two stage least squares (2SLS) estimator**. The reason is simple. Using the algebra of OLS, it can be shown that when we use  $\hat{y}_2$  as the IV for  $y_2$ , the IV estimates  $\hat{\beta}_0$ ,  $\hat{\beta}_1$ , and  $\hat{\beta}_2$  are *identical* to the OLS estimates from the regression of

$$y_1 \text{ on } \hat{y}_2 \text{ and } z_1. \quad [15.38]$$

In other words, we can obtain the 2SLS estimator in two stages. The **first stage** is to run the regression in (15.36), where we obtain the fitted values  $\hat{y}_2$ . The second stage is the OLS regression (15.38). Because we use  $\hat{y}_2$  in place of  $y_2$ , the 2SLS estimates can differ substantially from the OLS estimates.

Some economists like to interpret the regression in (15.38) as follows. The fitted value,  $\hat{y}_2$ , is the estimated version of  $y_2^*$ , and  $y_2^*$  is uncorrelated with  $u_1$ . Therefore, 2SLS first “purges”  $y_2$  of its correlation with  $u_1$  before doing the OLS regression in (15.38). We can show this by plugging  $y_2 = y_2^* + v_2$  into (15.22):

$$y_1 = \beta_0 + \beta_1 y_2^* + \beta_2 z_1 + u_1 + \beta_1 v_2. \quad [15.39]$$

Now, the composite error  $u_1 + \beta_1 v_2$  has zero mean and is uncorrelated with  $y_2^*$  and  $z_1$ , which is why the OLS regression in (15.38) works.

Most econometrics packages have special commands for 2SLS, so there is no need to perform the two stages explicitly. In fact, in most cases you should avoid doing the second stage manually, as the standard errors and test statistics obtained in this way are *not* valid. [The reason is that the error term in (15.39) includes  $v_2$ , but the standard errors involve the variance of  $u_1$  only.] Any regression software that supports 2SLS asks for the dependent variable, the list of explanatory variables (both exogenous and endogenous), and the entire list of instrumental variables (that is, all exogenous variables). The output is typically quite similar to that for OLS.

In model (15.28) with a single IV for  $y_2$ , the IV estimator from Section 15-2 is identical to the 2SLS estimator. Therefore, when we have one IV for each endogenous explanatory variable, we can call the estimation method IV or 2SLS.

Adding more exogenous variables changes very little. For example, suppose the wage equation is

$$\log(\text{wage}) = \beta_0 + \beta_1 \text{educ} + \beta_2 \text{exper} + \beta_3 \text{exper}^2 + u_1, \quad [15.40]$$

where  $u_1$  is uncorrelated with both  $\text{exper}$  and  $\text{exper}^2$ . Suppose that we also think mother’s and father’s educations are uncorrelated with  $u_1$ . Then, we can use both of these as IVs for  $\text{educ}$ . The reduced form (or first stage equation) equation for  $\text{educ}$  is

$$\text{educ} = \pi_0 + \pi_1 \text{exper} + \pi_2 \text{exper}^2 + \pi_3 \text{motheduc} + \pi_4 \text{fatheduc} + v_2, \quad [15.41]$$

and identification requires that  $\pi_3 \neq 0$  or  $\pi_4 \neq 0$  (or both, of course).

**EXAMPLE 15.5****Return to Education for Working Women**

We estimate equation (15.40) using the data in MROZ. First, we test  $H_0: \pi_3 = 0, \pi_4 = 0$  in (15.41) using an  $F$  test. The result is  $F = 124.76$ , and  $p\text{-value} = .0000$ . As expected,  $educ$  is (partially) correlated with parents' education.

When we estimate (15.40) by 2SLS, we obtain, in equation form,

$$\begin{aligned} \widehat{\log(\text{wage})} &= .048 + .061 \text{educ} + .044 \text{exper} - .0009 \text{exper}^2 \\ &\quad (.400) \quad (.031) \quad (.013) \quad (.0004) \\ n &= 428, R^2 = .136. \end{aligned}$$

The estimated return to education is about 6.1%, compared with an OLS estimate of about 10.8%. Because of its relatively large standard error, the 2SLS estimate is barely statistically significant at the 5% level against a two-sided alternative.

The assumptions needed for 2SLS to have the desired large sample properties are given in the chapter appendix, but it is useful to briefly summarize them here. If we write the structural equation as in (15.28),

$$y_1 = \beta_0 + \beta_1 y_2 + \beta_2 z_1 + \cdots + \beta_k z_{k-1} + u_1, \quad [15.42]$$

then we assume each  $z_j$  to be uncorrelated with  $u_1$ . In addition, we need at least one exogenous variable *not* in (15.42) that is partially correlated with  $y_2$ . This ensures consistency. For the usual 2SLS standard errors and  $t$  statistics to be asymptotically valid, we also need a homoskedasticity assumption: the variance of the structural error,  $u_1$ , cannot depend on any of the exogenous variables. For time series applications, we need more assumptions, as we will see in Section 15-7.

**15-3b Multicollinearity and 2SLS**

In Chapter 3, we introduced the problem of multicollinearity and showed how correlation among regressors can lead to large standard errors for the OLS estimates. Multicollinearity can be even more serious with 2SLS. To see why, the (asymptotic) variance of the 2SLS estimator of  $\beta_1$  can be approximated as

$$\sigma^2 / [\widehat{\text{SST}}_2 (1 - \hat{R}_2^2)], \quad [15.43]$$

where  $\sigma^2 = \text{Var}(u_1)$ ,  $\widehat{\text{SST}}_2$  is the total variation in  $\hat{y}_2$ , and  $\hat{R}_2^2$  is the  $R$ -squared from a regression of  $\hat{y}_2$  on all other exogenous variables appearing in the structural equation. There are two reasons why the variance of the 2SLS estimator is larger than that for OLS. First,  $\hat{y}_2$ , by construction, has less variation than  $y_2$ . (Remember: Total sum of squares = explained sum of squares + residual sum of squares; the variation in  $y_2$  is the total sum of squares, while the variation in  $\hat{y}_2$  is the explained sum of squares from the first stage regression.) Second, the correlation between  $\hat{y}_2$  and the exogenous variables in (15.42) is often much higher than the correlation between  $y_2$  and these variables. This essentially defines the multicollinearity problem in 2SLS.

As an illustration, consider Example 15.4. When  $educ$  is regressed on the exogenous variables in Table 15.1 (not including  $nearc4$ ),  $R$ -squared = .475; this is a moderate degree of multicollinearity, but the important thing is that the OLS standard error on  $\hat{\beta}_{educ}$  is quite small. When we obtain the first stage fitted values,  $\widehat{educ}$ , and regress these on the exogenous variables in Table 15.1,  $R$ -squared = .995, which indicates a very high degree of multicollinearity between  $\widehat{educ}$  and the remaining exogenous variables in the table. (This high  $R$ -squared is not too surprising because  $\widehat{educ}$  is a function of all the exogenous variables in Table 15.1, plus  $nearc4$ .) Equation (15.43) shows that an  $\hat{R}_2^2$  close to one can result in a very large standard error for the 2SLS estimator. But as with OLS, a large sample size can help offset a large  $\hat{R}_2^2$ .

## 15-3c Detecting Weak Instruments

In Section 15-1 we briefly discussed the problem of weak instruments. We focused on equation (15.19), which demonstrates how a small correlation between the instrument and error can lead to very large inconsistency (and therefore bias) if the instrument,  $z$ , also has little correlation with the explanatory variable,  $x$ . The same problem can arise in the context of the multiple equation model in equation (15.42), whether we have one instrument for  $y_2$  or more instruments than we need.

We also mentioned the findings of Staiger and Stock (1997), and we now discuss the practical implications of this research in a bit more depth. Importantly, Staiger and Stock study the case of where all instrumental variables are exogenous. With the exogeneity requirement satisfied by the instruments, they focus on the case where the instruments are weakly correlated with  $y_2$ , and they study the validity of standard errors, confidence intervals, and  $t$  statistics involving the coefficient  $\beta_1$  on  $y_2$ . The mechanism they used to model weak correlation led to an important finding: even with very large sample sizes the 2SLS estimator can be biased and a distribution that is very different from standard normal.

Building on Staiger and Stock (1997), Stock and Yogo (2005) (SY for short) proposed methods for detecting situations where weak instruments will lead to substantial bias and distorted statistical inference. Conveniently, Stock and Yogo obtained rules concerning the size of the  $t$  statistic (with one instrument) or the  $F$  statistic (with more than one instrument) from the first-stage regression. The theory is much too involved to pursue here. Instead, we describe some simple rules of thumb proposed by Stock and Yogo that are easy to implement.

The key implication of the SY work is that one needs more than just a statistical rejection of the null hypothesis in the first stage regression at the usual significance levels. For example, in equation (15.6), it is not enough to reject the null hypothesis stated in (15.7) at the 5% significance level. Using bias calculations for the instrumental variables estimator, SY recommend that one can proceed with the usual IV inference if the first-stage  $t$  statistic has absolute value larger than  $\sqrt{10} \approx 3.2$ . Readers will recognize this value as being well above the 95<sup>th</sup> percentile of the standard normal distribution, 1.96, which is what we would use for a standard 5% significance level. This same rule of thumb applies in the multiple regression model with a single endogenous explanatory variable,  $y_2$ , and a single instrumental variable,  $z_k$ . In particular, the  $t$  statistic in testing hypothesis (15.31) should be at least 3.2 in absolute value.

SY cover the case of 2SLS, too. In this case, we must focus on the first-stage  $F$  statistic for exclusion of the instrumental variables for  $y_2$ , and the SY rule is  $F > 10$ . (Notice this is the same rule based on the  $t$  statistic when there is only one instrument, as  $t^2 = F$ .) For example, consider equation (15.34), where we have two instruments for  $y_2$ ,  $z_2$  and  $z_3$ . Then the  $F$  statistic for the null hypothesis

$$H_0: \pi_2 = 0, \pi_3 = 0$$

should have  $F > 10$ . Remember, this is not the overall  $F$  statistic for all of the exogenous variables in (15.34). We test only the coefficients on the proposed IVs for  $y_2$ , that is, the exogenous variables that do not appear in (15.22). In Example 15.5 the relevant  $F$  statistic is 124.76, which is well above 10, implying that we do not have to worry about weak instruments. (Of course, the exogeneity of the parents' education variables is in doubt.)

The rule of thumb of requiring the  $F$  statistic to be larger than 10 tends to work well and is easy to remember. However, like all rules of thumb involving statistical inference, it makes no sense to use 10 as a knife-edge cutoff. For example, one can probably proceed if  $F = 9.94$ , as it is pretty close to 10. The rule of thumb should be used as a guideline. SY have more detailed suggestions for cases where there are many instruments for  $y_2$ , say five or more.

A more complicated issue is what happens if there is heteroskedasticity in either the equation of interest, (15.28), or the reduced form (first stage) for the endogenous explanatory variables, (15.30). Stock and Yogo (2005) did not allow for heteroskedasticity in either equation (or, in a time series or panel context, serial correlation). It makes sense that the requirements for the first-stage  $t$  or  $F$  statistic would be more stringent. Work by Olea and Pflueger (2013) suggests this is the case: the first-stage  $F$  might need to be more like 20 rather than 10 in order to ensure the instruments are sufficiently strong. This is an ongoing area of research.

### 15-3d Multiple Endogenous Explanatory Variables

Two stage least squares can also be used in models with more than one endogenous explanatory variable. For example, consider the model

$$y_1 = \beta_0 + \beta_1 y_2 + \beta_2 y_3 + \beta_3 z_1 + \beta_4 z_2 + \beta_5 z_3 + u_1, \quad [15.44]$$

where  $E(u_1) = 0$  and  $u_1$  is uncorrelated with  $z_1$ ,  $z_2$ , and  $z_3$ . The variables  $y_2$  and  $y_3$  are endogenous explanatory variables: each may be correlated with  $u_1$ .

To estimate (15.44) by 2SLS, we need *at least two* exogenous variables that do not appear in (15.44) but that are correlated with  $y_2$  and  $y_3$ . Suppose we have two excluded exogenous variables, say  $z_4$  and  $z_5$ . Then, from our analysis of a single endogenous explanatory variable, we need either  $z_4$  or  $z_5$  to appear in each reduced form for  $y_2$  and  $y_3$ . (As before, we can use  $F$  statistics to test this.) Although this is necessary for identification, unfortunately, it is not sufficient. Suppose that  $z_4$  appears in each reduced form, but  $z_5$  appears in neither. Then, we do not really have two exogenous variables partially correlated with  $y_2$  and  $y_3$ . Two stage least squares will not produce consistent estimators of the  $\beta_j$ .

Generally, when we have more than one endogenous explanatory variable in a regression model, identification can fail in several complicated ways. But we can easily state a necessary condition for identification, which is called the **order condition**.

#### GOING FURTHER 15.3

The following model explains violent crime rates, at the city level, in terms of a binary variable for whether gun control laws exist and other controls:

$$\begin{aligned} \text{violent} = & \beta_0 + \beta_1 \text{guncontrol} + \beta_2 \text{unem} \\ & + \beta_3 \text{popul} + \beta_4 \text{percblick} \\ & + \beta_5 \text{age18\_21} + \dots \end{aligned}$$

Some researchers have estimated similar equations using variables such as the number of National Rifle Association members in the city and the number of subscribers to gun magazines as instrumental variables for *gun-control* [see, for example, Kleck and Patterson (1993)]. Are these convincing instruments?

**Order Condition for Identification of an Equation.** We need at least as many excluded exogenous variables as there are included endogenous explanatory variables in the structural equation. The order condition is simple to check, as it only involves counting endogenous and exogenous variables. The sufficient condition for identification is called the **rank condition**. We have seen special cases of the rank condition before—for example, in the discussion surrounding equation (15.35). A general statement of the rank condition requires matrix algebra and is beyond the scope of this text. [See Wooldridge (2010, Chapter 5).] It is even more difficult to obtain diagnostics for weak instruments.

### 15-3e Testing Multiple Hypotheses after 2SLS Estimation

We must be careful when testing multiple hypotheses in a model estimated by 2SLS. It is tempting to use either the sum of squared residuals or the  $R$ -squared form of the  $F$  statistic, as we learned with OLS in Chapter 4. The fact that the  $R$ -squared in 2SLS can be negative suggests that the usual way of computing  $F$  statistics might not be appropriate; this is the case. In fact, if we use the 2SLS residuals to compute the SSRs for both the restricted and unrestricted models, there is no guarantee that  $SSR_r \geq SSR_{ur}$ ; if the reverse is true, the  $F$  statistic would be negative.

It is possible to combine the sum of squared residuals from the second stage regression [such as (15.38)] with  $SSR_{ur}$  to obtain a statistic with an approximate  $F$  distribution in large samples. Because many econometrics packages have simple-to-use test commands that can be used to test multiple hypotheses after 2SLS estimation, we omit the details. Davidson and MacKinnon (1993) and Wooldridge (2010, Chapter 5) contain discussions of how to compute  $F$ -type statistics for 2SLS.

## 15-4 IV Solutions to Errors-in-Variables Problems

In the previous sections, we presented the use of instrumental variables as a way to solve the omitted variables problem, but they can also be used to deal with the measurement error problem. As an illustration, consider the model

$$y = \beta_0 + \beta_1 x_1^* + \beta_2 x_2 + u, \quad [15.45]$$

where  $y$  and  $x_2$  are observed but  $x_1^*$  is not. Let  $x_1$  be an observed measurement of  $x_1^*$ :  $x_1 = x_1^* + e_1$ , where  $e_1$  is the measurement error. In Chapter 9, we showed that correlation between  $x_1$  and  $e_1$  causes OLS, where  $x_1$  is used in place of  $x_1^*$ , to be biased and inconsistent. We can see this by writing

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + (u - \beta_1 e_1). \quad [15.46]$$

If the classical errors-in-variables (CEV) assumptions hold, the bias in the OLS estimator of  $\beta_1$  is toward zero. Without further assumptions, we can do nothing about this.

In some cases, we can use an IV procedure to solve the measurement error problem. In (15.45), we assume that  $u$  is uncorrelated with  $x_1^*$ ,  $x_1$ , and  $x_2$ ; in the CEV case, we assume that  $e_1$  is uncorrelated with  $x_1^*$  and  $x_2$ . These imply that  $x_2$  is exogenous in (15.46), but that  $x_1$  is correlated with  $e_1$ . What we need is an IV for  $x_1$ . Such an IV must be correlated with  $x_1$ , uncorrelated with  $u$ —so that it can be excluded from (15.45)—and uncorrelated with the measurement error,  $e_1$ .

One possibility is to obtain a second measurement on  $x_1^*$ , say,  $z_1$ . Because it is  $x_1^*$  that affects  $y$ , it is only natural to assume that  $z_1$  is uncorrelated with  $u$ . If we write  $z_1 = x_1^* + a_1$ , where  $a_1$  is the measurement error in  $z_1$ , then we must assume that  $a_1$  and  $e_1$  are uncorrelated. In other words,  $x_1$  and  $z_1$  both mismeasure  $x_1^*$ , but their measurement errors are uncorrelated. Certainly,  $x_1$  and  $z_1$  are correlated through their dependence on  $x_1^*$ , so we can use  $z_1$  as an IV for  $x_1$ .

Where might we get two measurements on a variable? Sometimes, when a group of workers is asked for their annual salary, their employers can provide a second measure. For married couples, each spouse can independently report the level of savings or family income. In the Ashenfelter and Krueger (1994) study cited in Section 14-3, each twin was asked about his or her sibling's years of education; this gives a second measure that can be used as an IV for self-reported education in a wage equation. (Ashenfelter and Krueger combined differencing and IV to account for the omitted ability problem as well; more on this in Section 15-8.) Generally, though, having two measures of an explanatory variable is rare.

An alternative is to use other exogenous variables as IVs for a potentially mismeasured variable. For example, our use of *motheduc* and *fatheduc* as IVs for *educ* in Example 15.5 can serve this purpose. If we think that  $educ = educ^* + e_1$ , then the IV estimates in Example 15.5 do not suffer from measurement error if *motheduc* and *fatheduc* are uncorrelated with the measurement error,  $e_1$ . This is probably more reasonable than assuming *motheduc* and *fatheduc* are uncorrelated with ability, which is contained in  $u$  in (15.45).

IV methods can also be adopted when using things like test scores to control for unobserved characteristics. In Section 9-2, we showed that, under certain assumptions, proxy variables can be used to solve the omitted variables problem. In Example 9.3, we used IQ as a proxy variable for unobserved ability. This simply entails adding IQ to the model and performing an OLS regression. But there is an alternative that works when IQ does not fully satisfy the proxy variable assumptions. To illustrate, write a wage equation as

$$\log(\text{wage}) = \beta_0 + \beta_1 \text{educ} + \beta_2 \text{exper} + \beta_3 \text{exper}^2 + \text{abil} + u, \quad [15.47]$$

where we again have the omitted ability problem. But we have two test scores that are *indicators* of ability. We assume that the scores can be written as

$$\text{test}_1 = \gamma_1 \text{abil} + e_1$$

and

$$\text{test}_2 = \delta_1 \text{abil} + e_2,$$



where  $\gamma_1 > 0$ ,  $\delta_1 > 0$ . Since it is ability that affects wage, we can assume that  $test_1$  and  $test_2$  are uncorrelated with  $u$ . If we write  $abil$  in terms of the first test score and plug the result into (15.47), we get

$$\log(wage) = \beta_0 + \beta_1 educ + \beta_2 exper + \beta_3 exper^2 + \alpha_1 test_1 + (u - \alpha_1 e_1), \quad [15.48]$$

where  $\alpha_1 = 1/\gamma_1$ . Now, if we assume that  $e_1$  is uncorrelated with all the explanatory variables in (15.47), including  $abil$ , then  $e_1$  and  $test_1$  must be correlated. [Notice that  $educ$  is not endogenous in (15.48); however,  $test_1$  is.] This means that estimating (15.48) by OLS will produce inconsistent estimators of the  $\beta_j$  (and  $\alpha_1$ ). Under the assumptions we have made,  $test_1$  does not satisfy the proxy variable assumptions.

If we assume that  $e_2$  is also uncorrelated with all the explanatory variables in (15.47) and that  $e_1$  and  $e_2$  are uncorrelated, then  $e_1$  is uncorrelated with the second test score,  $test_2$ . Therefore,  $test_2$  can be used as an IV for  $test_1$ .

### EXAMPLE 15.6 Using Two Test Scores as Indicators of Ability

We use the data in WAGE2 to implement the preceding procedure, where  $IQ$  plays the role of the first test score and  $KWW$  (knowledge of the world of work) is the second test score. The explanatory variables are the same as in Example 9.3: *educ*, *exper*, *tenure*, *married*, *south*, *urban*, and *black*. Rather than adding  $IQ$  and doing OLS, as in column (2) of Table 9.2, we add  $IQ$  and use  $KWW$  as its instrument. The coefficient on *educ* is .025 (se = .017). This is a low estimate, and it is not statistically different from zero. This is a puzzling finding, and it suggests that one of our assumptions fails; perhaps  $e_1$  and  $e_2$  are correlated.

## 15-5 Testing for Endogeneity and Testing Overidentifying Restrictions

In this section, we describe two important tests in the context of instrumental variables estimation.

### 15-5a Testing for Endogeneity

The 2SLS estimator is less efficient than OLS when the explanatory variables are exogenous; as we have seen, the 2SLS estimates can have very large standard errors. Therefore, it is useful to have a test for endogeneity of an explanatory variable that shows whether 2SLS is even necessary. Obtaining such a test is rather simple.

To illustrate, suppose we have a single suspected endogenous variable,

$$y_1 = \beta_0 + \beta_1 y_2 + \beta_2 z_1 + \beta_3 z_2 + u_1, \quad [15.49]$$

where  $z_1$  and  $z_2$  are exogenous. We have two additional exogenous variables,  $z_3$  and  $z_4$ , which do not appear in (15.49). If  $y_2$  is uncorrelated with  $u_1$ , we should estimate (15.49) by OLS. How can we test this? Hausman (1978) suggested directly comparing the OLS and 2SLS estimates and determining whether the differences are statistically significant. After all, both OLS and 2SLS are consistent if all variables are exogenous. If 2SLS and OLS differ significantly, we conclude that  $y_2$  must be endogenous (maintaining that the  $z_j$  are exogenous).

It is a good idea to compute OLS and 2SLS to see if the estimates are practically different. To determine whether the differences are statistically significant, it is easier to use a regression test. This is based on estimating the reduced form for  $y_2$ , which in this case is

$$y_2 = \pi_0 + \pi_1 z_1 + \pi_2 z_2 + \pi_3 z_3 + \pi_4 z_4 + v_2. \quad [15.50]$$

Now, since each  $z_j$  is uncorrelated with  $u_1$ ,  $y_2$  is uncorrelated with  $u_1$  if, and only if,  $v_2$  is uncorrelated with  $u_1$ ; this is what we wish to test. Write  $u_1 = \delta_1 v_2 + e_1$ , where  $e_1$  is uncorrelated with  $v_2$  and has zero mean. Then,  $u_1$  and  $v_2$  are uncorrelated if, and only if,  $\delta_1 = 0$ . The easiest way to test this is to include  $v_2$  as an additional regressor in (15.49) and to do a  $t$  test. There is only one problem with implementing this:  $v_2$  is not observed, because it is the error term in (15.50). Because we can estimate the reduced form for  $y_2$  by OLS, we can obtain the reduced form residuals,  $\hat{v}_2$ . Therefore, we estimate

$$y_1 = \beta_0 + \beta_1 y_2 + \beta_2 z_1 + \beta_3 z_2 + \delta_1 \hat{v}_2 + \text{error} \quad [15.51]$$

by OLS and test  $H_0: \delta_1 = 0$  using a  $t$  statistic. If we reject  $H_0$  at a small significance level, we conclude that  $y_2$  is endogenous because  $v_2$  and  $u_1$  are correlated.

#### Testing for Endogeneity of a Single Explanatory Variable:

(i) Estimate the reduced form for  $y_2$  by regressing it on *all* exogenous variables (including those in the structural equation and the additional IVs). Obtain the residuals,  $\hat{v}_2$ .

(ii) Add  $\hat{v}_2$  to the structural equation (which includes  $y_2$ ) and test for significance of  $\hat{v}_2$  using an OLS regression. If the coefficient on  $\hat{v}_2$  is statistically different from zero, we conclude that  $y_2$  is indeed endogenous. We might want to use a heteroskedasticity-robust  $t$  test.

#### EXAMPLE 15.7 Return to Education for Working Women

We can test for endogeneity of *educ* in (15.40) by obtaining the residuals  $\hat{v}_2$  from estimating the reduced form (15.41)—using only working women—and including these in (15.40). When we do this, the coefficient on  $\hat{v}_2$  is  $\hat{\delta}_1 = .058$ , and  $t = 1.67$ . This is moderate evidence of positive correlation between  $u_1$  and  $v_2$ . It is probably a good idea to report both estimates because the 2SLS estimate of the return to education (6.1%) is well below the OLS estimate (10.8%).

An interesting feature of the regression from step (ii) of the test for endogeneity is that the coefficient estimates on all explanatory variables (except, of course,  $\hat{v}_2$ ) are identical to the 2SLS estimates. For example, estimating (15.51) by OLS produces the same  $\hat{\beta}_j$  as estimating (15.49) by 2SLS. One benefit of this equivalence is that it provides an easy check on whether you have done the proper regression in testing for endogeneity. But it also gives a different, useful interpretation of 2SLS: adding  $\hat{v}_2$  to the original equation as an explanatory variable, and applying OLS, clears up the endogeneity of  $y_2$ . So, when we start by estimating (15.49) by OLS, we can quantify the importance of allowing  $y_2$  to be endogenous by seeing how much  $\hat{\beta}_1$  changes when  $\hat{v}_2$  is added to the equation. Irrespective of the outcome of the statistical tests, we can see whether the change in  $\hat{\beta}_1$  is expected and is practically significant.

If, in the end, the 2SLS estimates are chosen, one should obtain the standard errors using built-in 2SLS routines rather than those from regression (15.51). The standard errors obtained from the OLS regression (15.51) are valid only under the null hypothesis  $\delta_1 = 0$ .

We can also test for endogeneity of multiple explanatory variables. For each suspected endogenous variable, we obtain the reduced form residuals, as in part (i). Then, we test for joint significance of these residuals in the structural equation, using an  $F$  test. Joint significance indicates that at least one suspected explanatory variable is endogenous. The number of exclusion restrictions tested is the number of suspected endogenous explanatory variables.

### 15-5b Testing Overidentification Restrictions

When we introduced the simple instrumental variables estimator in Section 15-1, we emphasized that the instrument must satisfy two requirements: it must be uncorrelated with the error (exogeneity) and

correlated with the endogenous explanatory variable (relevance). We have now seen that, even in models with additional explanatory variables, the second requirement can be tested using a  $t$  test (with just one instrument) or an  $F$  test (when there are multiple instruments). In the context of the simple IV estimator, we noted that the exogeneity requirement cannot be tested. However, if we have more instruments than we need, we can effectively test whether some of them are uncorrelated with the structural error.

As a specific example, again consider equation (15.49) with two instrumental variables for  $y_2$ ,  $z_3$ , and  $z_4$ . Remember,  $z_1$  and  $z_2$  essentially act as their own instruments. Because we have two instruments for  $y_2$ , we can estimate (15.49) using, say, only  $z_3$  as an IV for  $y_2$ ; let  $\check{\beta}_1$  be the resulting IV estimator of  $\beta_1$ . Then, we can estimate (15.49) using only  $z_4$  as an IV for  $y_2$ ; call this IV estimator  $\tilde{\beta}_1$ . If all  $z_j$  are exogenous, and if  $z_3$  and  $z_4$  are each partially correlated with  $y_2$ , then  $\check{\beta}_1$  and  $\tilde{\beta}_1$  are both consistent for  $\beta_1$ . Therefore, if our logic for choosing the instruments is sound,  $\check{\beta}_1$  and  $\tilde{\beta}_1$  should differ only by sampling error. Hausman (1978) proposed basing a test of whether  $z_3$  and  $z_4$  are both exogenous on the difference,  $\check{\beta}_1 - \tilde{\beta}_1$ . Shortly, we will provide a simpler way to obtain a valid test, but, before doing so, we should understand how to interpret the outcome of the test.

If we conclude that  $\check{\beta}_1$  and  $\tilde{\beta}_1$  are statistically different from one another, then we have no choice but to conclude that either  $z_3$ ,  $z_4$ , or both fail the exogeneity requirement. Unfortunately, we cannot know which is the case (unless we simply assert from the beginning that, say,  $z_3$  is exogenous). For example, if  $y_2$  denotes years of schooling in a log wage equation,  $z_3$  is mother's education, and  $z_4$  is father's education, a statistically significant difference in the two IV estimators implies that one or both of the parents' education variables are correlated with  $u_1$  in (15.54).

Certainly, rejecting that one's instruments are exogenous is serious and requires a new approach. But the more serious, and subtle, problem in comparing IV estimates is that they may be similar even though both instruments fail the exogeneity requirement. In the previous example, it seems likely that if mother's education is positively correlated with  $u_1$ , then so is father's education. Therefore, the two IV estimates may be similar even though each is inconsistent. In effect, because the IVs in this example are chosen using similar reasoning, their separate use in IV procedures may very well lead to similar estimates that are nevertheless both inconsistent. The point is that we should not feel especially comfortable if our IV procedures pass the Hausman test.

Another problem with comparing two IV estimates is that often they may seem practically different yet, statistically, we cannot reject the null hypothesis that they are consistent for the same population parameter. For example, in estimating (15.40) by IV using *motheduc* as the only instrument, the coefficient on *educ* is .049 (.037). If we use only *fatheduc* as the IV for *educ*, the coefficient on *educ* is .070 (.034). [Perhaps not surprisingly, the estimate using both parents' education as IVs is in between these two, .061 (.031).] For policy purposes, the difference between 5% and 7% for the estimated return to a year of schooling is substantial. Yet, as shown in Example 15.8, the difference is not statistically significant.

The procedure of comparing different IV estimates of the same parameter is an example of testing **overidentifying restrictions**. The general idea is that we have more instruments than we need to estimate the parameters consistently. In the previous example, we had one more instrument than we need, and this results in one overidentifying restriction that can be tested. In the general case, suppose that we have  $q$  more instruments than we need. For example, with one endogenous explanatory variable,  $y_2$ , and three proposed instruments for  $y_2$ , we have  $q = 3 - 1 = 2$  overidentifying restrictions. When  $q$  is two or more, comparing several IV estimates is cumbersome. Instead, we can easily compute a test statistic based on the 2SLS residuals. The idea is that, if all instruments are exogenous, the 2SLS residuals should be uncorrelated with the instruments, up to sampling error. But if there are  $k + 1$  parameters and  $k + 1 + q$  instruments, the 2SLS residuals have a zero mean and are identically uncorrelated with  $k$  linear combinations of the instruments. (This algebraic fact contains, as a special case, the fact that the OLS residuals have a zero mean and are uncorrelated with the  $k$  explanatory variables.) Therefore, the test checks whether the 2SLS residuals are correlated with  $q$  linear functions of the instruments, and we need not decide on the functions; the test does that for us automatically.

The following regression-based test is valid when the homoskedasticity assumption, listed as Assumption 2SLS.5 in the chapter appendix, holds.

#### Testing Overidentifying Restrictions:

- (i) Estimate the structural equation by 2SLS and obtain the 2SLS residuals,  $\hat{u}_1$ .
- (ii) Regress  $\hat{u}_1$  on *all exogenous* variables. Obtain the  $R$ -squared, say,  $R_1^2$ .
- (iii) Under the null hypothesis that all IVs are uncorrelated with  $u_1$ ,  $nR_1^2 \stackrel{a}{\sim} \chi_q^2$ , where  $q$  is the number of instrumental variables from outside the model minus the total number of endogenous explanatory variables. If  $nR_1^2$  exceeds (say) the 5% critical value in the  $\chi_q^2$  distribution, we reject  $H_0$  and conclude that at least some of the IVs are not exogenous.

### EXAMPLE 15.8 Return to Education for Working Women

When we use *motheduc* and *fatheduc* as IVs for *educ* in (15.40), we have a single overidentifying restriction. Regressing the 2SLS residuals  $\hat{u}_1$  on *exper*, *exper*<sup>2</sup>, *motheduc*, and *fatheduc* produces  $R_1^2 = .0009$ . Therefore,  $nR_1^2 = 428(.0009) = .3852$ , which is a very small value in a  $\chi_1^2$  distribution ( $p$ -value = .535). Therefore, the parents' education variables pass the overidentification test. When we add husband's education to the IV list, we get two overidentifying restrictions, and  $nR_1^2 = 1.11$  ( $p$ -value = .574). Subject to the preceding cautions, it seems reasonable to add *huseduc* to the IV list, as this reduces the standard error of the 2SLS estimate: the 2SLS estimate on *educ* using all three instruments is .080 ( $se = .022$ ), so this makes *educ* much more significant than when *huseduc* is not used as an IV ( $\hat{\beta}_{educ} = .061$ ,  $se = .031$ ).

When  $q = 1$ , a natural question is: How does the test obtained from the regression-based procedure compare with a test based on directly comparing the estimates? In fact, the two procedures are asymptotically the same. As a practical matter, it makes sense to compute the two IV estimates to see how they differ. More generally, when  $q \geq 2$ , one can compare the 2SLS estimates using all IVs to the IV estimates using single instruments. By doing so, one can see if the various IV estimates are practically different, whether or not the overidentification test rejects or fails to reject.

In the previous example, we alluded to a general fact about 2SLS: under the standard 2SLS assumptions, adding instruments to the list improves the asymptotic efficiency of the 2SLS. But this requires that any new instruments are in fact exogenous—otherwise, 2SLS will not even be consistent—and it is only an asymptotic result. With the typical sample sizes available, adding too many instruments—that is, increasing the number of overidentifying restrictions—can cause severe biases in 2SLS. A detailed discussion would take us too far afield. A nice illustration is given by Bound, Jaeger, and Baker (1995), who argue that the 2SLS estimates of the return to education obtained by Angrist and Krueger (1991), using many instrumental variables, are likely to be seriously biased (even with hundreds of thousands of observations!).

The overidentification test can be used whenever we have more instruments than we need. If we have just enough instruments, the model is said to be *just identified*, and the  $R$ -squared in part (ii) will be identically zero. As we mentioned earlier, we cannot test exogeneity of the instruments in the just identified case.

The test can be made robust to heteroskedasticity of arbitrary form; for details, see Wooldridge (2010, Chapter 5).

## 15-6 2SLS with Heteroskedasticity

Heteroskedasticity in the context of 2SLS raises essentially the same issues as with OLS. Most importantly, it is possible to obtain standard errors and test statistics that are (asymptotically) robust to heteroskedasticity of arbitrary and unknown form. In fact, expression (8.4) continues to be valid if the  $\hat{r}_{ij}$  are obtained as the residuals from regressing  $\hat{x}_{ij}$  on the other  $\hat{x}_{ih}$ , where the “ $\hat{\cdot}$ ” denotes fitted values

from the first stage regressions (for endogenous explanatory variables). Wooldridge (2010, Chapter 5) contains more details. Some software packages do this routinely.

We can also test for heteroskedasticity, using an analog of the Breusch-Pagan test that we covered in Chapter 8. Let  $\hat{u}$  denote the 2SLS residuals and let  $z_1, z_2, \dots, z_m$  denote all the exogenous variables (including those used as IVs for the endogenous explanatory variables). Then, under reasonable assumptions [spelled out, for example, in Wooldridge (2010, Chapter 5)], an asymptotically valid statistic is the usual  $F$  statistic for joint significance in a regression of  $\hat{u}^2$  on  $z_1, z_2, \dots, z_m$ . The null hypothesis of homoskedasticity is rejected if the  $z_j$  are jointly significant.

If we apply this test to Example 15.8, using *motheduc*, *fatheduc*, and *huseduc* as instruments for *educ*, we obtain  $F_{5,422} = 2.53$  and  $p\text{-value} = .029$ . This is evidence of heteroskedasticity at the 5% level. We might want to compute heteroskedasticity-robust standard errors to account for this.

If we know how the error variance depends on the exogenous variables, we can use a weighted 2SLS procedure, essentially the same as in Section 8-4. After estimating a model for  $\text{Var}(u|z_1, z_2, \dots, z_m)$ , we divide the dependent variable, the explanatory variables, and all the instrumental variables for observation  $i$  by  $\sqrt{\hat{h}_i}$ , where  $\hat{h}_i$  denotes the estimated variance. (The constant, which is both an explanatory variable and an IV, is divided by  $\sqrt{\hat{h}_i}$ ; see Section 8-4.) Then, we apply 2SLS on the transformed equation using the transformed instruments.

## 15-7 Applying 2SLS to Time Series Equations

When we apply 2SLS to time series data, many of the considerations that arose for OLS in Chapters 10, 11, and 12 are relevant. Write the structural equation for each time period as

$$y_t = \beta_0 + \beta_1 x_{t1} + \dots + \beta_k x_{tk} + u_t, \quad [15.52]$$

where one or more of the explanatory variables  $x_{tj}$  might be correlated with  $u_t$ . Denote the set of exogenous variables by  $z_{t1}, \dots, z_{tm}$ :

$$E(u_t) = 0, \text{Cov}(z_{tj}, u_t) = 0, \quad j = 1, \dots, m.$$

### GOING FURTHER 15.4

A model to test the effect of growth in government spending on growth in output is

$$gGDP_t = \beta_0 + \beta_1 gGOV_t + \beta_2 INVRAT_t + \beta_3 gLAB_t + u_t,$$

where  $g$  indicates growth,  $GDP$  is real gross domestic product,  $GOV$  is real government spending,  $INVRAT$  is the ratio of gross domestic investment to GDP, and  $LAB$  is the size of the labor force. [See equation (6) in Ram (1986).] Under what assumptions would a dummy variable indicating whether the president in year  $t - 1$  is a Republican be a suitable IV for  $gGOV_t$ ?

Any exogenous explanatory variable is also a  $z_{tj}$ . For identification, it is necessary that  $m \geq k$  (we have as many exogenous variables as explanatory variables).

The mechanics of 2SLS are identical for time series or cross-sectional data, but for time series data the statistical properties of 2SLS depend on the trending and correlation properties of the underlying sequences. In particular, we must be careful to include trends if we have trending dependent or explanatory variables. Since a time trend is exogenous, it can always serve as its own instrumental variable. The same is true of seasonal dummy variables, if monthly or quarterly data are used.

Series that have strong persistence (have unit roots) must be used with care, just as with OLS. Often, differencing the equation is warranted before estimation, and this applies to the instruments as well.

Under analogs of the assumptions in Chapter 11 for the asymptotic properties of OLS, 2SLS using time series data is consistent and asymptotically normally distributed. In fact, if we replace the explanatory variables with the instrumental variables in stating the assumptions, we only need to add the identification assumptions for 2SLS. For example, the homoskedasticity assumption is stated as

$$E(u_t^2 | z_{t1}, \dots, z_{tm}) = \sigma^2, \quad [15.53]$$

and the no serial correlation assumption is stated as

$$E(u_t u_s | \mathbf{z}_t, \mathbf{z}_s) = 0 \quad \text{for all } t \neq s, \quad [15.54]$$

where  $\mathbf{z}_t$  denotes all exogenous variables at time  $t$ . A full statement of the assumptions is given in the chapter appendix. We will provide examples of 2SLS for time series problems in Chapter 16; see also Computer Exercise C4.

As in the case of OLS, the no serial correlation assumption can often be violated with time series data. Fortunately, it is very easy to test for AR(1) serial correlation. If we write  $u_t = \rho u_{t-1} + e_t$  and plug this into equation (15.52), we get

$$y_t = \beta_0 + \beta_1 x_{t1} + \cdots + \beta_k x_{tk} + \rho u_{t-1} + e_t, \quad t \geq 2. \quad [15.55]$$

To test  $H_0: \rho = 0$ , we must replace  $u_{t-1}$  with the 2SLS residuals,  $\hat{u}_{t-1}$ . Further, if  $x_{ij}$  is endogenous in (15.52), then it is endogenous in (15.55), so we still need to use an IV. Because  $e_t$  is uncorrelated with all past values of  $u_t$ ,  $\hat{u}_{t-1}$  can be used as its own instrument.

### Testing for AR(1) Serial Correlation after 2SLS:

- (i) Estimate (15.52) by 2SLS and obtain the 2SLS residuals,  $\hat{u}_t$ .
- (ii) Estimate

$$y_t = \beta_0 + \beta_1 x_{t1} + \cdots + \beta_k x_{tk} + \rho \hat{u}_{t-1} + \text{error}_t, \quad t = 2, \dots, n$$

by 2SLS, using the same instruments from part (i), in addition to  $\hat{u}_{t-1}$ . Use the  $t$  statistic on  $\hat{\rho}$  to test  $H_0: \rho = 0$ .

As with the OLS version of this test from Chapter 12, the  $t$  statistic only has asymptotic justification, but it tends to work well in practice. A heteroskedasticity-robust version can be used to guard against heteroskedasticity. Further, lagged residuals can be added to the equation to test for higher forms of serial correlation using a joint  $F$  test.

What happens if we detect serial correlation? Some econometrics packages will compute standard errors that are robust to fairly general forms of serial correlation and heteroskedasticity. This is a nice, simple way to go if your econometrics package does this. The computations are very similar to those in Section 12-5 for OLS. [See Wooldridge (1995) for formulas and other computational methods.]

An alternative is to use the AR(1) model and correct for serial correlation. The procedure is similar to that for OLS and places additional restrictions on the instrumental variables. The quasi-differenced equation is the same as in equation (12.32):

$$\tilde{y}_t = \beta_0(1 - \rho) + \beta_1 \tilde{x}_{t1} + \cdots + \beta_k \tilde{x}_{tk} + e_t, \quad t \geq 2, \quad [15.56]$$

where  $\tilde{x}_{ij} = x_{ij} - \rho x_{i-1,j}$ . (We can use the  $t = 1$  observation just as in Section 12-3, but we omit that for simplicity here.) The question is: What can we use as instrumental variables? It seems natural to use the quasi-differenced instruments,  $\tilde{z}_{ij} = z_{ij} - \rho z_{i-1,j}$ . This only works, however, if in (15.52) the original error  $u_t$  is uncorrelated with the instruments at times  $t$ ,  $t - 1$ , and  $t + 1$ . That is, the instrumental variables must be strictly exogenous in (15.52). This rules out lagged dependent variables as IVs, for example. It also eliminates cases where future movements in the IVs react to current and past changes in the error,  $u_t$ .

### 2SLS with AR(1) Errors:

- (i) Estimate (15.52) by 2SLS and obtain the 2SLS residuals,  $\hat{u}_t$ ,  $t = 1, 2, \dots, n$ .
- (ii) Obtain  $\hat{\rho}$  from the regression of  $\hat{u}_t$  on  $\hat{u}_{t-1}$ ,  $t = 2, \dots, n$  and construct the quasi-differenced variables  $\tilde{y}_t = y_t - \hat{\rho} y_{t-1}$ ,  $\tilde{x}_{ij} = x_{ij} - \hat{\rho} x_{i-1,j}$ , and  $\tilde{z}_{ij} = z_{ij} - \hat{\rho} z_{i-1,j}$  for  $t \geq 2$ . (Remember, in most cases, some of the IVs will also be explanatory variables.)
- (iii) Estimate (15.56) (where  $\rho$  is replaced with  $\hat{\rho}$ ) by 2SLS, using the  $\tilde{z}_{ij}$  as the instruments. Assuming that (15.56) satisfies the 2SLS assumptions in the chapter appendix, the usual 2SLS test statistics are asymptotically valid.



We can also use the first time period as in Prais-Winsten estimation of the model with exogenous explanatory variables. The transformed variables in the first time period—the dependent variable, explanatory variables, and instrumental variables—are obtained simply by multiplying all first-period values by  $(1 - \hat{\rho})^{1/2}$ . (See also Section 12-3.)

## 15-8 Applying 2SLS to Pooled Cross Sections and Panel Data

Applying instrumental variables methods to independently pooled cross sections raises no new difficulties. As with models estimated by OLS, we should often include time period dummy variables to allow for aggregate time effects. These dummy variables are exogenous—because the passage of time is exogenous—and so they act as their own instruments.

### EXAMPLE 15.9 Effect of Education on Fertility

In Example 13.1, we used the pooled cross section in FERTIL1 to estimate the effect of education on women's fertility, controlling for various other factors. As in Sander (1992), we allow for the possibility that *educ* is endogenous in the equation. As instrumental variables for *educ*, we use mother's and father's education levels (*meduc*, *feduc*). The 2SLS estimate of  $\beta_{educ}$  is  $-.153$  ( $se = .039$ ), compared with the OLS estimate  $-.128$  ( $se = .018$ ). The 2SLS estimate shows a somewhat larger effect of education on fertility, but the 2SLS standard error is over twice as large as the OLS standard error. (In fact, the 95% confidence interval based on 2SLS easily contains the OLS estimate.) The OLS and 2SLS estimates of  $\beta_{educ}$  are *not statistically* different, as can be seen by testing for endogeneity of *educ* as in Section 15-5: when the reduced form residual,  $\hat{v}_2$ , is included with the other regressors in Table 13.1 (including *educ*), its *t* statistic is  $.702$ , which is not significant at any reasonable level. Therefore, in this case, we conclude that the difference between 2SLS and OLS could be entirely due to sampling error.

Instrumental variables estimation can be combined with panel data methods, particularly first differencing, to estimate parameters consistently in the presence of unobserved effects and endogeneity in one or more time-varying explanatory variables. The following simple example illustrates this combination of methods.

### EXAMPLE 15.10 Job Training and Worker Productivity

Suppose we want to estimate the effect of another hour of job training on worker productivity. For the two years 1987 and 1988, consider the simple panel data model

$$\log(scrap_{it}) = \beta_0 + \delta_0 d88_t + \beta_1 hrsemp_{it} + a_i + u_{it}, \quad t = 1, 2,$$

where *scrap<sub>it</sub>* is firm *i*'s scrap rate in year *t* and *hrsemp<sub>it</sub>* is hours of job training per employee. As usual, we allow different year intercepts and a constant, unobserved firm effect, *a<sub>i</sub>*.

For the reasons discussed in Section 13-2, we might be concerned that *hrsemp<sub>it</sub>* is correlated with *a<sub>i</sub>*, the latter of which contains unmeasured worker ability. As before, we difference to remove *a<sub>i</sub>*:

$$\Delta \log(scrap_i) = \delta_0 + \beta_1 \Delta hrsemp_i + \Delta u_i. \quad [15.57]$$

Normally, we would estimate this equation by OLS. But what if  $\Delta u_i$  is correlated with  $\Delta hrsemp_i$ ? For example, a firm might hire more skilled workers, while at the same time reducing the level of job training. In this case, we need an instrumental variable for  $\Delta hrsemp_i$ . Generally, such an IV would be hard to find, but we can exploit the fact that some firms received job training grants in 1988. If we assume that grant designation is uncorrelated with  $\Delta u_i$ —something that is reasonable, because the grants were given

at the beginning of 1988—then  $\Delta grant_i$  is valid as an IV, provided  $\Delta hrsemp$  and  $\Delta grant$  are correlated. Using the data in JTRAIN differenced between 1987 and 1988, the first stage regression is

$$\begin{aligned}\widehat{\Delta hrsemp} &= .51 + 27.88 \Delta grant \\ (1.56) \quad (3.13) \\ n &= 45, R^2 = .392.\end{aligned}$$

This confirms that the change in hours of job training per employee is strongly positively related to receiving a job training grant in 1988. In fact, receiving a job training grant increased per-employee training by almost 28 hours, and grant designation accounted for almost 40% of the variation in  $\Delta hrsemp$ . Two stage least squares estimation of (15.57) gives

$$\begin{aligned}\Delta \log(scrap) &= -.033 - .014 \Delta hrsemp \\ (.127) \quad (.008) \\ n &= 45, R^2 = .016.\end{aligned}$$

This means that 10 more hours of job training per worker are estimated to reduce the scrap rate by about 14%. For the firms in the sample, the average amount of job training in 1988 was about 17 hours per worker, with a minimum of zero and a maximum of 88.

For comparison, OLS estimation of (15.57) gives  $\hat{\beta}_1 = -.0076$  ( $se = .0045$ ), so the 2SLS estimate of  $\beta_1$  is almost twice as large in magnitude and is slightly more statistically significant.

When  $T \geq 3$ , the differenced equation may contain serial correlation. The same test and correction for AR(1) serial correlation from Section 15-7 can be used, where all regressions are pooled across  $i$  as well as  $t$ . Because we do not want to lose an entire time period, the Prais-Winsten transformation should be used for the initial time period.

Unobserved effects models containing lagged dependent variables also require IV methods for consistent estimation. The reason is that, after differencing,  $\Delta y_{i,t-1}$  is correlated with  $\Delta u_{it}$  because  $y_{i,t-1}$  and  $u_{i,t-1}$  are correlated. We can use two or more lags of  $y$  as IVs for  $\Delta y_{i,t-1}$ . [See Wooldridge (2010, Chapter 11) for details.]

Instrumental variables after differencing can be used on matched pairs samples as well. Ashenfelter and Krueger (1994) differenced the wage equation across twins to eliminate unobserved ability:

$$\log(wage_2) - \log(wage_1) = \delta_0 + \beta_1(educ_{2,2} - educ_{1,1}) + (u_2 - u_1),$$

where  $educ_{1,1}$  is years of schooling for the first twin as reported by the first twin and  $educ_{2,2}$  is years of schooling for the second twin as reported by the second twin. To account for possible measurement error in the self-reported schooling measures, Ashenfelter and Krueger used  $(educ_{2,1} - educ_{1,2})$  as an IV for  $(educ_{2,2} - educ_{1,1})$ , where  $educ_{2,1}$  is years of schooling for the second twin as reported by the first twin and  $educ_{1,2}$  is years of schooling for the first twin as reported by the second twin. The IV estimate of  $\beta_1$  is .167 ( $t = 3.88$ ), compared with the OLS estimate on the first differences of .092 ( $t = 3.83$ ) [see Ashenfelter and Krueger (1994, Table 3)].

## Summary

In Chapter 15, we have introduced the method of instrumental variables as a way to estimate the parameters in a linear model consistently when one or more explanatory variables are endogenous. An instrumental variable must have two properties: (1) it must be exogenous, that is, uncorrelated with the error term of the structural equation; (2) it must be partially correlated with the endogenous explanatory variable. Finding a variable with these two properties is usually challenging.

The method of two stage least squares, which allows for more instrumental variables than we have explanatory variables, is used routinely in the empirical social sciences. When used properly, it can allow us to estimate *ceteris paribus* effects in the presence of endogenous explanatory variables. This is true in cross-sectional, time series, and panel data applications. But when instruments are poor—which means they are correlated with the error term, only weakly correlated with the endogenous explanatory variable, or both—then 2SLS can be worse than OLS.

When we have valid instrumental variables, we can test whether an explanatory variable is endogenous, using the test in Section 15-5. In addition, though we can never test whether all IVs are exogenous, we can test that at least some of them are—assuming that we have more instruments than we need for consistent estimation (that is, the model is overidentified). Heteroskedasticity and serial correlation can be tested for and dealt with using methods similar to the case of models with exogenous explanatory variables.

In this chapter, we used omitted variables and measurement error to illustrate the method of instrumental variables. IV methods are also indispensable for simultaneous equations models, which we will cover in Chapter 16.

Key Terms

Endogenous Explanatory Variables	Instrument	Order Condition
Errors-in-Variables	Instrumental Variable	Overidentifying Restrictions
Exclusion Restrictions	Instrumental Variables (IV) Estimator	Rank Condition
Exogenous Explanatory Variables	Instrument Exogeneity	Reduced Form Equation
Exogenous Variables	Instrument Relevance	Structural Equation
First Stage	Natural Experiment	Two Stage Least Squares (2SLS) Estimator
Identification	Omitted Variables	Weak Instruments

Problems

- 1 Consider a simple model to estimate the effect of personal computer (PC) ownership on college grade point average for graduating seniors at a large public university:
- $$GPA = \beta_0 + \beta_1 PC + u,$$
- where *PC* is a binary variable indicating PC ownership.
- (i) Why might PC ownership be correlated with *u*?

(ii) Explain why *PC* is likely to be related to parents’ annual income. Does this mean parental income is a good IV for *PC*? Why or why not?

(iii) Suppose that, four years ago, the university gave grants to buy computers to roughly one-half of the incoming students, and the students who received grants were randomly chosen. Carefully explain how you would use this information to construct an instrumental variable for *PC*.
- 2 Suppose that you wish to estimate the effect of class attendance on student performance, as in Example 6.3. A basic model is
- $$stdfnl = \beta_0 + \beta_1 atndrte + \beta_2 priGPA + \beta_3 ACT + u,$$
- where the variables are defined as in Chapter 6.
- (i) Let *dist* be the distance from the students’ living quarters to the lecture hall. Do you think *dist* is uncorrelated with *u*?

(ii) Assuming that *dist* and *u* are uncorrelated, what other assumption must *dist* satisfy to be a valid IV for *atndrte*?

(iii) Suppose, as in equation (6.18), we add the interaction term  $priGPA \cdot atndrte$ :

$$stndfnl = \beta_0 + \beta_1 atndrte + \beta_2 priGPA + \beta_3 ACT + \beta_4 priGPA \cdot atndrte + u.$$

If  $atndrte$  is correlated with  $u$ , then, in general, so is  $priGPA \cdot atndrte$ . What might be a good IV for  $priGPA \cdot atndrte$ ? [Hint: If  $E(u|priGPA, ACT, dist) = 0$ , as happens when  $priGPA$ ,  $ACT$ , and  $dist$  are all exogenous, then any function of  $priGPA$  and  $dist$  is uncorrelated with  $u$ .]

3 Consider the simple regression model

$$y = \beta_0 + \beta_1 x + u$$

and let  $z$  be a *binary* instrumental variable for  $x$ . Use (15.10) to show that the IV estimator  $\hat{\beta}_1$  can be written as

$$\hat{\beta}_1 = (\bar{y}_1 - \bar{y}_0) / (\bar{x}_1 - \bar{x}_0),$$

where  $\bar{y}_0$  and  $\bar{x}_0$  are the sample averages of  $y_i$  and  $x_i$  over the part of the sample with  $z_i = 0$ , and where  $\bar{y}_1$  and  $\bar{x}_1$  are the sample averages of  $y_i$  and  $x_i$  over the part of the sample with  $z_i = 1$ . This estimator, known as a *grouping estimator*, was first suggested by Wald (1940).

4 Suppose that, for a given state in the United States, you wish to use annual time series data to estimate the effect of the state-level minimum wage on the employment of those 18 to 25 years old ( $EMP$ ). A simple model is

$$gEMP_t = \beta_0 + \beta_1 gMIN_t + \beta_2 gPOP_t + \beta_3 gGSP_t + \beta_4 gGDP_t + u_t,$$

where  $MIN_t$  is the minimum wage, in real dollars;  $POP_t$  is the population from 18 to 25 years old;  $GSP_t$  is gross state product; and  $GDP_t$  is U.S. gross domestic product. The  $g$  prefix indicates the growth rate from year  $t - 1$  to year  $t$ , which would typically be approximated by the difference in the logs.

- (i) If we are worried that the state chooses its minimum wage partly based on unobserved (to us) factors that affect youth employment, what is the problem with OLS estimation?
- (ii) Let  $USMIN_t$  be the U.S. minimum wage, which is also measured in real terms. Do you think  $gUSMIN_t$  is uncorrelated with  $u_t$ ?
- (iii) By law, any state's minimum wage must be at least as large as the U.S. minimum. Explain why this makes  $gUSMIN_t$  a potential IV candidate for  $gMIN_t$ .

5 Refer to equations (15.19) and (15.20). Assume that  $\sigma_u = \sigma_x$ , so that the population variation in the error term is the same as it is in  $x$ . Suppose that the instrumental variable,  $z$ , is slightly correlated with  $u$ :  $\text{Corr}(z, u) = .1$ . Suppose also that  $z$  and  $x$  have a somewhat stronger correlation:  $\text{Corr}(z, x) = .2$ .

- (i) What is the asymptotic bias in the IV estimator?
- (ii) How much correlation would have to exist between  $x$  and  $u$  before OLS has more asymptotic bias than 2SLS?

6 (i) In the model with one endogenous explanatory variable, one exogenous explanatory variable, and one extra exogenous variable, take the reduced form for  $y_2$  (15.26), and plug it into the structural equation (15.22). This gives the reduced form for  $y_1$ :

$$y_1 = \alpha_0 + \alpha_1 z_1 + \alpha_2 z_2 + v_1.$$

Find the  $\alpha_j$  in terms of the  $\beta_j$  and the  $\pi_j$ .

- (ii) Find the reduced form error,  $v_1$ , in terms of  $u_1$ ,  $v_2$ , and the parameters.
- (iii) How would you consistently estimate the  $\alpha_j$ ?

7 The following is a simple model to measure the effect of a school choice program on standardized test performance [see Rouse (1998) for motivation and Computer Exercise C11 for an analysis of a subset of Rouse's data]:

$$score = \beta_0 + \beta_1 choice + \beta_2 faminc + u_1,$$

where *score* is the score on a statewide test, *choice* is a binary variable indicating whether a student attended a choice school in the last year, and *faminc* is family income. The IV for *choice* is *grant*, the dollar amount granted to students to use for tuition at choice schools. The grant amount differed by family income level, which is why we control for *faminc* in the equation.

- (i) Even with *faminc* in the equation, why might *choice* be correlated with  $u_1$ ?
  - (ii) If within each income class, the grant amounts were assigned randomly, is *grant* uncorrelated with  $u_1$ ?
  - (iii) Write the reduced form equation for *choice*. What is needed for *grant* to be partially correlated with *choice*?
  - (iv) Write the reduced form equation for *score*. Explain why this is useful. (*Hint*: How do you interpret the coefficient on *grant*?)
- 8 Suppose you want to test whether girls who attend a girls' high school do better in math than girls who attend coed schools. You have a random sample of senior high school girls from a state in the United States, and *score* is the score on a standardized math test. Let *girlshs* be a dummy variable indicating whether a student attends a girls' high school.
- (i) What other factors would you control for in the equation? (You should be able to reasonably collect data on these factors.)
  - (ii) Write an equation relating *score* to *girlshs* and the other factors you listed in part (i).
  - (iii) Suppose that parental support and motivation are unmeasured factors in the error term in part (ii). Are these likely to be correlated with *girlshs*? Explain.
  - (iv) Discuss the assumptions needed for the number of girls' high schools within a 20-mile radius of a girl's home to be a valid IV for *girlshs*.
  - (v) Suppose that, when you estimate the reduced form for *girlshs*, you find that the coefficient on *numghs* (the number of girls' high schools within a 20-mile radius) is negative and statistically significant. Would you feel comfortable proceeding with IV estimation where *numghs* is used as an IV for *girlshs*? Explain.
- 9 Suppose that, in equation (15.8), you do not have a good instrumental variable candidate for *skipped*. But you have two other pieces of information on students: combined SAT score and cumulative GPA prior to the semester. What would you do instead of IV estimation?
- 10 In a recent article, Evans and Schwab (1995) studied the effects of attending a Catholic high school on the probability of attending college. For concreteness, let *college* be a binary variable equal to unity if a student attends college, and zero otherwise. Let *CathHS* be a binary variable equal to one if the student attends a Catholic high school. A linear probability model is

$$\text{college} = \beta_0 + \beta_1 \text{CathHS} + \text{other factors} + u,$$

where the other factors include gender, race, family income, and parental education.

- (i) Why might *CathHS* be correlated with  $u$ ?
  - (ii) Evans and Schwab have data on a standardized test score taken when each student was a sophomore. What can be done with this variable to improve the ceteris paribus estimate of attending a Catholic high school?
  - (iii) Let *CathRel* be a binary variable equal to one if the student is Catholic. Discuss the two requirements needed for this to be a valid IV for *CathHS* in the preceding equation. Which of these can be tested?
  - (iv) Not surprisingly, being Catholic has a significant positive effect on attending a Catholic high school. Do you think *CathRel* is a convincing instrument for *CathHS*?
- 11 Consider a simple time series model where the explanatory variable has classical measurement error:

$$\begin{aligned} y_t &= \beta_0 + \beta_1 x_t^* + u_t \\ x_t &= x_t^* + e_t, \end{aligned} \quad [15.58]$$

where  $u_t$  has zero mean and is uncorrelated with  $x_t^*$  and  $e_t$ . We observe  $y_t$  and  $x_t$  only. Assume that  $e_t$  has zero mean and is uncorrelated with  $x_t^*$  and that  $x_t^*$  also has a zero mean (this last assumption is only to simplify the algebra).

- (i) Write  $x_t^* = x_t - e_t$  and plug this into (15.58). Show that the error term in the new equation, say,  $v_t$ , is negatively correlated with  $x_t$  if  $\beta_1 > 0$ . What does this imply about the OLS estimator of  $\beta_1$  from the regression of  $y_t$  on  $x_t$ ?
- (ii) In addition to the previous assumptions, assume that  $u_t$  and  $e_t$  are uncorrelated with all past values of  $x_t^*$  and  $e_t$ ; in particular, with  $x_{t-1}^*$  and  $e_{t-1}$ . Show that  $E(x_{t-1}v_t) = 0$  where  $v_t$  is the error term in the model from part (i).
- (iii) Are  $x_t$  and  $x_{t-1}$  likely to be correlated? Explain.
- (iv) What do parts (ii) and (iii) suggest as a useful strategy for consistently estimating  $\beta_0$  and  $\beta_1$ ?

## Computer Exercises

**C1** Use the data in WAGE2 for this exercise.

- (i) In Example 15.2, if *sibs* is used as an instrument for *educ*, the IV estimate of the return to education is .122. To convince yourself that using *sibs* as an IV for *educ* is *not* the same as just plugging *sibs* in for *educ* and running an OLS regression, run the regression of  $\log(\text{wage})$  on *sibs* and explain your findings.
- (ii) The variable *brthord* is birth order (*brthord* is one for a first-born child, two for a second-born child, and so on). Explain why *educ* and *brthord* might be negatively correlated. Regress *educ* on *brthord* to determine whether there is a statistically significant negative correlation.
- (iii) Use *brthord* as an IV for *educ* in equation (15.1). Report and interpret the results.
- (iv) Now, suppose that we include number of siblings as an explanatory variable in the wage equation; this controls for family background, to some extent:

$$\log(\text{wage}) = \beta_0 + \beta_1 \text{educ} + \beta_2 \text{sibs} + u.$$

Suppose that we want to use *brthord* as an IV for *educ*, assuming that *sibs* is exogenous. The reduced form for *educ* is

$$\text{educ} = \pi_0 + \pi_1 \text{sibs} + \pi_2 \text{brthord} + v.$$

State and test the identification assumption.

- (v) Estimate the equation from part (iv) using *brthord* as an IV for *educ* (and *sibs* as its own IV). Comment on the standard errors for  $\hat{\beta}_{\text{educ}}$  and  $\hat{\beta}_{\text{sibs}}$ .
- (vi) Using the fitted values from part (iv),  $\widehat{\text{educ}}$ , compute the correlation between  $\widehat{\text{educ}}$  and *sibs*. Use this result to explain your findings from part (v).

**C2** The data in FERTIL2 include, for women in Botswana during 1988, information on number of children, years of education, age, and religious and economic status variables.

- (i) Estimate the model

$$\text{children} = \beta_0 + \beta_1 \text{educ} + \beta_2 \text{age} + \beta_3 \text{age}^2 + u$$

by OLS and interpret the estimates. In particular, holding *age* fixed, what is the estimated effect of another year of education on fertility? If 100 women receive another year of education, how many fewer children are they expected to have?

- (ii) The variable *frsthalf* is a dummy variable equal to one if the woman was born during the first six months of the year. Assuming that *frsthalf* is uncorrelated with the error term from part (i), show that *frsthalf* is a reasonable IV candidate for *educ*. (Hint: You need to do a regression.)
- (iii) Estimate the model from part (i) by using *frsthalf* as an IV for *educ*. Compare the estimated effect of education with the OLS estimate from part (i).



- (iv) Add the binary variables *electric*, *tv*, and *bicycle* to the model and assume these are exogenous. Estimate the equation by OLS and 2SLS and compare the estimated coefficients on *educ*. Interpret the coefficient on *tv* and explain why television ownership has a negative effect on fertility.

**C3** Use the data in CARD for this exercise.

- (i) The equation we estimated in Example 15.4 can be written as

$$\log(\text{wage}) = \beta_0 + \beta_1 \text{educ} + \beta_2 \text{exper} + \cdots + u,$$

where the other explanatory variables are listed in Table 15.1. In order for IV to be consistent, the IV for *educ*, *nearc4*, must be uncorrelated with *u*. Could *nearc4* be correlated with things in the error term, such as unobserved ability? Explain.

- (ii) For a subsample of the men in the data set, an IQ score is available. Regress *IQ* on *nearc4* to check whether average IQ scores vary by whether the man grew up near a four-year college. What do you conclude?
- (iii) Now, regress *IQ* on *nearc4*, *smsa66*, and the 1966 regional dummy variables *reg662*, ..., *reg669*. Are *IQ* and *nearc4* related after the geographic dummy variables have been partialled out? Reconcile this with your findings from part (ii).
- (iv) From parts (ii) and (iii), what do you conclude about the importance of controlling for *smsa66* and the 1966 regional dummies in the  $\log(\text{wage})$  equation?

**C4** Use the data in INTDEF for this exercise. A simple equation relating the three-month T-bill rate to the inflation rate (constructed from the Consumer Price Index) is

$$i3_t = \beta_0 + \beta_1 \text{inf}_t + u_t.$$

- (i) Estimate this equation by OLS, omitting the first time period for later comparisons. Report the results in the usual form.
- (ii) Some economists feel that the Consumer Price Index mismeasures the true rate of inflation, so that the OLS from part (i) suffers from measurement error bias. Reestimate the equation from part (i), using  $\text{inf}_{t-1}$  as an IV for  $\text{inf}_t$ . How does the IV estimate of  $\beta_1$  compare with the OLS estimate?
- (iii) Now, first difference the equation:

$$\Delta i3_t = \beta_0 + \beta_1 \Delta \text{inf}_t + \Delta u_t.$$

Estimate this by OLS and compare the estimate of  $\beta_1$  with the previous estimates.

- (iv) Can you use  $\Delta \text{inf}_{t-1}$  as an IV for  $\Delta \text{inf}_t$  in the differenced equation in part (iii)? Explain. (Hint: Are  $\Delta \text{inf}_t$  and  $\Delta \text{inf}_{t-1}$  sufficiently correlated?)

**C5** Use the data in CARD for this exercise.

- (i) In Table 15.1, the difference between the IV and OLS estimates of the return to education is economically important. Obtain the reduced form residuals,  $\hat{v}_2$ , from the reduced form regression *educ* on *nearc4*, *exper*, *exper*<sup>2</sup>, *black*, *smsa*, *south*, *smsa66*, *reg662*, ..., *reg669*—see Table 15.1. Use these to test whether *educ* is exogenous; that is, determine if the difference between OLS and IV is statistically significant.
- (ii) Estimate the equation by 2SLS, adding *nearc2* as an instrument. Does the coefficient on *educ* change much?
- (iii) Test the single overidentifying restriction from part (ii).

**C6** Use the data in MURDER for this exercise. The variable *mrdrte* is the murder rate, that is, the number of murders per 100,000 people. The variable *exec* is the total number of prisoners executed for the current and prior two years; *unem* is the state unemployment rate.

- (i) How many states executed at least one prisoner in 1991, 1992, or 1993? Which state had the most executions?
- (ii) Using the two years 1990 and 1993, do a pooled regression of *mrdrte* on *d93*, *exec*, and *unem*. What do you make of the coefficient on *exec*?

- (iii) Using the changes from 1990 to 1993 only (for a total of 51 observations), estimate the equation

$$\Delta mrd rte = \delta_0 + \beta_1 \Delta exec + \beta_2 \Delta unem + \Delta u$$

by OLS and report the results in the usual form. Now, does capital punishment appear to have a deterrent effect?

- (iv) The change in executions may be at least partly related to changes in the expected murder rate, so that  $\Delta exec$  is correlated with  $\Delta u$  in part (iii). It might be reasonable to assume that  $\Delta exec_{-1}$  is uncorrelated with  $\Delta u$ . (After all,  $\Delta exec_{-1}$  depends on executions that occurred three or more years ago.) Regress  $\Delta exec$  on  $\Delta exec_{-1}$  to see if they are sufficiently correlated; interpret the coefficient on  $\Delta exec_{-1}$ .
- (v) Reestimate the equation from part (iii), using  $\Delta exec_{-1}$  as an IV for  $\Delta exec$ . Assume that  $\Delta unem$  is exogenous. How do your conclusions change from part (iii)?

**C7** Use the data in PHILLIPS for this exercise.

- (i) In Example 11.5, we estimated an expectations augmented Phillips curve of the form

$$\Delta inf_t = \beta_0 + \beta_1 unem_t + e_t,$$

where  $\Delta inf_t = inf_t - inf_{t-1}$ . In estimating this equation by OLS, we assumed that the supply shock,  $e_t$ , was uncorrelated with  $unem_t$ . If this is false, what can be said about the OLS estimator of  $\beta_1$ ?

- (ii) Suppose that  $e_t$  is unpredictable given all past information:  $E(e_t | inf_{t-1}, unem_{t-1}, \dots) = 0$ . Explain why this makes  $unem_{t-1}$  a good IV candidate for  $unem_t$ .
- (iii) Regress  $unem_t$  on  $unem_{t-1}$ . Are  $unem_t$  and  $unem_{t-1}$  significantly correlated?
- (iv) Estimate the expectations augmented Phillips curve by IV. Report the results in the usual form and compare them with the OLS estimates from Example 11.5.

**C8** Use the data in 401KSUBS for this exercise. The equation of interest is a linear probability model:

$$pira = \beta_0 + \beta_1 p401k + \beta_2 inc + \beta_3 inc^2 + \beta_4 age + \beta_5 age^2 + u.$$

The goal is to test whether there is a tradeoff between participating in a 401(k) plan and having an individual retirement account (IRA). Therefore, we want to estimate  $\beta_1$ .

- (i) Estimate the equation by OLS and discuss the estimated effect of  $p401k$ .
- (ii) For the purposes of estimating the ceteris paribus tradeoff between participation in two different types of retirement savings plans, what might be a problem with ordinary least squares?
- (iii) The variable  $e401k$  is a binary variable equal to one if a worker is *eligible* to participate in a 401(k) plan. Explain what is required for  $e401k$  to be a valid IV for  $p401k$ . Do these assumptions seem reasonable?
- (iv) Estimate the reduced form for  $p401k$  and verify that  $e401k$  has significant partial correlation with  $p401k$ . Since the reduced form is also a linear probability model, use a heteroskedasticity-robust standard error.
- (v) Now, estimate the structural equation by IV and compare the estimate of  $\beta_1$  with the OLS estimate. Again, you should obtain heteroskedasticity-robust standard errors.
- (vi) Test the null hypothesis that  $p401k$  is in fact exogenous, using a heteroskedasticity-robust test.

**C9** The purpose of this exercise is to compare the estimates and standard errors obtained by correctly using 2SLS with those obtained using inappropriate procedures. Use the data file WAGE2.

- (i) Use a 2SLS routine to estimate the equation

$$\log(wage) = \beta_0 + \beta_1 educ + \beta_2 exper + \beta_3 tenure + \beta_4 black + u,$$

where  $sibs$  is the IV for  $educ$ . Report the results in the usual form.

- (ii) Now, manually carry out 2SLS. That is, first regress  $educ_i$  on  $sibs_i$ ,  $exper_i$ ,  $tenure_i$ , and  $black_i$  and obtain the fitted values,  $\widehat{educ}_i$ ,  $i = 1, \dots, n$ . Then, run the second stage regression  $\log(wage_i)$  on  $\widehat{educ}_i$ ,  $exper_i$ ,  $tenure_i$ , and  $black_i$ ,  $i = 1, \dots, n$ . Verify that the  $\hat{\beta}_j$  are identical to those obtained

from part (i), but that the standard errors are somewhat different. The standard errors obtained from the second stage regression when manually carrying out 2SLS are generally inappropriate.

- (iii) Now, use the following two-step procedure, which generally yields inconsistent parameter estimates of the  $\beta_j$ , and not just inconsistent standard errors. In step one, regress  $educ_i$  on  $sibs_i$  only and obtain the fitted values, say  $\widehat{educ}_i$ . (Note that this is an incorrect first stage regression.) Then, in the second step, run the regression of  $\log(wage_i)$  on  $\widehat{educ}_i$ ,  $exper_i$ ,  $tenure_i$ , and  $black_i$ ,  $i = 1, \dots, n$ . How does the estimate from this incorrect, two-step procedure compare with the correct 2SLS estimate of the return to education?

**C10** Use the data in HTV for this exercise.

- (i) Run a simple OLS regression of  $\log(wage)$  on  $educ$ . Without controlling for other factors, what is the 95% confidence interval for the return to another year of education?
- (ii) The variable  $ctuit$ , in thousands of dollars, is the change in college tuition facing students from age 17 to age 18. Show that  $educ$  and  $ctuit$  are essentially uncorrelated. What does this say about  $ctuit$  as a possible IV for  $educ$  in a simple regression analysis?
- (iii) Now, add to the simple regression model in part (i) a quadratic in experience and a full set of regional dummy variables for current residence and residence at age 18. Also include the urban indicators for current and age 18 residences. What is the estimated return to a year of education?
- (iv) Again using  $ctuit$  as a potential IV for  $educ$ , estimate the reduced form for  $educ$ . [Naturally, the reduced form for  $educ$  now includes the explanatory variables in part (iii).] Show that  $ctuit$  is now statistically significant in the reduced form for  $educ$ .
- (v) Estimate the model from part (iii) by IV, using  $ctuit$  as an IV for  $educ$ . How does the confidence interval for the return to education compare with the OLS CI from part (iii)?
- (vi) Do you think the IV procedure from part (v) is convincing?

**C11** The data set in VOUCHER, which is a subset of the data used in Rouse (1998), can be used to estimate the effect of school choice on academic achievement. Attendance at a choice school was paid for by a voucher, which was determined by a lottery among those who applied. The data subset was chosen so that any student in the sample has a valid 1994 math test score (the last year available in Rouse's sample). Unfortunately, as pointed out by Rouse, many students have missing test scores, possibly due to attrition (that is, leaving the Milwaukee public school district). These data include students who applied to the voucher program and were accepted, students who applied and were not accepted, and students who did not apply. Therefore, even though the vouchers were chosen by lottery among those who applied, we do not necessarily have a random sample from a population where being selected for a voucher has been randomly determined. (An important consideration is that students who never applied to the program may be systematically different from those who did—and in ways that we cannot know based on the data.)

Rouse (1998) uses panel data methods of the kind we discussed in Chapter 14 to allow student fixed effects; she also uses instrumental variables methods. This problem asks you to do a cross-sectional analysis which winning the lottery for a voucher acts as an instrumental variable for attending a choice school. Actually, because we have multiple years of data on each student, we construct two variables. The first, *choicelyrs*, is the number of years from 1991 to 1994 that a student attended a choice school; this variable ranges from zero to four. The variable *selectyrs* indicates the number of years a student was selected for a voucher. If the student applied for the program in 1990 and received a voucher then *selectyrs* = 4; if he or she applied in 1991 and received a voucher then *selectyrs* = 3; and so on. The outcome of interest is *mnce*, the student's percentile score on a math test administered in 1994.

- (i) Of the 990 students in the sample, how many were never awarded a voucher? How many had a voucher available for four years? How many students actually attended a choice school for four years?
- (ii) Run a simple regression of *choicelyrs* on *selectyrs*. Are these variables related in the direction you expected? How strong is the relationship? Is *selectyrs* a sensible IV candidate for *choicelyrs*?
- (iii) Run a simple regression of *mnce* on *choicelyrs*. What do you find? Is this what you expected? What happens if you add the variables *black*, *hispanic*, and *female*?

- (iv) Why might *choiceyrs* be endogenous in an equation such as

$$mnce = \beta_0 + \beta_1 \text{choiceyrs} + \beta_2 \text{black} + \beta_3 \text{hispanic} + \beta_4 \text{female} + u_1$$

- (v) Estimate the equation in part (iv) by instrumental variables, using *selectyrs* as the IV for *choiceyrs*. Does using IV produce a positive effect of attending a choice school? What do you make of the coefficients on the other explanatory variables?
- (vi) To control for the possibility that prior achievement affects participating in the lottery (as well as predicting attrition), add *mnce90*—the math score in 1990—to the equation in part (iv). Estimate the equation by OLS and IV, and compare the results for  $\beta_1$ . For the IV estimate, how much is each year in a choice school worth on the math percentile score? Is this a practically large effect?
- (vii) Why is the analysis from part (vi) not entirely convincing? [Hint: Compared with part (v), what happens to the number of observations, and why?]
- (viii) The variables *choiceyrs1*, *choiceyrs2*, and so on are dummy variables indicating the different number of years a student could have been in a choice school (from 1991 to 1994). The dummy variables *selectyrs1*, *selectyrs2*, and so on have a similar definition, but for being selected from the lottery. Estimate the equation

$$\begin{aligned} mnce = & \beta_0 + \beta_1 \text{choiceyrs1} + \beta_2 \text{choiceyrs2} + \beta_3 \text{choiceyrs3} + \beta_4 \text{choiceyrs4} \\ & + \beta_5 \text{black} + \beta_6 \text{hispanic} + \beta_7 \text{female} + \beta_8 \text{mnce90} + u_1 \end{aligned}$$

by IV, using as instruments the four *selectyrs* dummy variables. (As before, the variables *black*, *hispanic*, and *female* act as their own IVs.) Describe your findings. Do they make sense?

- C12** Use the data in CATHOLIC to answer this question. The model of interest is

$$\text{math12} = \beta_0 + \beta_1 \text{cathhs} + \beta_2 \text{lfaminc} + \beta_3 \text{motheduc} + \beta_4 \text{fatheduc} + u,$$

where *cathhs* is a binary indicator for whether a student attends a Catholic high school.

- (i) How many students are in the sample? What percentage of these students attend a Catholic high school?
- (ii) Estimate the above equation by OLS. What is the estimate of  $\beta_1$ ? What is its 95% confidence interval?
- (iii) Using *parcath* as an instrument for *cathhs*, estimate the reduced form for *cathhs*. What is the *t* statistic for *parcath*? Is there evidence of a weak instrument problem?
- (iv) Estimate the above equation by IV, using *parcath* as an IV for *cathhs*. How does the estimate and 95% CI compare with the OLS quantities?
- (v) Test the null hypothesis that *cathhs* is exogenous. What is the *p*-value of the test?
- (vi) Suppose you add the interaction between *cathhs* • *motheduc* to the above model. Why is it generally endogenous? Why is *pareduc* • *motheduc* a good IV candidate for *cathhs* • *motheduc*?
- (vii) Before you create the interactions in part (vi), first find the sample average of *motheduc* and create *cathhs* • (*motheduc* − *motheduc*) and *parcath* • (*motheduc* − *motheduc*). Add the first interaction to the model and use the second as an IV. Of course, *cathhs* is also instrumented. Is the interaction term statistically significant?
- (viii) Compare the coefficient on *cathhs* in (vii) to that in part (iv). Is including the interaction important for estimating the average partial effect?

- C13** Use the data in LABSUP to answer the following questions. These are data on almost 32,000 black or Hispanic women. Every woman in the sample is married. It is a subset of the data used in Angrist and Evans (1998). Our interest here is in determining how weekly hours worked, *hours*, changes with number of children (*kids*). All women in the sample have at least two children. The two potential

instrumental variables for *kids*, which is suspected as being endogenous, work to generate exogenous variation starting with two children. See the original article for further discussion.

- (i) Estimate the equation

$$\text{hours} = \beta_0 + \beta_1 \text{kids} + \beta_2 \text{nonmomi} + \beta_3 \text{educ} + \beta_4 \text{age} + \beta_5 \text{age}^2 + \beta_6 \text{black} + \beta_7 \text{hispan} + u$$

by OLS and obtain the heteroskedasticity-robust standard errors. Interpret the coefficient on *kids*. Discuss its statistical significance.

- (ii) A variable that Angrist and Evans propose as an instrument is *samesex*, a binary variable equal to one if the first two children are the same biological sex. What do you think is the argument for why it is a relevant instrument for *kids*?
- (iii) Run the regression

$$\text{kids}_i \text{ on } \text{samesex}_i, \text{nonmomi}_i, \text{educ}_i, \text{age}_i, \text{age}_i^2, \text{black}_i, \text{hispan}_i$$

and see if the story from part (ii) holds up. In particular, interpret the coefficient on *samesex*. How statistically significant is *samesex*?

- (iv) Can you think of mechanisms by which *samesex* is correlated with *u* in the equation in part (i)? (It is fine to assume that biological sex is randomly determined.) [Hint: How might a family's finances be affected based on whether they have two children of the same sex or two children of opposite sex?]
- (v) Is it legitimate to check for exogeneity of *samesex* by adding it to the regression in part (i) and testing its significance? Explain.
- (vi) Using *samesex* as an IV for *kids*, obtain the IV estimates of the equation in part (i). How does the *kids* coefficient compare with the OLS estimate? Is the IV estimate precise?
- (vii) Now add *multi2nd* as an instrument. Obtain the *F* statistic from the first stage regression and determining whether *samesex* and *multi2nd* are sufficiently strong.
- (viii) Using *samesex* and *multi2nd* both as instruments for *kids*, how does the 2SLS estimate compare with the OLS and IV estimates from the previous parts?
- (ix) Using the estimation from part (viii), is there strong evidence that *kids* is endogenous in the *hours* equation?
- (x) In part (viii), how many overidentification restrictions are there? Does the overidentification test pass?

## APPENDIX 15A

### 15A.1 Assumptions for Two Stage Least Squares

This appendix covers the assumptions under which 2SLS has desirable large sample properties. We first state the assumptions for cross-sectional applications under random sampling. Then, we discuss what needs to be added for them to apply to time series and panel data.

### 15A.2 Assumption 2SLS.1 (Linear in Parameters)

The model in the population can be written as

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_k x_k + u,$$

where  $\beta_0, \beta_1, \dots, \beta_k$  are the unknown parameters (constants) of interest and *u* is an unobserved random error or random disturbance term. The instrumental variables are denoted as  $z_j$ .

It is worth emphasizing that Assumption 2SLS.1 is virtually identical to MLR.1 (with the minor exception that 2SLS.1 mentions the notation for the instrumental variables,  $z_j$ ). In other words, the model we are interested in is the same as that for OLS estimation of the  $\beta_j$ . Sometimes it is easy to lose sight of the fact that we can apply different estimation methods to the same model. Unfortunately, it is not uncommon to hear researchers say “I estimated an OLS model” or “I used a 2SLS model.” Such statements are meaningless. OLS and 2SLS are different *estimation* methods that are applied to the *same* model. It is true that they have desirable statistical properties under different sets of assumptions on the model, but the relationship they are estimating is given by the equation in 2SLS.1 (or MLR.1). The point is similar to that made for the unobserved effects panel data model covered in Chapters 13 and 14: pooled OLS, first differencing, fixed effects, and random effects are different estimation methods for the same model.

### 15A.3 Assumption 2SLS.2 (Random Sampling)

We have a random sample on  $y$ , the  $x_j$ , and the  $z_j$ .

### 15A.4 Assumption 2SLS.3 (Rank Condition)

(i) There are no perfect linear relationships among the instrumental variables. (ii) The rank condition for identification holds.

With a single endogenous explanatory variable, as in equation (15.42), the rank condition is easily described. Let  $z_1, \dots, z_m$  denote the exogenous variables, where  $z_k, \dots, z_m$  do not appear in the structural model (15.42). The reduced form of  $y_2$  is

$$y_2 = \pi_0 + \pi_1 z_1 + \pi_2 z_2 + \dots + \pi_{k-1} z_{k-1} + \pi_k z_k + \dots + \pi_m z_m + v_2.$$

Then, we need at least one of  $\pi_k, \dots, \pi_m$  to be nonzero. This requires at least one exogenous variable that does not appear in (15.42) (the order condition). Stating the rank condition with two or more endogenous explanatory variables requires matrix algebra. [See Wooldridge (2010, Chapter 5).]

### 15A.5 Assumption 2SLS.4 (Exogenous Instrumental Variables)

The error term  $u$  has zero mean, and each IV is uncorrelated with  $u$ . Remember that any  $x_j$  that is uncorrelated with  $u$  also acts as an IV.

### 15A.6 Theorem 15A.1

Under Assumptions 2SLS.1 through 2SLS.4, the 2SLS estimator is consistent.

### 15A.7 Assumption 2SLS.5 (Homoskedasticity)

Let  $\mathbf{z}$  denote the collection of all instrumental variables. Then,  $E(u^2|\mathbf{z}) = \sigma^2$ .

### 15A.8 Theorem 15A.2

Under Assumptions 2SLS.1 through 2SLS.5, the 2SLS estimators are asymptotically normally distributed. Consistent estimators of the asymptotic variance are given as in equation (15.43), where  $\sigma^2$  is replaced with  $\hat{\sigma}^2 = (n - k - 1)^{-1} \sum_{i=1}^n \hat{u}_i^2$ , and the  $\hat{u}_i$  are the 2SLS residuals.

The 2SLS estimator is also the best IV estimator under the five assumptions given. We state the result here. A proof can be found in Wooldridge (2010, Chapter 5).



**15A.9 Theorem 15A.3**

Under Assumptions 2SLS.1 through 2SLS.5, the 2SLS estimator is asymptotically efficient in the class of IV estimators that uses linear combinations of the exogenous variables as instruments.

If the homoskedasticity assumption does not hold, the 2SLS estimators are still asymptotically normal, but the standard errors (and  $t$  and  $F$  statistics) need to be adjusted; many econometrics packages do this routinely. Moreover, the 2SLS estimator is no longer the asymptotically efficient IV estimator, in general. We will not study more efficient estimators here [see Wooldridge (2010, Chapter 8)].

For time series applications, we must add some assumptions. First, as with OLS, we must assume that all series (including the IVs) are weakly dependent: this ensures that the law of large numbers and the central limit theorem hold. For the usual standard errors and test statistics to be valid, as well as for asymptotic efficiency, we must add a no serial correlation assumption.

**15A.10 Assumption 2SLS.6 (No Serial Correlation)**

Equation (15.54) holds.

A similar no serial correlation assumption is needed in panel data applications. Tests and corrections for serial correlation were discussed in Section 15-7.