

## Section 11 Endogenous Regressors and Instrumental Variables

*If  $x$  is correlated with  $u$*

- If  $\text{cov}(x, u) \neq 0$ , then OLS is biased and inconsistent
  - Coefficient on  $x$  will pick up the effects of the parts of  $u$  that are correlated with it in addition to the direct effects of  $x$
  - Direction of bias depends on sign of correlation between  $x$  and  $u$
  - This is exactly analogous to omitted-variables bias
- **Measurement error** (discussed above under internal validity)
  - Suppose that the dependent variable is measured accurately but that we measure  $x$  with error:  $\tilde{x}_i = x_i + \eta_i$ .
  - The estimated model is  $y_i = \beta_0 + \beta_1 \tilde{x}_i + (u_i - \beta_1 \eta_i)$ .
  - Because  $\eta$  is part of  $\tilde{x}$  and therefore correlated with it, the composite error term is now correlated with the actual regressor, meaning that the OLS slope estimator is biased and inconsistent.
    - If  $u$  and  $\eta$  are independent and normal, then  $\text{plim } \hat{\beta}_1 = \frac{\sigma_x^2}{\sigma_x^2 + \sigma_\eta^2} \beta_1$ .
    - The estimator is biased toward zero.
    - If most of the variation in  $\tilde{x}$  comes from  $x$ , then the bias will be small.
    - As the variance of the measurement error grows in relation to the variation in the true variable, the magnitude of the bias increases.
    - As a worst-case limit, if the true  $x$  doesn't vary across our sample of observations and all of the variation in our measure  $\tilde{x}$  is random noise, then the expected value of our coefficient is zero.
  - Best solution is getting a better measure.
  - Alternatives are instrumental variables or direct measurement of degree of measurement error.
    - For example, if an alternative, precise measure is available for some arguably random sub-sample of observations, then we can calculate the variance of the true variable and the variance of the measurement error and correct the estimate.
- **Omitted-variables bias**
  - We derived this result at the beginning of the multiple regression analysis
  - Omitted variable is included in error. If omitted variable is correlated with included variable, then OLS estimator of coefficient on included variable is biased and inconsistent.

- **Simultaneous-equations bias (simultaneity bias)**

- Suppose that  $y$  and  $x$  are part of a larger theoretical system of equations:
 
$$y = \beta_0 + \beta_1 x + \dots + u$$

$$x = \gamma_0 + \gamma_1 y + \dots + v$$
- The two variables are “jointly determined” and both are endogenous.
  - There is “feedback” from  $y$  to  $x$ , or “reverse causality” (actually, bidirectional)
- $u \Rightarrow y \Rightarrow x$ , so  $u$  and  $x$  are correlated
- Supply and demand curves are difficult to estimate because both  $q$  and  $p$  are endogenous
- We will study a simple supply-demand model in great detail when we talk about simultaneous-equation models shortly

### *Instrumental variables*

- Recall the method of moments analysis by which we derived the OLS estimators
  - We used the assumed population moment conditions
 
$$E(u) = 0, \text{cov}(x, u) = 0$$
 to derive the OLS normal equations as sample moment conditions:  $\frac{1}{n} \sum_{i=1}^n \hat{u}_i^2 = 0, \frac{1}{n} \sum_{i=1}^n x_i \hat{u}_i = 0$
  - If  $\text{cov}(x, u) \neq 0$ , then the population moment conditions are invalid and we will get biased and inconsistent estimators from the OLS sample moment conditions.
- The **instrumental-variables estimator** can be derived from the method of moments.
- As usual, suppose that  $y = \beta_0 + \beta_1 x + u$ , but suppose that  $\text{cov}(x, u) \neq 0$ .
- Let  $z$  be a variable with the following properties:
  - $z$  does not have a direct effect on  $y$ . It does not belong in the equation alongside  $x$ . ( $z$  affects  $y$  only through  $x$ , not independently.)
  - $z$  is exogenous. It is not correlated with  $u$ .
  - $z$  is strongly correlated with  $x$ , the endogenous regressor.
- This makes  $z$  a valid instrumental variable.
- We can exploit  $\text{cov}(z, u) = 0$  as our second moment condition in place of  $\text{cov}(x, u) = 0$ , which is not true for this model.
- The sample moment conditions are

$$\sum_{i=1}^n \hat{u}_i = \sum_{i=1}^n (y_i - \hat{\beta}_0^{IV} - \hat{\beta}_1^{IV} x_i) = 0$$

$$\sum_{i=1}^n \hat{u}_i z_i = \sum_{i=1}^n z_i (y_i - \hat{\beta}_0^{IV} - \hat{\beta}_1^{IV} x_i) = 0.$$

- Solving the normal equations yields  $\hat{\beta}_1^{IV} = \frac{\sum_{i=1}^n (z_i - \bar{z})(y_i - \bar{y})}{\sum_{i=1}^n (z_i - \bar{z})(x_i - \bar{x})}$ .
- Compare this to the standard OLS slope estimator  $\hat{\beta}_1^{OLS} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})}$ .
- In matrix terms,  $\hat{\beta}^{IV} = (\mathbf{Z}'\mathbf{X})^{-1} \mathbf{Z}'\mathbf{y}$  vs.  $\hat{\beta}^{OLS} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}$
- Properties of IV estimator:
  - Consistent as long as  $z$  is exogenous
  - Asymptotically normal
  - $\hat{\beta}_1^{IV} \sim N \left( \beta_1, \frac{\sigma^2}{r_{xz}^2 \sum_{i=1}^n (x_i - \bar{x})^2} \right), r_{xz} \equiv \widehat{\text{corr}}(x, z)$
  - As usual, we estimate  $\sigma^2$  by  $\sigma_{IV}^2 = \frac{\sum_{i=1}^n (y_i - \hat{\beta}_0^{IV} - \hat{\beta}_1^{IV} x_i)^2}{N - 2}$
- Weak instruments: If  $r_{xz}$  is near zero, then the variance of  $\hat{\beta}_1^{IV}$  is large and the IV estimator is unreliable.

## Two-stage least squares

- What if we have multiple strong instruments and/or multiple endogenous regressors in a multiple regression?
- With more instruments than endogenous regressors, we have an “overidentified” system with alternative choices of instruments.
  - Suppose that  $x_k$  is endogenous but the first  $k - 1$  regressors are exogenous
  - Suppose that  $z_1$  through  $z_L$  are  $L$  valid instruments
  - Any linear combination of the instruments is admissible
  - Let's choose the one that is as highly correlated with  $x_k$  as possible:
    - To get that, we regress  $x_k = \gamma_0 + \gamma_1 x_1 + \dots + \gamma_{k-1} x_{k-1} + \theta_1 z_1 + \dots + \theta_L z_L + v_k$  and use the fitted values  $\hat{x}_k$  as the instrument for  $x_k$
- This amounts to doing two separate regressions, the **first-stage** regression of  $x_k$  on the exogenous  $x$  variables and the instruments  $z$ , then a **second-stage** regression of  $y = \beta_0 + \beta_1 x_1 + \dots + \beta_{k-1} x_{k-1} + \beta_k \hat{x}_k + u^*$
- The estimators of  $\beta$  from the second-stage regression are called 2SLS estimators.

- But it's not exactly like doing two separate regressions because our estimator of the error variance uses the actual values of  $x_K$  rather than the fitted values:

$$\hat{\sigma}_{IV}^2 = \frac{\sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i,1} - \dots - \hat{\beta}_k x_{i,k})^2}{n - k}$$

- (If you do the second regression manually substituting in the fitted values, Stata will use the fitted values to calculate the residuals rather than the actual.)
- 2SLS easily extends to multiple endogenous regressors, as long as there are more independent instruments than endogenous regressors.
  - Suppose there are  $G$  “good” exogenous regressors,  $B = k - G$  “bad” endogenous regressors, and  $L$  “lucky” instrumental variables.
  - $L > B$  means overidentified,  $L = B$  is just identified,  $L < B$  means underidentified (and can't be estimated by IV)
  - $y = \beta_0 + \beta_1 x_1 + \dots + \beta_G x_G + \beta_{G+1} x_{G+1} + \dots + \beta_k x_k + u$
  - First-stage regressions:  

$$x_{G+j} = \gamma_{0j} + \gamma_{1j} x_1 + \dots + \gamma_{Gj} x_G + \theta_{1j} z_1 + \dots + \theta_{Lj} z_L + v_j, \quad j = 1, \dots, B$$
  - Get fitted values:  

$$\hat{x}_{G+j} = \hat{\gamma}_{0j} + \hat{\gamma}_{1j} x_1 + \dots + \hat{\gamma}_{Gj} x_G + \hat{\theta}_{1j} z_1 + \dots + \hat{\theta}_{Lj} z_L, \quad j = 1, \dots, B$$
  - Regress original equation replacing endogenous regressors with fitted values  

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_G x_G + \beta_{G+1} \hat{x}_{G+1} + \dots + \beta_k \hat{x}_K + u^*$$
- To implement 2SLS in Stata, use `ivregress 2sls depvar exvars (endvars = instvars) , options`

### *Overidentification and generalized method of moments*

- If we have additional instruments beyond the minimum (i.e., an overidentified system), then we have more information than we need to estimate the model.
- Suppose that  $z_1$  and  $z_2$  are both valid instruments for endogenous  $x$
- All three moment conditions:

$$\sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = \hat{m}_1 = 0$$

$$\sum_{i=1}^n z_{i,1} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = \hat{m}_2 = 0$$

$$\sum_{i=1}^n z_{i,2} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = \hat{m}_3 = 0$$

are theoretically true.

- This opens the door for two possibilities:
  - We can determine the degree to which we cannot satisfy all three of these conditions simultaneously and use that as evidence of whether the model's

assumptions are valid. (If the model is perfect, then all three should be zero except for sampling error.)

- These are **specification tests** discussed below
- We can think about alternative estimators (called **GMM estimators**) that would minimize a weighted average of the squares of the  $m$  moments.
- 2SLS is a GMM estimator with a particular weighting of the moment conditions.

### *Instrument strength*

- A strong instrument must provide correlation with part of the endogenous regressor that is *not* explained by the other (exogenous) regressors.
- Regression of  $x_K = \gamma_0 + \gamma_1 x_1 + \dots + \gamma_{k-1} x_{k-1} + \theta_1 z_1 + \nu_k$  allows us to test  $\theta_1 = 0$  with a standard  $F = t^2$  test.
  - However, conventional wisdom says that the instrument is weak unless  $F > 10$  rather than the standard critical values for testing this hypothesis.
  - This test can be applied with multiple instruments and one endogenous regressor, with 10 still being the traditional threshold for weak instruments.
  - (There are appropriate criteria for more complicated cases.)

### *Specification tests*

- If the model is overidentified, then we can do two kinds of tests:
  - A **Hausman test** of whether the  $x$  variables that we are treating as endogenous truly are endogenous
  - A test of the overidentifying restrictions, which can be interpreted as a test of instrument validity
- **Hausman test**
  - $H_0 : \text{cov}(x, u) = 0, H_1 : \text{cov}(x, u) \neq 0$
  - Under null hypothesis, OLS is consistent and efficient, IV is consistent but inefficient. Since both are consistent,  $q \equiv \hat{\beta}^{OLS} - \hat{\beta}^{IV} \rightarrow 0$  in large samples
  - Under alternative hypothesis, OLS is inconsistent but IV is consistent, so  $q = \hat{\beta}^{OLS} - \hat{\beta}^{IV} \rightarrow c \neq 0$  in large samples.
  - Stata command `hausman` implements the procedure
  - Wooldridge gives alternative implementation **adding residuals from first-stage regression to OLS of original equation** and testing whether they are significant with a standard  $t$  test
    - If they are significant, then the  $\hat{\nu}$  residual is correlated with  $y$ . This is the part of  $x_k$  that is not exogenous, so if it is correlated with  $y$  then there is endogeneity and we need 2SLS because OLS is biased and inconsistent.
- **Tests for instrument validity**

- Is  $z$  correlated with  $u$ ?
  - Can't do direct test because we can't get consistent estimators for  $u$  without valid instruments, and we can't know whether instruments are valid without consistent estimator of  $u$ .
  - With extra instruments (overidentified model), we can use some to test the others.
- **LM test:** Do 2SLS/IV, get residuals, regress  $\hat{u}$  on all  $z$  instruments and exogenous regressors. Under null hypothesis that all instruments (and the overall model specification) are valid,  $nR^2$  from this regression  $\sim \chi^2$  with  $L - B$  degrees of freedom.
- The  **$J$  statistic** is another common test of overidentifying restrictions:
  - As above, regress the 2SLS/IV residuals on the exogenous variables in the equation and all the instruments.
  - Compute the  $F$  statistic for the null hypothesis that the coefficients on the instruments are zero.
  - The test statistic  $LF$  (where  $L$  is the number of instruments) is asymptotically distributed as a  $\chi^2$  with  $L - B$  degrees of freedom (number of instruments – number of endogenous regressors = number of overidentifying restrictions to be tested).
  - Why does the  $J$  test or the LM test work?
    - If the instruments are exogenous, then they should not be correlated with  $y$  **except through their effects on  $x$** .
    - The 2SLS residuals are the part of  $y$  that is orthogonal to the part of  $z$  that works through  $x$ .
    - If that is the only correlation that  $z$  has with  $y$  (there is no direct effect either direction), then the residuals should be uncorrelated with  $z$ , conditional on the other  $x$  variables, the included exogenous variables.
- Rejection of the null hypothesis tells us that at least one of the overidentifying restrictions does not hold, which may mean that one or more of the instruments is invalid.

### *Examples of IV regression*

- **Article about Waldman's autism analysis in WSJ**
  - Waldman
    - Unit of observation = state
    - Dependent variable = autism prevalence
    - Endogenous regressor = TV watching
    - Instrument = rainfall

- Angrist
  - Unit = individual male
  - Dependent = future earnings
  - Endogenous regressor = military service during Vietnam War
  - Instrument = draft lottery
- Angrist & Krueger
  - Unit = individual student
  - Dependent = future earnings
  - Endogenous regressor = age of HS graduation
  - Instrument = birthdate & restrictions on entry
- Levitt
  - Unit = city?
  - Dependent = violent crime
  - Endogenous regressor = number of police
  - Instrument = elections
- Hoxby
  - Unit = school
  - Dependent = test scores
  - Endogenous regressor = competition between school districts
  - Instrument = streams
- Iyer
  - Unit = Indian region
  - Dependent = quality of public goods
  - Endogenous regressor = British rule
  - Instrument = death of heirless indigenous ruler
- Feyrer & Sacerdote
  - Unit = island country
  - Dependent = wealth
  - Endogenous regressor = date of colonization
  - Instrument = wind patterns
- Oster
  - Unit = ?
  - Dependent = risky sex in presence of HIV
  - Endogenous regressor = degree of confrontation of HIV
  - Instrument = distance from origin of HIV
- Bannedsen et al.
  - Unit = Danish family firm
  - Dependent = firm profitability
  - Endogenous regressor = nepotism (did firm stay in family)
  - Instrument = first-born son in family
- Olken

- Unit = region of Indonesia
- Dependent = social capital
- Endogenous regressor = amount of TV/radio consumed
- Instrument = signal strength of reception in region

• **HGL Problem 10.1**

- Now done as daily problem earlier.

	(1) <i>RENT</i>	(2) <i>MDHOUSE</i>	(3) <i>MDHOUSE</i>	(4) <i>RENT</i>	(5) <i>RENT</i>	(6) <i>EHAT</i>
<i>C</i>	125.9 (14.19)	-18.67 (12.00)	7.225 (8.936)	120.7 (12.43)	120.7 (15.71)	-62.85 (26.95)
<i>PCTURBAN</i>	0.525 (0.249)	0.182 (0.115)	0.616 (0.131)	0.0815 (0.244)	0.0815 (0.305)	-0.283 (0.258)
<i>MDHOUSE</i>	1.521 (0.228)			2.240 (0.268)	2.240 (0.339)	
<i>FAMINC</i>		2.731 (0.682)				4.448 (1.532)
<i>REG2</i>		-5.095 (4.122)				-6.768 (9.262)
<i>REG3</i>		-1.778 (4.073)				4.847 (9.151)
<i>REG4</i>		13.41 (4.048)				-18.77 (9.096)
<i>VHAT</i>				-1.589 (0.398)		
<i>N</i>	50	50	50	50	50	50
<i>R</i> <sup>2</sup>	0.669	0.691	0.317	0.754	0.599	0.226
<i>SSE</i>	20259.6	3767.6	8322.2	15054.0	24565.7	19019.9

○ *Standard errors in parentheses*

○ **Description**

- State data
- RENT paid on houses is DV
- MDHOUSE is median house price (in equation)
- PCTURBAN is % of population in urban areas (in equation)
- FAMINC is median family income (instrument)
- REG1-REG3 are regional dummies (instruments)
- (1) is OLS, (2) is first-stage regression, (3) is regression on exogenous  $x$  without instruments
- (a) Why might MDHOUSE be endogenous?
- (b) Test if instruments are weak
  - Instruments should have  $F > 10$

$$F = \frac{(SSE_R - SSE_U)}{J\hat{\sigma}_U^2} = \frac{(8322.2 - 3767.6)}{4(3767.6/(50 - 6))} = \frac{4554.6}{4(85.6)} = 13.3$$

- 
- (c) Hausman test using (4)
  - VHAT is the residual from the first-stage regression (2)
  - (4) adds VHAT to second-stage regression
    - Can use either fitted or actual MDHOUSE values here because actual = fitted + VHAT
  - $H_0$  is that MDHOUSE is exogenous, in which case VHAT should not add anything to the OLS structural regression
  - $t = -1.589/0.389 = -3.99$ . Critical value = 2.01, so reject exogeneity of MDHOUSE
  - We can't use OLS and need to do 2SLS
- (d) Note that (4) and (5) have identical coefficients because VHAT is orthogonal to regressors (PCTURBAN) by construction, so no omitted-variable bias
  - Comparing 2SLS (5) to OLS (1)
    - PCTURBAN coefficient is smaller and se is larger
    - MDHOUSE coefficient is larger and very significant
- (e) Testing overidentifying restrictions:
  - Regress EHAT (2SLS residuals) on instruments and exogenous regressors
  - If only effect of instruments on RENT is through MDHOUSE, then this regression should have no correlation:
    - Effects of PCTURBAN have already been taken out in 2SLS regression and can't be in EHAT
    - Instruments should have no independent effect if they are valid
  - $NR^2 = 50 \times 0.226 = 11.3$ . Critical value for  $\chi^2(3) = 7.815$ , so we reject the null hypothesis of instrument validity.
  - This test does not indicate which of the instruments might be at fault, or indeed where in the collective hypotheses of the model (exogeneity, absences from equation, correct specification) the problem lies