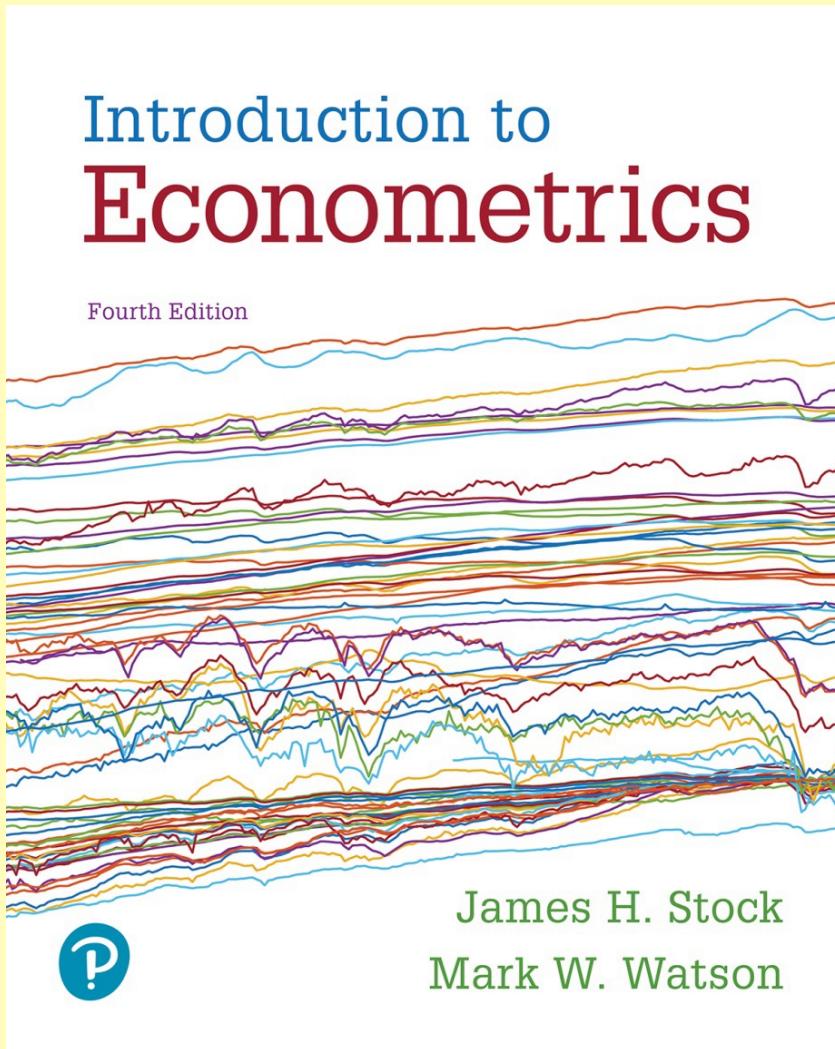


# Introduction to Econometrics

- Fourth Edition



- Chapter 12

Instrumental Variables  
Regressions

# Outline

1. IV Regression: Why and What; Two Stage Least Squares
2. The General IV Regression Model
3. Checking Instrument Validity
  - a) Weak and strong instruments
  - b) Instrument exogeneity
4. Application: Demand for cigarettes
5. Examples: Where Do Instruments Come From?

Instrumental variables (IV) regression is a general way to obtain a consistent estimator of the unknown causal coefficients when the regressor,  $X$ , is correlated with the error term,  $u$ .

To understand how IV regression works, think of the variation in  $X$  as having two parts: one part that, for whatever reason, is correlated with  $u$  (this is the part that causes the problems) and a second part that is uncorrelated with  $u$ .

If you had information that allowed you to isolate the second part, you could focus on those variations in  $X$  that are uncorrelated with  $u$  and disregard the variations in  $X$  that bias the OLS estimates.

This is, in fact, what IV regression does.

The information about the movements in  $X$  that are uncorrelated with  $u$  is gleaned from one or more additional variables, called **instrumental variables** or simply **instruments**.

# IV Regression: Why?

Three important threats to internal validity are:

- Omitted variable bias from a variable that is correlated with  $X$  but is unobserved (so cannot be included in the regression) and for which there are inadequate control variables;
- Simultaneous causality bias ( $X$  causes  $Y$ ,  $Y$  causes  $X$ );
- Errors-in-variables bias ( $X$  is measured with error) All three problems result in  $E(u|X) \neq 0$ .
- Instrumental variables regression can eliminate bias when  $E(u|X) \neq 0$  – using an *instrumental variable* (IV),  $Z$ .

# The IV Estimator with a Single Regressor and a Single Instrument (SW Section 12.1)

$$Y_i = \beta_0 + \beta_1 X_i + u_i$$

- The goal is an estimate of the causal effect  $\beta_1$ . However,  $X$  is correlated with the error term, and we cannot solve the problem simply by including control variables.
- Instrumental variables (IV) regression breaks  $X$  into two parts: a part that might be correlated with  $u$ , and a part that is not. By isolating the part that is not correlated with  $u$ , it is possible to estimate  $\beta_1$ .
- This is done using an ***instrumental variable***,  $Z_i$ , which is correlated with  $X_i$  but uncorrelated with  $u_i$ .

# Terminology: Endogeneity and Exogeneity

An **endogenous** variable is one that is correlated with  $u$

An **exogenous** variable is one that is uncorrelated with  $u$

In IV regression, we focus on the case that  $X$  is endogenous and there is an instrument,  $Z$ , which is exogenous.

*Digression on terminology:* “Endogenous” literally means “determined within the system.” If  $X$  is jointly determined with  $Y$ , then a regression of  $Y$  on  $X$  is subject to simultaneous causality bias. But this definition of endogeneity is too narrow because IV regression can be used to address OV bias and errors-in-variable bias. Thus we use the broader definition of endogeneity above.

# Two Conditions for a Valid Instrument

$$Y_i = \beta_0 + \beta_1 X_i + u_i$$

For an instrumental variable (an “*instrument*”)  $Z$  to be valid, it must satisfy two conditions:

1. ***Instrument relevance***:  $\text{corr}(Z_i, X_i) \neq 0$
2. ***Instrument exogeneity***:  $\text{corr}(Z_i, u_i) = 0$

Suppose for now that you have such a  $Z_i$  (we’ll discuss how to find instrumental variables later). How can you use  $Z_i$  to estimate  $\beta_1$ ?

# The IV estimator with one $X$ and one $Z$ (1 of 7)

## Explanation #1: Two Stage Least Squares (TSLS)

As it sounds, TSLS has two stages – two regressions:

- (1) Isolate the part of  $X$  that is uncorrelated with  $u$  by regressing  $X$  on  $Z$  using OLS:

$$X_i = \pi_0 + \pi_1 Z_i + v_i \quad (1)$$

- Because  $Z_i$  is uncorrelated with  $u_i$ ,  $\pi_0 + \pi_1 Z_i$  is uncorrelated with  $u_i$ . We don't know  $\pi_0$  or  $\pi_1$  but we have estimated them, so...
- Compute the predicted values of  $X_i$ , where  $\hat{X}_i = \hat{\pi}_0 + \hat{\pi}_1 Z_i$ ,  $i = 1, \dots, n$ .

# The IV estimator with one $X$ and one $Z$ (2 of 7)

- (2) Replace  $\underline{X}_i$  by  $\hat{X}_i$  in the regression of interest: regress  $\underline{Y}$  on  $\hat{X}_i$  using OLS:

$$Y_i = \beta_0 + \beta_1 \hat{X}_i + u_i \quad (2)$$

- Because  $\hat{X}_i$  is uncorrelated with  $u_i$ , the first least squares assumption holds for regression (2). (This requires  $n$  to be large so that  $\pi_0$  and  $\pi_1$  are precisely estimated.)
- Thus, in large samples,  $\beta_1$  can be estimated by OLS using regression (2)
- The resulting estimator is called the Two Stage Least Squares (TSLS) estimator,  $\hat{\beta}_1^{TSLS}$ .

# Two Stage Least Squares: Summary

Suppose  $Z_i$ , satisfies the two conditions for a valid instrument:

1. **Instrument relevance**:  $\text{corr}(Z_i, X_i) \neq 0$
2. **Instrument exogeneity**:  $\text{corr}(Z_i, u_i) = 0$

Two-stage least squares:

Stage 1: Regress  $X_i$  on  $Z_i$  (including an intercept), obtain the predicted values  $\hat{X}_i$

Stage 2: Regress  $Y_i$  on  $\hat{X}_i$  (including an intercept); the coefficient on  $\hat{X}_i$  is the TSLS estimator,  $\hat{\beta}_1^{\text{TSLS}}$ .

$\hat{\beta}_1^{\text{TSLS}}$  is a consistent estimator of  $\beta_1$ .

# The IV estimator with one $X$ and one $Z$ (3 of 7)

## Explanation #2: A direct algebraic derivation

$$Y_i = \beta_0 + \beta_1 X_i + u_i$$

Thus,

$$\begin{aligned}\text{cov}(Y_i, Z_i) &= \text{cov}(\beta_0 + \beta_1 X_i + u_i, Z_i) \\ &= \text{cov}(\beta_0, Z_i) + \text{cov}(\beta_1 X_i, Z_i) + \text{cov}(u_i, Z_i) \\ &= 0 + \text{cov}(\beta_1 X_i, Z_i) + 0 \\ &= \beta_1 \text{cov}(X_i, Z_i)\end{aligned}$$

where  $\text{cov}(u_i, Z_i) = 0$  by instrument exogeneity; thus

$$\beta_1 = \frac{\text{cov}(Y_i, Z_i)}{\text{cov}(X_i, Z_i)}$$

# The IV estimator with one $X$ and one $Z$ (4 of 7)

$$\beta_1 = \frac{\text{cov}(Y_i, Z_i)}{\text{cov}(X_i, Z_i)}$$

The IV estimator replaces these population covariances with sample covariances:

$$\hat{\beta}_1^{TSLS} = \frac{s_{YZ}}{s_{XZ}},$$

where  $s_{YZ}$  and  $s_{XZ}$  are the sample covariances. This is the TSLS estimator – just a different derivation!

# The IV estimator with one $X$ and one $Z$ (5 of 7)

## Explanation #3: Derivation from the “reduced form”

The “reduced form” relates  $Y$  to  $Z$  and  $X$  to  $Z$ :

$$X_i = \pi_0 + \pi_1 Z_i + \nu_i$$

$$Y_i = \gamma_0 + \gamma_1 Z_i + w_i$$

where  $w_i$  is an error term. Because  $Z$  is exogenous,  $Z$  is uncorrelated with both  $\nu_i$  and  $w_i$ .

*The idea:* A unit change in  $Z_i$  results in a change in  $X_i$  of  $\pi_1$  and a change in  $Y_i$  of  $\gamma_1$ . Because that change in  $X_i$  arises from the exogenous change in  $Z_i$ , that change in  $X_i$  is exogenous. Thus an exogenous change in  $X_i$  of  $\pi_1$  units is associated with a change in  $Y_i$  of  $\gamma_1$  units – so the effect on  $Y$  of an exogenous change in  $X$  is  $\beta_1 = \gamma_1/\pi_1$ .

# The IV estimator with one $X$ and one $Z$ (6 of 7)

The math:

$$X_i = \pi_0 + \pi_1 Z_i + v_i$$

$$Y_i = \gamma_0 + \gamma_1 Z_i + w_i$$

Solve the  $X$  equation for  $Z$ :

$$Z_i = -\pi_0/\pi_1 + (1/\pi_1)X_i - (1/\pi_1)v_i$$

Substitute this into the  $Y$  equation and collect terms:

$$\begin{aligned} Y_i &= \gamma_0 + \gamma_1 Z_i + w_i \\ &= \gamma_0 + \gamma_1 [-\pi_0/\pi_1 + (1/\pi_1)X_i - (1/\pi_1)v_i] + w_i \\ &= [\gamma_0 - \pi_0\gamma_1/\pi_1] + (\gamma_1/\pi_1)X_i + [w_i - (\gamma_1/\pi_1)v_i] \\ &= \beta_0 + \beta_1 X_i + u_i, \end{aligned}$$

where  $\beta_0 = \gamma_0 - \pi_0\gamma_1/\pi_1$ ,  $\beta_1 = \gamma_1/\pi_1$ , and  $u_i = w_i - (\gamma_1/\pi_1)v_i$ .

# The IV estimator with one $X$ and one $Z$ (7 of 7)

$$X_i = \pi_0 + \pi_1 Z_i + v_i$$

$$Y_i = \gamma_0 + \gamma_1 Z_i + w_i$$

yields

$$Y_i = \beta_0 + \beta_1 X_i + u_i,$$

where  $\beta_1 = \gamma_1 / \pi_1$

*Interpretation:* An exogenous change in  $X_i$  of  $\pi_1$  units is associated with a change in  $Y_i$  of  $\gamma_1$  units – so the effect on  $Y$  of an exogenous unit change in  $X$  is  $\beta_1 = \gamma_1 / \pi_1$ .

# Example #1: Effect of Studying on Grades (1 of 6)

What is the effect on grades of studying for an additional hour per day?

$$Y = \text{GPA}$$

$$X = \text{study time (hours per day)}$$

Data: grades and study hours of college freshmen.

*Would you expect the OLS estimator of  $\beta_1$  (the effect on GPA of studying an extra hour per day) to be unbiased? Why or why not?*

# Example #1: Effect of Studying on Grades (2 of 6)

Stinebrickner, Ralph and Stinebrickner, Todd R. (2008) "The Causal Effect of Studying on Academic Performance," *The B.E. Journal of Economic Analysis & Policy*: Vol. 8: Iss. 1 (Frontiers), Article 14.

- $n = 210$  freshman at Berea College (Kentucky) in 2001
- $Y =$  first-semester GPA
- $X =$  average study hours per day (time use survey)
- Roommates were randomly assigned
- $Z = 1$  if roommate brought video game,  $= 0$  otherwise

Do you think  $Z_i$  (whether a roommate brought a video game) is a valid instrument?

1. Is it relevant (correlated with  $X$ )?
2. Is it exogenous (uncorrelated with  $u$ )?

# Example #1: Effect of Studying on Grades (3 of 6)

$$X = \pi_0 + \pi_1 Z + v_i$$

$$Y = \gamma_0 + \gamma_1 Z + w_i$$

$Y = GPA$  (4 point scale)

$X = time spent studying$  (hours per day)

$Z = 1$  if roommate brought video game,  $= 0$  otherwise

## Stinebrinckner and Stinebrinckner's findings

$$\hat{\pi}_1 = -.668$$

$$\hat{\gamma}_1 = -.241$$

$$\hat{\beta}_1^{IV} = \frac{\hat{\gamma}_1}{\hat{\pi}_1} = \frac{-.241}{-.668} = 0.360$$

What are the units? Do these estimates make sense in a real-world way? (Note: They actually ran the regressions including additional regressors – more on this later.)

# Consistency of the TSLS estimator

$$\hat{\beta}_1^{TSLS} = \frac{s_{YZ}}{s_{XZ}}$$

The sample covariances are consistent:  $s_{YZ} \xrightarrow{p} \text{cov}(Y, Z)$

and  $s_{XZ} \xrightarrow{p} \text{cov}(X, Z)$ . Thus,

$$\hat{\beta}_1^{TSLS} = \frac{s_{YZ}}{s_{XZ}} \xrightarrow{p} \frac{\text{cov}(Y, Z)}{\text{cov}(X, Z)} = \beta_1$$

- The instrument relevance condition,  $\text{cov}(X, Z) \neq 0$ , ensures that you don't divide by zero.

## Example #2: Supply and demand for butter (1 of 2)

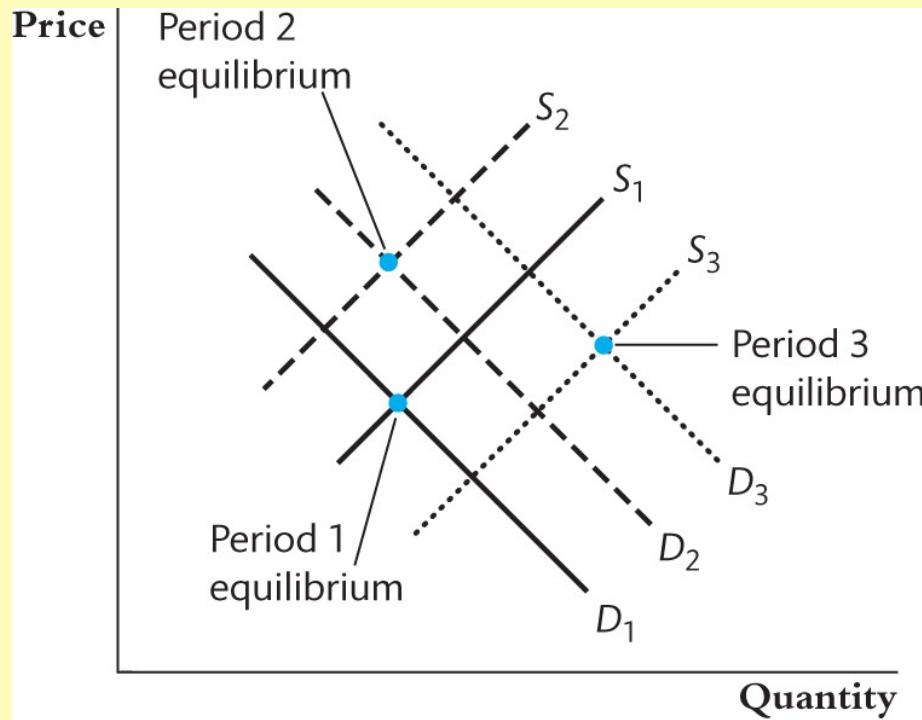
IV regression was first developed to estimate demand elasticities for agricultural goods, for example, butter:

$$\ln(Q_i^{butter}) = \beta_0 + \beta_1 \ln(P_i^{butter}) + u_i$$

- $\beta_1$  = price elasticity of demand for butter = percent change in quantity for a 1% change in price (recall log-log specification discussion)
- Data: observations on price and quantity of butter for different years
- The OLS regression of  $\ln(Q_i^{butter})$  on  $\ln(P_i^{butter})$  suffers from simultaneous causality bias (*why?*)

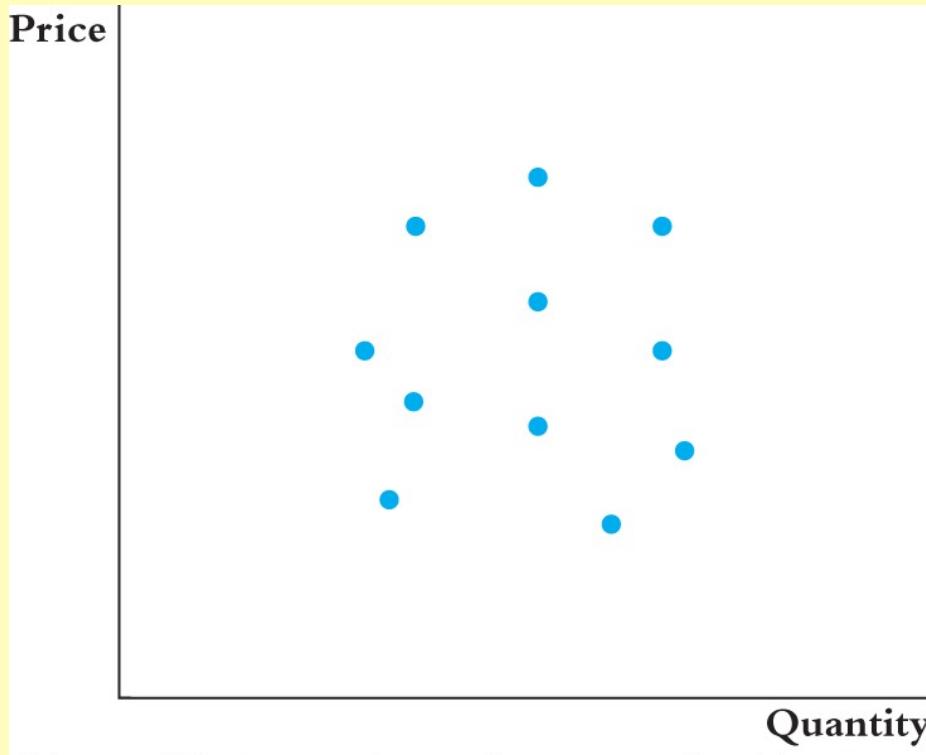
## Example #2: Supply and demand for butter (2 of 2)

Simultaneous causality bias in the OLS regression of  $\ln(Q_i^{butter})$  on  $\ln(P_i^{butter})$  arises because price and quantity are determined by the interaction of demand *and* supply:



(a) Demand and supply in three time periods

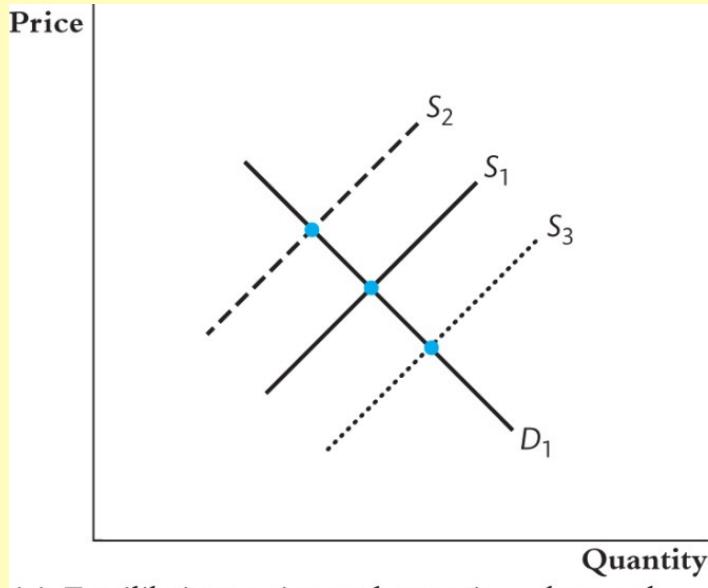
# This interaction of demand and supply produces data like:



(b) Equilibrium price and quantity for 11 time periods

*Would a regression using these data produce the demand curve?*

# But...what would you get if only supply shifted?



(c) Equilibrium price and quantity when only the supply curve shifts

- TSLS estimates the demand curve by isolating shifts in price and quantity that arise from shifts in supply.
- Z is a variable that shifts supply but not demand.

# TSLS in the supply-demand example (1 of 2)

$$\ln(Q_i^{butter}) = \beta_0 + \beta_1 \ln(P_i^{butter}) + u_i$$

Let  $Z$  = rainfall in dairy-producing regions.

Is  $Z$  a valid instrument?

- (1) Relevant?  $\text{corr}(\text{rain}_i, \ln(P_i^{butter})) \neq 0$ ?

*Plausibly:* insufficient rainfall means less grazing means less butter means higher prices

- (2) Exogenous?  $\text{corr}(\text{rain}_i, u_i) = 0$ ?

*Plausibly:* whether it rains in dairy-producing regions shouldn't affect demand for butter

# TSLS in the supply-demand example (2 of 2)

$$\ln(Q_i^{butter}) = \beta_0 + \beta_1 \ln(P_i^{butter}) + u_i$$

$Z_i = rain_i$  = rainfall in dairy-producing regions.

Stage 1: regress  $\ln(P_i^{butter})$  on rain,  $\widehat{\ln(P_i^{butter})}$

$\widehat{\ln(P_i^{butter})}$  isolates changes in log price that arise from supply (part of supply, at least)

Stage 2: regress  $\ln(Q_i^{butter})$  on  $\widehat{\ln(P_i^{butter})}$

The regression counterpart of using shifts in the supply curve to trace out the demand curve.

## Example #3: Test scores and class size (1 of 2)

- The California test score/class size regressions still could have OV bias (e.g. parental involvement).
- In principle, this bias can be eliminated by IV regression (TSLS).
- IV regression requires a valid instrument, that is, an instrument that is:
  1. relevant:  $\text{corr}(Z_i, STR_i) \neq 0$
  2. exogenous:  $\text{corr}(Z_i, u_i) = 0$

## Example #3: Test scores and class size (2 of 2)

Here is a (hypothetical) instrument:

- some districts, randomly hit by an earthquake, “double up” classrooms:  
 $Z_i = \text{Quake}_i = 1 \text{ if hit by quake, } = 0 \text{ otherwise}$
- *Do the two conditions for a valid instrument hold?*
- The earthquake makes it as *if* the districts were in a random assignment experiment. Thus, the variation in *STR* arising from the earthquake is exogenous.
- The first stage of TSLS regresses *STR* against *Quake*, thereby isolating the part of *STR* that is exogenous (the part that is “as if ” randomly assigned)

# Inference using TSLS (1 of 5)

- In large samples, the sampling distribution of the TSLS estimator is normal
- Inference (hypothesis tests, confidence intervals) proceeds in the usual way, e.g.  $\pm 1.96SE$
- The idea behind the large-sample normal distribution of the TSLS estimator is that – like all the other estimators we have considered – it involves an average of mean zero i.i.d. random variables, to which we can apply the CLT.
- Here is the math (SW App. 12.3)...

# Inference using TSLS (2 of 5)

$$\hat{\beta}_1^{TSLS} = \frac{S_{YZ}}{S_{XZ}} = \frac{\frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})(Z_i - \bar{Z})}{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Z_i - \bar{Z})} = \frac{\sum_{i=1}^n Y_i(Z_i - \bar{Z})}{\sum_{i=1}^n X_i(Z_i - \bar{Z})}$$

Substitute in  $Y_i = \beta_0 + \beta_1 X_i + u_i$  and simplify:

$$\hat{\beta}_1^{TSLS} = \frac{\beta_1 \sum_{i=1}^n X_i(Z_i - \bar{Z}) + \sum_{i=1}^n u_i(Z_i - \bar{Z})}{\sum_{i=1}^n X_i(Z_i - \bar{Z})}$$

SO...

# Inference using TSLS (3 of 5)

$$\hat{\beta}_1^{TSLS} = \beta_1 + \frac{\sum_{i=1}^n u_i (Z_i - \bar{Z})}{\sum_{i=1}^n X_i (Z_i - \bar{Z})}.$$

So

$$\hat{\beta}_1^{TSLS} - \beta_1 = \frac{\sum_{i=1}^n u_i (Z_i - \bar{Z})}{\sum_{i=1}^n X_i (Z_i - \bar{Z})}$$

Multiply through by  $\sqrt{n}$ :

$$\sqrt{n}(\hat{\beta}_1^{TSLS} - \beta_1) = \frac{\frac{1}{\sqrt{n}} \sum_{i=1}^n (Z_i - \bar{Z}) u_i}{\frac{1}{n} \sum_{i=1}^n X_i (Z_i - \bar{Z})}$$

# Inference using TSLS (4 of 5)

$$\sqrt{n}(\hat{\beta}_1^{TSLS} - \beta_1) = \frac{\frac{1}{\sqrt{n}} \sum_{i=1}^n (Z_i - \bar{Z}) u_i}{\frac{1}{n} \sum_{i=1}^n X_i (Z_i - \bar{Z})}$$

$$\frac{1}{n} \sum_{i=1}^n X_i (Z_i - \bar{Z}) = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Z_i - \bar{Z}) \xrightarrow{p} \text{cov}(X, Z) \neq 0$$

$\frac{1}{\sqrt{n}} \sum_{i=1}^n (Z_i - \bar{Z}) u_i$  is distributed  $N(0, \text{var}[(Z - \mu_Z)u])$  (CLT)

so:  $\hat{\beta}_1^{TSLS}$  is approx. distributed  $N(\beta_1, \sigma_{\hat{\beta}_1^{TSLS}}^2)$ ,

where  $\sigma_{\hat{\beta}_1^{TSLS}}^2 = \frac{1}{n} \frac{\text{var}[(Z_i - \mu_Z)u_i]}{[\text{cov}(Z_i, X_i)]^2}$ .

where  $\text{cov}(X, Z) \neq 0$  because the instrument is relevant

# Inference using TSLS (5 of 5)

$\hat{\beta}_1^{TSLS}$  is approx. distributed  $N(\beta_1, \sigma_{\hat{\beta}_1^{TSLS}}^2)$ ,

- Statistical inference proceeds in the usual way.
- The justification is (as usual) based on large samples
- This all assumes that the instruments are valid – we'll discuss what happens if they aren't valid shortly.
- ***Important note on standard errors:***
  - The OLS standard errors from the second stage regression aren't right – they don't take into account the estimation in the first stage ( $\hat{X}_i$  is estimated).
  - Instead, use a single specialized command that computes the TSLS estimator and the correct SEs.
  - As usual, use heteroskedasticity-robust SEs

# Example #4: Demand for Cigarettes (1 of 3)

$$\ln(Q_i^{\text{cigarettes}}) = \beta_0 + \beta_1 \ln(P_i^{\text{cigarettes}}) + u_i$$

Why is the OLS estimator of  $\beta_1$  likely to be biased?

- Data set: Panel data on annual cigarette consumption and average prices paid (including tax), by state, for the 48 continental US states, 1985–1995.
- Proposed instrumental variable:
  - $Z_i$  = general sales tax per pack in the state =  $SalesTax_i$ ,
  - Do you think this instrument is plausibly valid?
    1. Relevant?  $\text{corr}(SalesTax_i, \ln(P_i^{\text{cigarettes}})) \neq 0$ ?
    2. Exogenous?  $\text{corr}(SalesTax_i, u_i) = 0$ ?

# Example #4: Demand for Cigarettes (2 of 3)

For now, use data from 1995 only.

First stage OLS regression:

$$\ln(\overline{P}_i^{\text{cigarettes}}) = 4.63 + .031 \text{SalesTax}_i, n = 48$$

Second stage OLS regression:

$$\ln(\overline{Q}_i^{\text{cigarettes}}) = 9.72 - 1.08 \ln(\overline{P}_i^{\text{cigarettes}}), n = 48$$

Combined TSLS regression with correct, heteroskedasticity-robust standard errors:

$$\ln(\overline{Q}_i^{\text{cigarettes}}) = 9.72 - 1.08, \ln(\overline{P}_i^{\text{cigarettes}}) n = 48$$
$$(1.53) (0.32)$$

# **STATA Example: Cigarette demand, First stage**

## **Instrument = Z = rtaxso = general sales tax (real \$/pack)**

**x            z**

```
. reg lravgprs rtaxso if year==1995, r
```

Regression with robust standard errors

Number of obs = 48  
F( 1, 46) = 40.39  
Prob > F = 0.0000  
R-squared = 0.4710  
Root MSE = .09394

---

	Robust					
lravgprs	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
<hr/>						
rtaxso	.0307289	.0048354	6.35	0.000	.0209956	.0404621
_cons	4.616546	.0289177	159.64	0.000	4.558338	4.674755

---

**X-hat**

```
. predict lravphat
```

**Now we have the predicted values from the 1<sup>st</sup> stage**



# Second stage

**y       $X\text{-hat}$**

```
. reg lpackpc lravphat if year==1995, r
```

Regression with robust standard errors

Number of obs = 48  
F( 1, 46) = 10.54  
Prob > F = 0.0022  
R-squared = 0.1525  
Root MSE = .22645

		Robust				
	lpackpc	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
-----+-----						
	lravphat	-1.083586	.3336949	-3.25	0.002	-1.755279 -.4118932
	_cons	9.719875	1.597119	6.09	0.000	6.505042 12.93471

-----+-----

- These coefficients are the TSLS estimates
- The standard errors are wrong because they ignore the fact that the first stage was estimated

# Combined into a single command

**Y**            **X**            **Z**

```
. ivregress 2sls lpackpc (lravgprs = rtaxso) if year==1995, vce(robust);
```

Instrumental variables (2SLS) regression

Number of obs = 48  
Wald chi2(1) = 12.05  
Prob > chi2 = 0.0005  
R-squared = 0.4011  
Root MSE = .18635

Robust						
lpackpc	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
lravgprs	-1.083587	.3122035	-3.47	0.001	-1.695494	-.471679
_cons	9.719876	1.496143	6.50	0.000	6.78749	12.65226

Instrumented: lravgprs *This is the endogenous regressor*

Instruments: rtaxso *This is the instrumental variable*

Estimated cigarette demand equation:

$$\ln(Q_i^{\text{cigarettes}}) = 9.72 - 1.08 \ln(P_i^{\text{cigarettes}}), n = 48$$

(1.53) (0.31)

# Summary of IV Regression with a Single $X$ and $Z$

- A valid instrument  $Z$  must satisfy two conditions:
  1. *relevance*:  $\text{corr}(Z_i, X_i) \neq 0$
  2. *exogeneity*:  $\text{corr}(Z_i, u_i) = 0$
- TSLS proceeds by first regressing  $X$  on  $Z$  to get  $\hat{X}$ , then regressing  $Y$  on  $\hat{X}$
- The key idea is that the first stage isolates part of the variation in  $X$  that is uncorrelated with  $u$
- If the instrument is valid, then the large-sample sampling distribution of the TSLS estimator is normal, so inference proceeds as usual

# The General IV Regression Model (SW Section 12.2)

- So far we have considered IV regression with a single endogenous regressor ( $X$ ) and a single instrument ( $Z$ ).
- We need to extend this to:
  - multiple endogenous regressors ( $X_1, \dots, X_k$ )
  - multiple included exogenous variables ( $W_1, \dots, W_r$ ) or control variables
  - multiple instrumental variables ( $Z_1, \dots, Z_m$ ). Having more (relevant) instruments can produce a smaller variance of TSLS: the  $R^2$  of the first stage increases, so you have more variation in  $\hat{X}$ .
- **New terminology:** identification & overidentification

# Identification (1 of 2)

- In general, a parameter is said to be ***identified*** if different values of the parameter produce different distributions of the data.
- In IV regression, whether the coefficients are identified depends on the relation between the number of instruments ( $m$ ) and the number of endogenous regressors ( $k$ )
  - Intuitively, if there are fewer instruments than endogenous regressors, we can't estimate  $\beta_1, \dots, \beta_k$ 
    - For example, suppose  $k = 1$  but  $m = 0$  (no instruments)!

# Identification (2 of 2)

The coefficients  $\beta_1, \dots, \beta_k$  are said to be:

- ***exactly identified*** if  $m = k$ .

There are just enough instruments to estimate  $\beta_1, \dots, \beta_k$ .

- ***overidentified*** if  $m > k$ .

There are more than enough instruments to estimate  $\beta_1, \dots, \beta_k$ .  
If so, you can test whether the instruments are valid (a test of the “overidentifying restrictions”) – we’ll return to this later

- ***underidentified*** if  $m < k$ .

There are too few instruments to estimate  $\beta_1, \dots, \beta_k$ . If so, you need to get more instruments!

# The General IV Regression Model: Summary

$$Y_i = \beta_0 + \beta_1 X_{1i} + \dots + \beta_k X_{ki} + \beta_{k+1} W_{1i} + \dots + \beta_{k+r} W_{ri} + u_i$$

- $Y_i$  is the **dependent variable**
- $X_{1i}, \dots, X_{ki}$  are the **endogenous regressors** (potentially correlated with  $u_i$ )
- $W_{1i}, \dots, W_{ri}$  are the **included exogenous regressors** (uncorrelated with  $u_i$ ) or **control variables** (included so that  $Z_i$  is uncorrelated with  $u_i$ , once the  $W$ 's are included)
- $\beta_0, \beta_1, \dots, \beta_{k+r}$  are the unknown regression coefficients
- $Z_{1i}, \dots, Z_{mi}$  are the  $m$  **instrumental variables** (the **excluded exogenous variables**)
- The coefficients are **overidentified** if  $m > k$ ; **exactly identified** if  $m = k$ ; and **underidentified** if  $m < k$ .

# TSLS with a Single Endogenous Regressor

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 W_{1i} + \dots + \beta_{1+r} W_{ri} + u_i$$

- $m$  instruments:  $Z_{1i}, \dots, Z_m$
- First stage
  - Regress  $X_1$  on *all* the exogenous regressors: regress  $X_1$  on  $W_1, \dots, W_r, Z_1, \dots, Z_m$ , and an intercept, by OLS
  - Compute predicted values  $\hat{X}_{1i}, i = 1, \dots, n$
- Second stage
  - Regress  $Y$  on  $\hat{X}_{1i}, W_1, \dots, W_r$ , and an intercept, by OLS
  - The coefficients from this second stage regression are the TSLS estimators, but *SEs* are wrong
- To get correct *SEs*, do this in a single step in your regression software

# Example #4: Demand for Cigarettes (3 of 3)

Suppose income is exogenous (this is plausible – *why?*), and we also want to estimate the income elasticity:

$$\ln(Q_i^{\text{cigarettes}}) = \beta_0 + \beta_1 \ln(P_i^{\text{cigarettes}}) + \beta_2 \ln(\text{income}_i) + u_i$$

We actually have two instruments:

$Z_{1i}$  = general sales tax<sub>*i*</sub>

$Z_{2i}$  = cigarette-specific tax<sub>*i*</sub>

- Endogenous variable:  $\ln(P_i^{\text{cigarettes}})$ , (“one X”)
- Included exogenous variable:  $\ln(\text{Income}_i)$  (“one W”)
- Instruments (excluded endogenous variables): general sales tax, cigarette-specific tax (“two Zs”)
- *Is  $\beta_1$  over-, under-, or exactly identified?*

# Example: Cigarette demand, one instrument

IV: rtaxso = real overall sales tax in state

Y W X Z

. ivreg lpackpc lperinc (lravgprs = rtaxso) if year==1995, r

IV (2SLS) regression with robust standard errors

Number of obs = 48  
F( 2, 45) = 8.19  
Prob > F = 0.0009  
R-squared = 0.4189  
Root MSE = .18957

	Robust				
lpackpc	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
lravgprs	-1.143375	.3723025	-3.07	0.004	-1.893231 -.3935191
lperinc	.214515	.3117467	0.69	0.495	-.413375 .842405
_cons	9.430658	1.259392	7.49	0.000	6.894112 11.9672

Instrumented: lravgprs

Instruments: lperinc rtaxso

STATA lists ALL the exogenous regressors  
as instruments - slightly different  
terminology than we have been using

- Running IV as a single command yields the correct SEs
- Use , r for heteroskedasticity-robust SEs

# Example: Cigarette demand, two instruments (1 of 2)

<i>Y</i>	<i>W</i>	<i>X</i>	<i>Z</i> <sub>1</sub>	<i>Z</i> <sub>2</sub>	
. ivreg lpackpc lperinc (lravgprs = rtaxso rtax) if year==1995, r;					
IV (2SLS) regression with robust standard errors					Number of obs = 48
					F( 2, 45) = 16.17
					Prob > F = 0.0000
					R-squared = 0.4294
					Root MSE = .18786

---

Robust						
lpackpc	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
lravgprs	-1.277424	.2496099	-5.12	0.000	-1.780164	-.7746837
lperinc	.2804045	.2538894	1.10	0.275	-.230955	.7917641
_cons	9.894955	.9592169	10.32	0.000	7.962993	11.82692

---

Instrumented: lravgprs

Instruments: lperinc rtaxso rtax    *STATA lists ALL the exogenous regressors as "instruments" - slightly different terminology than we have been using*

---

# *Example: Cigarette demand, two instruments (2 of 2)*

TSLS estimates,  $Z$  = sales tax ( $m = 1$ )

$$\ln(Q_i^{\text{cigarettes}}) = 9.43 - 1.14 \ln(P_i^{\text{cigarettes}}) + 0.21 \ln(\text{Income}_i)$$

(1.26) (0.37) (0.31)

TSLS estimates,  $Z$  = sales tax & cig-only tax ( $m = 2$ )

$$\ln(Q_i^{\text{cigarettes}}) = 9.89 - 1.28 \ln(P_i^{\text{cigarettes}}) + 0.28 \ln(\text{Income}_i)$$

(0.96) (0.25) (0.25)

- Smaller *SEs* for  $m = 2$ . Using 2 instruments gives more information – more “as-if random variation.”
- Low income elasticity (not a luxury good); income elasticity not statistically significantly different from 0
- Surprisingly high price elasticity

# The General Instrument Validity Assumptions

$$Y_i = \beta_0 + \beta_1 X_{1i} + \dots + \beta_k X_{ki} + \beta_{k+1} W_{1i} + \dots + \beta_{k+r} W_{ri} + u_i$$

(1) **Instrument exogeneity**:  $\text{corr}(Z_{1i}, u_i) = 0, \dots, \text{corr}(Z_{mi}, u_i) = 0$

(2) **Instrument relevance**: General case, multiple  $X$ 's

Suppose the second stage regression could be run using the predicted values from the *population* first stage regression. Then: there is no perfect multicollinearity in this (infeasible) second stage regression.

- Special case of one  $X$ : the general assumption is equivalent to (a) at least one instrument must enter the population counterpart of the first stage regression, and (b) the  $W$ 's are not perfectly multicollinear.

# The IV Regression Assumptions

$$Y_i = \beta_0 + \beta_1 X_{1i} + \dots + \beta_k X_{ki} + \beta_{k+1} W_{1i} + \dots + \beta_{k+r} W_{ri} + u_i$$

1.  $E(u_i | W_{1i}, \dots, W_{ri}) = 0$

- #1 says “the exogenous regressors are exogenous.”

2.  $(Y_i, X_{1i}, \dots, X_{ki}, W_{1i}, \dots, W_{ri}, Z_{1i}, \dots, Z_{mi})$  are i.i.d.

- #2 is not new

3. The  $X$ 's,  $W$ 's,  $Z$ 's, and  $Y$  have nonzero, finite 4<sup>th</sup> moments

- #3 is not new

4. The instruments  $(Z_{1i}, \dots, Z_{mi})$  are valid.

- We have discussed this

- Under 1–4, TSLS and its  $t$ -statistic are normally distributed

- The critical requirement is that the instruments be valid

# *W*'s as control variables (1 of 2)

- In many cases, the purpose of including the *W*'s is to control for omitted factors, so that once the *W*'s are included, *Z* is uncorrelated with  $u$ . If so, *W*'s don't need to be exogenous; instead, the *W*'s need to be effective control variables in the sense discussed in Chapter 6 – except now with a focus on producing an exogenous instrument.
- Technically, the condition for *W*'s being effective control variables is that the conditional mean of  $u_i$  does not depend on  $Z_i$ , given  $W_i$ :

$$E(u_i | W_i, Z_i) = E(u_i | W_i)$$

# *W*'s as control variables (2 of 2)

- Thus an alternative to IV regression assumption #1 is that conditional mean independence holds:

$$E(u_i | W_i, Z_i) = E(u_i | W_i)$$

This is the IV version of the conditional mean independence assumption in Chapter 6.

- *Here is the key idea:* in many applications you need to include control variables (*W*'s) so that *Z* is plausibly exogenous (uncorrelated with *u*).
- For the math, see SW Appendix 12.6. For an example, see...

# Example #1: Effect of Studying on Grades (4 of 6)

$$Y_i = \beta_0 + \beta_1 X_i + u_i$$

$Y$  = first-semester GPA

$X$  = average study hours per day

$Z$  = 1 if roommate brought video game, = 0 otherwise

Roommates were randomly assigned

Can you think of a reason that  $Z$  might be correlated with  $u$  – even though it is randomly assigned? What else enters the error term – what are other determinants of grades, beyond time spent studying?

# Example #1: Effect of Studying on Grades (5 of 6)

$$Y_i = \beta_0 + \beta_1 X_i + u_i$$

$$\text{corr}(Z_i, u_i) < 0,$$

Why might  $Z$  be correlated with  $u$ ?

- Here's a hypothetical possibility: the student's sex. Suppose:
  - Roommates are randomly assigned – except always men with men and women with women.
  - Women get better grades than men, holding constant hour spent studying
  - Men are more likely to bring a video game than women
  - Then  $\text{corr}(Z_i, u_i) < 0$  (males are more likely to have a [male] roommate who brings a video game – but males also tend to have lower grades, holding constant the amount of studying).
- Because  $\text{corr}(Z_i, u_i) < 0$  the IV (roommate brings video game) isn't valid.
  - This is the IV version of OV bias.
  - The solution to OV bias is to control for (or include) the OV – in this case, sex.

# Example #1: Effect of Studying on Grades (6 of 6)

- This logic leads you to include  $W$  = student's sex as a control variable in the IV regression:

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 W_i + u_i$$

- The TSLS estimate reported above is from a regression that included gender as a  $W$  variable – along with other variables such as individual  $i$ 's major.
- The conditional mean independence condition for an exogenous instrument is,  $E(u_i | Z_i, W_i) = E(u_i | W_i)$ .
  - In words: among men (conditional on  $W$  = male), roommates are randomly assigned, so whether your roommate brings a video game is random. Same thing among women (conditional on  $W$  = female).
  - The instrument is not exogenous if  $W$  isn't included in the regression.
  - But when  $W$  is included, the conditional mean independence condition  $E(u_i | Z_i, W_i) = E(u_i | W_i)$  holds, and the instrument is valid.

# Checking Instrument Validity (SW Section 12.3)

Recall the two requirements for valid instruments:

1. *Relevance* (special case of one X)

At least one instrument must enter the population counterpart of the first stage regression.

2. *Exogeneity*

**All** the instruments must be uncorrelated with the error term:  $\text{corr}(Z_{1i}, u_i) = 0, \dots, \text{corr}(Z_{mi}, u_i) = 0$

*What happens if one of these requirements isn't satisfied? How can you check? What do you do?*

*If you have multiple instruments, which should you use?*

# Checking Assumption #1: Instrument Relevance

We will focus on a single included endogenous regressor:

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 W_{1i} + \dots + \beta_{1+r} W_{ri} + u_i$$

First stage regression:

$$X_i = \pi_0 + \pi_1 Z_{1i} + \dots + \pi_m Z_{mi} + \pi_{m+1} W_{1i} + \dots + \pi_{m+k} W_{ki} + u_i$$

- The instruments are relevant if at least one of  $\pi_1, \dots, \pi_m$  are nonzero.
- The instruments are said to be **weak** if all the  $\pi_1, \dots, \pi_m$  are either zero or nearly zero.
- **Weak instruments** explain very little of the variation in  $X$ , beyond that explained by the  $W$ 's

# What are the consequences of weak instruments?

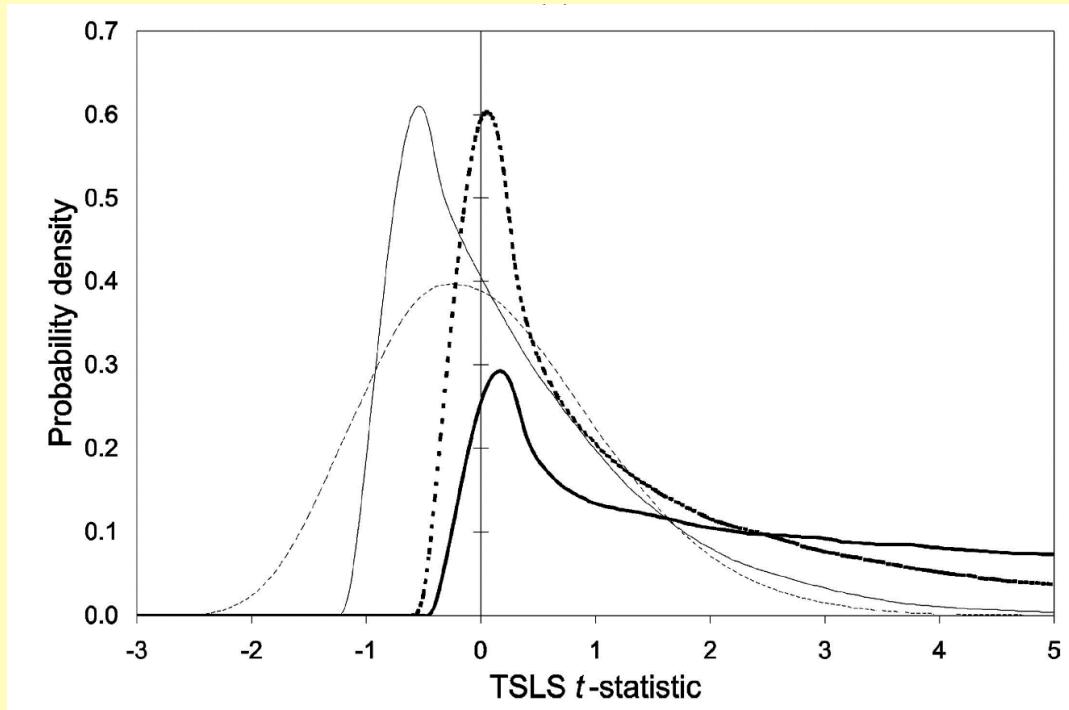
If instruments are weak, the sampling distribution of TSLS and its  $t$ -statistic are not (at all) normal, even with  $n$  large.

Consider the simplest case of 1  $X$ , 1  $Z$ , no control variables:

$$\begin{aligned}Y_i &= \beta_0 + \beta_1 X_i + u_i \\X_i &= \pi_0 + \pi_1 Z_i + u_i\end{aligned}$$

- The IV estimator is  $\hat{\beta}_1^{TSLS} = \frac{s_{YZ}}{s_{XZ}}$
- If cov( $X, Z$ ) is zero or small, then  $s_{XZ}$  will be small: With weak instruments, the denominator is nearly zero.
- If so, the sampling distribution of  $\hat{\beta}_1^{TSLS}$  (and its  $t$ -statistic) is not well approximated by its large- $n$  normal approximation...

# *An example: The sampling distribution of the TSLS $t$ -statistic with weak instruments*



Dark line = irrelevant instruments

Dashed light line = strong instruments

# *Why does our trusty normal approximation fail us?*

$$\hat{\beta}_1^{TSLS} = \frac{s_{YZ}}{s_{XZ}}$$

- If  $\text{cov}(X, Z)$  is small, small changes in  $s_{XZ}$  (from one sample to the next) can induce big changes in  $\hat{\beta}_1^{TSLS}$
- Suppose in one sample you calculate  $s_{XZ} = .00001\dots$
- Thus the large- $n$  normal approximation is a poor approximation to the sampling distribution of  $\hat{\beta}_1^{TSLS}$
- A better approximation is that  $\hat{\beta}_1^{TSLS}$  is distributed as the *ratio* of two correlated normal random variables (see SW App. 12.4)
- If instruments are weak, the usual methods of inference are unreliable – potentially very unreliable!

# Measuring the Strength of Instruments in Practice: The First-Stage $F$ -statistic

The first stage regression (one  $X$ ):

- Regress  $X$  on  $Z_1, \dots, Z_m, W_1, \dots, W_k$ .
- Totally irrelevant instruments  $\leftrightarrow$  all the coefficients on  $Z_1, \dots, Z_m$  are zero.
- The ***first-stage F-statistic*** tests the hypothesis that  $Z_1, \dots, Z_m$  do not enter the first stage regression.
- Weak instruments imply a small first stage  $F$ -statistic.

# Checking for Weak Instruments with a Single $X$ (1 of 2)

- Compute the first-stage  $F$ -statistic.

***Rule-of-thumb: If the first stage  $F$ -statistic is less than 10, then the set of instruments is weak.***

- If so, the TSLS estimator will be biased, and statistical inferences (standard errors, hypothesis tests, confidence intervals) can be misleading.
- *Which  $F$ -statistic?* Use the heteroskedasticity-robust  $F$ , for the usual reasons.

# Checking for Weak Instruments with a Single $X$ (2 of 2)

- Why compare the first-stage  $F$  to 10?
- Simply rejecting the null hypothesis that the coefficients on the  $Z$ 's are zero isn't enough – you need substantial predictive content for the normal approximation to be a good one.
- Comparing the first-stage  $F$  to 10 tests for whether the bias of TSLS, relative to OLS, is less than 10%. If  $F$  is smaller than 10, the relative bias exceeds 10%—that is, TSLS can have substantial bias (see SW App. 12.5).

# What to do if you have weak instruments

- Get better instruments (often easier said than done!)
- If you have many instruments, some are probably weaker than others and it's a good idea to drop the weaker ones (dropping an irrelevant instrument will increase the first-stage  $F$ )
- If you only have a few instruments, and all are weak, then you need to switch from TSLS to a method that is robust to weak instruments.
  - Separate the problem of estimation of  $\beta_1$  and construction of confidence intervals
  - This seems odd, but if TSLS isn't normally distributed, it makes sense (right?)

# Confidence Intervals with Weak Instruments (1 of 2)

- With weak instruments, TSLS confidence intervals are not valid – but some other confidence intervals *are*. Here are two ways to compute confidence intervals that are valid in large samples, even if instruments are weak:

## 1. The Anderson-Rubin confidence interval

- The Anderson-Rubin confidence interval is based on the Anderson-Rubin test statistic testing  $\beta_1 = \beta_{1,0}$ :
  - Compute  $\tilde{u}_i = Y_i - \beta_{1,0} X_i$
  - Regress  $\tilde{u}_i$  on  $W_{1i}, \dots, W_{ri}, Z_{1i}, \dots, Z_{mi}$
  - The AR test is the (heterosk-robust)  $F$ -statistic on  $Z_{1i}, \dots, Z_{mi}$
- Now invert this test: the 95% AR confidence interval is the set of  $\beta_1$  not rejected at the 5% level by the AR test.
- Computation: a pain by hand! use specialized software.

# Confidence Intervals with Weak Instruments (2 of 2)

2. Moreira's Conditional Likelihood Ratio confidence interval
  - This is relevant when you have more than one instrument.
  - The Conditional Likelihood Ratio (CLR) confidence interval is based on inverting Moreira's Conditional Likelihood Ratio test. Computing this test, its critical value, and the CLR confidence interval requires specialized software.
  - The CLR confidence interval tends to be tighter than the Anderson-Rubin confidence interval, especially when there are many instruments.
  - If your software produces the CLR confidence interval, **and if your errors are homoskedastic (or well approximated by homoskedasticity)**, the CLR interval is the one to use.

# Estimation with Weak Instruments

There are no unbiased estimators if instruments are weak or irrelevant. However, some estimators have a distribution more centered around  $\beta_1$  than TSLS.

- One such estimator is the limited information maximum likelihood estimator (LIML)
- The LIML estimator
  - can be derived as a maximum likelihood estimator
  - is the value of  $\beta_1$  that minimizes the  $p$ -value of the AR test(!)
- For more discussion about estimators, tests, and confidence intervals when you have weak instruments, see SW, App. 12.5

# Checking Assumption #2: Instrument Exogeneity

- Instrument exogeneity: **All** the instruments are uncorrelated with the error term:  $\text{corr}(Z_{1i}, u_i) = 0, \dots, \text{corr}(Z_{mi}, u_i) = 0$
- If the instruments are correlated with the error term, the first stage of TSLS cannot isolate a component of  $X$  that is uncorrelated with the error term, so  $\hat{X}$  is correlated with  $u$  and TSLS is inconsistent.
- If there are more instruments than endogenous regressors, it is possible to test – *partially* – for instrument exogeneity.

# Testing Overidentifying Restrictions

Consider the simplest case:

$$Y_i = \beta_0 + \beta_1 X_i + u_i$$

- Suppose there are two valid instruments:  $Z_{1i}, Z_{2i}$
- Then you could compute two separate TSLS estimates.
- Intuitively, if these 2 TSLS estimates are very different from each other, then something must be wrong: one or the other (or both) of the instruments must be invalid.
- The  $J$ -test of overidentifying restrictions makes this comparison in a statistically precise way.
- This can only be done if  $\#Z's > \#X's$  (overidentified).

# The *J*-test of Overidentifying Restrictions (1 of 2)

Suppose # instruments =  $m > \# X\text{'s} = k$  (overidentified)

$$Y_i = \beta_0 + \beta_1 X_{1i} + \dots + \beta_k X_{ki} + \beta_{k+r} W_{1i} + \dots + \beta_{k+r} W_{ri} + u_i$$

The *J*-test is the Anderson-Rubin test, using the TSLS estimator instead of the hypothesized value  $\beta_{1,0}$ . The recipe:

1. First estimate the equation of interest using TSLS and all  $m$  instruments; compute the predicted values  $\hat{Y}_i$ , using the *actual*  $X$ 's (not the  $\hat{X}$ 's used to estimate the second stage)
2. Compute the residuals  $\hat{u}_i = Y_i - \hat{Y}_i$
3. Regress  $\hat{u}_i$  against  $Z_{1i}, \dots, Z_{mi}, W_{1i}, \dots, W_{ri}$
4. Compute the *F*-statistic testing the hypothesis that the coefficients on  $Z_{1i}, \dots, Z_{mi}$  are all zero;
5. The ***J-statistic*** is  $J = mF$

# The $J$ -test of Overidentifying Restrictions (2 of 2)

$J = mF$ , where  $F$  = the  $F$ -statistic testing the coefficients on  $Z_{1i}, \dots, Z_{mi}$  in a regression of the TSLS residuals against  $Z_{1i}, \dots, Z_{mi}, W_{1i}, \dots, W_{ri}$ .

## Distribution of the $J$ -statistic

- Under the null hypothesis that all the instruments are exogenous,  $J$  has a chi-squared distribution with  $m-k$  degrees of freedom
- If  $m = k$ ,  $J = 0$  (*does this make sense?*)
- If some instruments are exogenous and others are endogenous, the  $J$  statistic will be large, and the null hypothesis that all instruments are exogenous will be rejected.

# Checking Instrument Validity: Summary (1 of 2)

This summary considers the case of a single  $X$ . The two requirements for valid instruments are:

## 1. *Relevance*

- At least one instrument must enter the population counterpart of the first stage regression.
- If instruments are weak, then the TSLS estimator is biased and the  $t$ -statistic has a non-normal distribution
- To check for weak instruments with a single included endogenous regressor, check the first-stage  $F$ 
  - If  $F > 10$ , instruments are strong – use TSLS
  - If  $F < 10$ , weak instruments – take some action.

# Checking Instrument Validity: Summary (2 of 2)

## 2. Exogeneity

- **All** the instruments must be uncorrelated with the error term:  $\text{corr}(Z_{1i}, u_i) = 0, \dots, \text{corr}(Z_{mi}, u_i) = 0$
- We can partially test for exogeneity: if  $m > 1$ , we can test the null hypothesis that all the instruments are exogenous, against the alternative that as many as  $m - 1$  are endogenous (correlated with  $u$ )
- The test is the  $J$ -test, which is constructed using the TSLS residuals.
- If the  $J$ -test rejects, then at least some of your instruments are endogenous – so you must make a difficult decision and jettison some (or all) of your instruments.

# Application to the Demand for Cigarettes (SW Section 12.4)

Why are we interested in knowing the elasticity of demand for cigarettes?

- Theory of optimal taxation. The optimal tax rate is inversely related to the price elasticity: the greater the elasticity, the less quantity is affected by a given percentage tax, so the smaller is the change in consumption and deadweight loss.
- Externalities of smoking – role for government intervention to discourage smoking
  - health effects of second-hand smoke? (non-monetary)
  - monetary externalities

# Panel data set

- Annual cigarette consumption, average prices paid by end consumer (including tax), personal income, and tax rates (cigarette-specific and general statewide sales tax rates)
- 48 continental US states, 1985–1995

## Estimation strategy

- We need to use IV estimation methods to handle the simultaneous causality bias that arises from the interaction of supply and demand.
- State binary indicators =  $W$  variables (control variables) which control for unobserved state-level characteristics that affect the demand for cigarettes and the tax rate, as long as those characteristics don't vary over time.

# Fixed-effects model of cigarette demand

$$\ln(Q_{it}^{\text{cigarettes}}) = \alpha_i + \beta_1 \ln(P_{it}^{\text{cigarettes}}) + \beta_2 \ln(\text{Income}_{it}) + u_{it}$$

- $i = 1, \dots, 48, t = 1985, 1986, \dots, 1995$
- $\text{corr}(\ln(P_{it}^{\text{cigarettes}}), u_{it})$  is plausibly nonzero because of supply/demand interactions
- $\alpha_i$  reflects unobserved omitted factors that vary across states but not over time, e.g. attitude towards smoking
- Estimation strategy:
  - Use panel data regression methods to eliminate  $\alpha_i$
  - Use TSLS to handle simultaneous causality bias
  - Use  $T = 2$  with 1985 – 1995 changes (“changes” method) – look at long-term response, not short-term dynamics (short- v. long-run elasticities)

# The “changes” method (when $T=2$ )

- One way to model long-term effects is to consider 10-year changes, between 1985 and 1995
- Rewrite the regression in “changes” form:

$$\begin{aligned}\ln(Q_{i1995}^{\text{cigarettes}}) - \ln(Q_{i1985}^{\text{cigarettes}}) &= \beta_1[\ln(P_{i1995}^{\text{cigarettes}}) - \ln(P_{i1985}^{\text{cigarettes}})] \\ &\quad + \beta_2[\ln(\text{Income}_{i1995}) - \ln(\text{Income}_{i1985})] + (u_{i1995} - u_{i1985})\end{aligned}$$

- Create “10-year change” variables, for example:
- 10-year change in log price =  $\ln(P_{i1995}) - \ln(P_{i1985})$
- Then estimate the demand elasticity by TSLS using 10-year changes in the instrumental variables
- This is equivalent to using the original data and including the state binary indicators (“ $W$ ” variables) in the regression

# STATA: Cigarette demand

First create “10-year change” variables

10-year change in log price

$$= \ln(P_{it}) - \ln(P_{it-10}) = \ln(P_{it} / P_{it-10})$$

```
. gen dlpackpc = log(packpc/packpc[_n-10])          _n-10 is the 10-yr lagged value
. gen dlavgprs = log(avgprs/avgprs[_n-10])
. gen dlperinc = log(perinc/perinc[_n-10])
. gen drtaxs  = rtaxs-rtaxs[_n-10]
. gen drtax   = rtax-rtax[_n-10]
. gen drtaxso = rtaxso-rtaxso[_n-10]
```

# Use TSLS to estimate the demand elasticity by using the “10-year changes” specification

	Y	W	X	Z	
. ivregress 2sls	dlpackpc	dlperinc	(dlavgprs = drtaxso)	,	r
IV (2SLS) regression with robust standard errors					Number of obs = 48
					F( 2, 45) = 12.31
					Prob > F = 0.0001
					R-squared = 0.5499
					Root MSE = .09092

---

Robust						
dlpackpc	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
<hr/>						
dlavgprs	-.9380143	.2075022	-4.52	0.000	-1.355945	-.5200834
dlperinc	.5259693	.3394942	1.55	0.128	-.1578071	1.209746
_cons	.2085492	.1302294	1.60	0.116	-.0537463	.4708446

---

Instrumented: dlavgprs

Instruments: dlperinc drtaxso

---

## NOTE:

- All the variables - Y, X, W, and Z's - are in 10-year changes
- Estimated elasticity = -.94 (SE = .21) - surprisingly elastic!
- Income elasticity small, not statistically different from zero
- Must check whether the instrument is relevant...

# Check instrument relevance: compute first-stage $F$

```
. reg dlavgprs drtaxso dlperinc
```

Source	SS	df	MS	Number of obs	=	48
Model	.191437213	2	.095718606	F( 2, 45)	=	23.86
Residual	.180549989	45	.004012222	Prob > F	=	0.0000
Total	.371987202	47	.007914621	R-squared	=	0.5146
				Adj R-squared	=	0.4931
				Root MSE	=	.06334
dlavgprs	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
drtaxso	.0254611	.0037374	6.81	0.000	.0179337	.0329885
dlperinc	-.2241037	.2119405	-1.06	0.296	-.6509738	.2027664
_cons	.5321948	.031249	17.03	0.000	.4692561	.5951334

```
. test drtaxso
```

```
( 1) drtaxso = 0
```

We didn't need to run "test" here!

With  $m=1$  instrument, the  $F$ -stat is  
the square of the  $t$ -stat:

$$6.81 \times 6.81 = 46.41$$

First stage  $F = 46.5 > 10$  so instrument is not weak

Can we check instrument exogeneity? **No**:  $m = k$

# Cigarette demand, 10 year changes – 2 IVs

	<i>Y</i>	<i>W</i>	<i>X</i>	<i>Z1</i>	<i>Z2</i>	
.	ivregress 2sls dlp packpc dlperinc (dlavgprs = drtaxso drtax) , vce(r)					
Instrumental variables (2SLS) regression						Number of obs = 48
						Wald chi2(2) = 45.44
						Prob > chi2 = 0.0000
						R-squared = 0.5466
						Root MSE = .08836

---

	Robust					
dlp packpc	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
dlavgprs	-1.202403	.1906896	-6.31	0.000	-1.576148	-.8286588
dlperinc	.4620299	.2995177	1.54	0.123	-.1250139	1.049074
_cons	.3665388	.1180414	3.11	0.002	.1351819	.5978957

---

Instrumented: dlavgprs

Instruments: dlperinc drtaxso drtax

---

drtaxso = general sales tax only

drtax = cigarette-specific tax only

Estimated elasticity is -1.2, even more elastic than using general sales tax only!

# First-stage F – both instruments

```
X      Z1      Z2      W
. reg dlavgprs drtaxso drtax dlperinc

Source |       SS          df         MS
-----+-----+-----+
Model | .289359873      3   .096453291
Residual | .082627329    44   .001877894
-----+-----+
Total | .371987202     47   .007914621

Number of obs =      48
F(  3,     44) =    51.36
Prob > F      = 0.0000
R-squared      = 0.7779
Adj R-squared = 0.7627
Root MSE       = .04333

-----+-----+
dlavgprs |   Coef.    Std. Err.      t    P>|t| [95% Conf. Interval]
-----+-----+
drtaxso |   .013457   .0030498    4.41  0.000   .0073106   .0196033
drtax |   .0075734   .0010488    7.22  0.000   .0054597   .009687
dlperinc |  -.0289943   .1474923   -0.20  0.845  -.3262455   .2682568
_cons |   .4919733   .0220923   22.27  0.000   .4474492   .5364973

. test drtaxso drtax
( 1) drtaxso = 0
( 2) drtax = 0
F(  2,     44) =    75.65      75.65 > 10 so instruments aren't weak
Prob > F = 0.0000
```

With  $m > k$ , we can test the overidentifying restrictions...

# Test the overidentifying restrictions (1 of 2)

- . predict e, resid *Computes predicted values for most recently estimated regression (the previous TSLS regression)*
- . reg e drtaxso drtax dlperinc *Regress e on Z's and W's*

Source	SS	df	MS	Number of obs =	48
Model	.037769176	3	.012589725	F( 3, 44) =	1.64
Residual	.336952289	44	.007658007	Prob > F =	0.1929
Total	.374721465	47	.007972797	R-squared =	0.1008
				Adj R-squared =	0.0395
				Root MSE =	.08751

e	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
drtaxso	.0127669	.0061587	2.07	0.044	.000355 .0251789
drtax	-.0038077	.0021179	-1.80	0.079	-.008076 .0004607
dlperinc	-.0934062	.2978459	-0.31	0.755	-.6936752 .5068627
_cons	.002939	.0446131	0.07	0.948	-.0869728 .0928509

- . test drtaxso drtax
    - ( 1) drtaxso = 0 *Compute J-statistic, which is m\*F, where F tests whether coefficients on the instruments are zero*
    - ( 2) drtax = 0
- $F( 2, 44) = 2.47$  so  $J = 2 \cdot 2.47 = 4.93$
- Prob > F = 0.0966 *\*\* WARNING - this uses the wrong d.f. \*\**

# Test the overidentifying restrictions (2 of 2)

The correct degrees of freedom for the  $J$ -statistic is  $m-k$ :

- $J = mF$ , where  $F$  = the  $F$ -statistic testing the coefficients on  $Z_{1i}, \dots, Z_{mi}$  in a regression of the TSLS residuals against  $Z_{1i}, \dots, Z_{mi}, W_{1i}, \dots, W_{mi}$ .
- Under the null hypothesis that all the instruments are exogenous,  $J$  has a chi-squared distribution with  $m-k$  degrees of freedom
- Here,  $J = 4.93$ , distributed chi-squared with d.f. = 1; the 5% critical value is 3.84, so reject at 5% sig. level.
- In STATA:

```
. dis "J-stat = " r(df)*r(F) " p-value = " chiprob(r(df)-1,r(df)*r(F))  
J-stat = 4.9319853 p-value = .02636401  
  
J = 2 * 2.47 = 4.93           p-value from chi-squared(1) distribution
```

*Now what???*

# Tabular summary of these results

**TABLE 12.1** Two Stage Least Squares Estimates of the Demand for Cigarettes Using Panel Data for 48 U.S. States

Dependent variable: $\ln(Q_{i,1995}^{\text{cigarettes}}) - \ln(Q_{i,1985}^{\text{cigarettes}})$			
Regressor	(1)	(2)	(3)
$\ln(P_{i,1995}^{\text{cigarettes}}) - \ln(P_{i,1985}^{\text{cigarettes}})$	-0.94 (0.21) [-1.36, -0.52]	-1.34 (0.23) [-1.80, -0.88]	-1.20 (0.20) [-1.60, -0.81]
$\ln(Inc_{i,1995}) - \ln(Inc_{i,1985})$	0.53 (0.34) [-0.16, 1.21]	0.43 (0.30) [-0.16, 1.02]	0.46 (0.31) [-0.16, 1.09]
Intercept	-0.12 (0.07)	-0.02 (0.07)	-0.05 (0.06)
Instrumental variable(s)	Sales tax	Cigarette-specific tax	Both sales tax and cigarette-specific tax
First-stage <i>F</i> -statistic	33.7	107.2	88.6
Overidentifying restrictions <i>J</i> -test and <i>p</i> -value	—	—	4.93 (0.026)

These regressions were estimated using data for 48 U.S. states (48 observations on the 10-year differences). The data are described in Appendix 12.1. The *J*-test of overidentifying restrictions is described in Key Concept 12.6 (its *p*-value is given in parentheses), and the first-stage *F*-statistic is described in Key Concept 12.5. Heteroskedasticity-robust standard errors are given in parentheses beneath coefficients, and 95% confidence intervals are given in brackets.

# How should we interpret the *J*-test rejection?

- *J*-test rejects the null hypothesis that both the instruments are exogenous
- This means that either *rtaxso* is endogenous, or *rtax* is endogenous, or both!
- The *J*-test doesn't tell us which! *You must exercise judgment...*
- Why might *rtax* (cig-only tax) be endogenous?
  - Political forces: history of smoking or lots of smokers ↗ political pressure for low cigarette taxes
  - If so, cig-only tax is endogenous
- This reasoning doesn't apply to general sales tax
- → use just one instrument, the general sales tax

# The Demand for Cigarettes: Summary of Empirical Results

- Use the estimated elasticity based on TSLS with the general sales tax as the only instrument:  
 $\text{Elasticity} = -.94, SE = .21$
- This elasticity is surprisingly large (not inelastic) – a 1% increase in prices reduces cigarette sales by nearly 1%. This is much more elastic than conventional wisdom in the health economics literature.
- This is a long-run (ten-year change) elasticity. *What would you expect a short-run (one-year change) elasticity to be – more or less elastic?*

# Assess the Validity of the Study (1 of 2)

## Remaining threats to internal validity?

1. Omitted variable bias?
  - *The fixed effects estimator controls for unobserved factors that vary across states but not over time*
2. Functional form mis-specification? (*could check this*)
3. Remaining simultaneous causality bias?
  - *Not if the general sales tax a valid instrument, once state fixed effects are included!*
4. Errors-in-variables bias?
5. Selection bias? (*no, we have all the states*)
6. An additional threat to internal validity of IV regression studies is whether the instrument is (1) relevant and (2) exogenous. *How significant are these threats in the cigarette elasticity application?*

# Assess the Validity of the Study (2 of 2)

## External validity?

- We have estimated a long-run elasticity – can it be generalized to a short-run elasticity? Why or why not?
- Suppose we want to use the estimated elasticity of  $-0.94$  to guide policy today. Here are two changes since the period covered by the data (1985–95) – do these changes pose a threat to external validity (generalization from 1985–95 to today)?
  - Levels of smoking today are lower than in 1985–1995
  - Cultural attitudes toward smoking have changed against smoking since 1985–95.

# Where Do Valid Instruments Come From? (SW Section 12.5)

## General comments

The hard part of IV analysis is finding valid instruments

- Method #1: “variables in another equation” (e.g. supply shifters that do not affect demand)
- Method #2: look for exogenous variation ( $Z$ ) that is “as if” randomly assigned (does not directly affect  $Y$ ) but affects  $X$ .
- These two methods are different ways to think about the same issues – see the link...
  - Rainfall shifts the supply curve for butter but not the demand curve; rainfall is “as if” randomly assigned
  - Sales tax shifts the supply curve for cigarettes but not the demand curve; sales taxes are “as if” randomly assigned

# Example: Cardiac Catheterization (1 of 3)

McClellan, Mark, Barbara J. McNeil, and Joseph P. Newhouse (1994), “Does More Intensive Treatment of Acute Myocardial Infarction in the Elderly Reduce Mortality?” *Journal of the American Medical Association*, vol. 272, no. 11, pp. 859 – 866.

Does cardiac catheterization improve longevity of heart attack patients?

$Y_i$  = survival time (in days) of heart attack patient

$X_i$  = 1 if patient receives cardiac catheterization,  
= 0 otherwise

- Clinical trials show that *CardCath* affects *SurvivalDays*.
- But is the treatment effective “in the field”?

# Example: Cardiac Catheterization (2 of 3)

$$SurvivalDays_i = \beta_0 + \beta_1 CardCath_i + u_i$$

- Is OLS unbiased? The decision to treat a patient by cardiac catheterization is endogenous – it is (was) made in the field by EMT technician and depends on  $u_i$  (unobserved patient health characteristics)
- If healthier patients are catheterized, then OLS has simultaneous causality bias and OLS overstates overestimates the CC effect
- Propose instrument: distance to the nearest CC hospital minus distance to the nearest “regular” hospital

# Example: Cardiac Catheterization (3 of 3)

- $Z$  = differential distance to CC hospital
  - Relevant? If a CC hospital is far away, patient won't be taken there and won't get CC
  - Exogenous? If distance to CC hospital doesn't affect survival, other than through effect on  $CardCath_i$ , then  $\text{corr}(\text{distance}, u_i) = 0$  so exogenous
  - If patients location is random, then differential distance is "as if" randomly assigned.
  - *The 1<sup>st</sup> stage is a linear probability model: distance affects the probability of receiving treatment*
- Results:
  - OLS estimates significant and large effect of CC
  - TSLS estimates a small, often insignificant effect

# Example: Crowding Out of Private Charitable Spending (1 of 4)

Gruber, Jonathan and Daniel M. Hungerman (2005),  
“Faith-Based Charity and Crowd Out During the Great Depression,” NBER Working Paper 11332.

Does government social service spending crowd out private (church, Red Cross, etc.) charitable spending?

$Y$  = private charitable spending (churches)

$X$  = government spending

What is the motivation for using instrumental variables?

Proposed instrument:

$Z$  = strength of Congressional delegation

# Example: Crowding Out of Private Charitable Spending (2 of 4)

- panel data, yearly, by state, 1929–1939, U.S.
- $Y$  = total benevolent spending by six church denominations (CCC, Lutheran, Northern Baptist, Presbyterian (2), Southern Baptist); benevolences =  $\frac{1}{4}$  of total church expenditures.
- $X$  = Federal relief spending under New Deal legislation (General Relief, Work Relief, Civil Works Administration, Aid to Dependent Children,...)
- $Z$  = tenure of state's representatives on House & Senate Appropriations Committees, in months
- $W$  = lots of fixed effects

# Example: Crowding Out of Private Charitable Spending (3 of 4)

Figure 1: Government and Church Relief during the Great Depression



# Example: Crowding Out of Private Charitable Spending (4 of 4)

## ***Assessment of validity:***

- Instrument validity:
  - Relevance?
  - Exogeneity?
- Other threats to internal validity:
  1. OV bias
  2. Functional form
  3. Measurement error
  4. Selection
  5. Simultaneous causality
- External validity to today in U.S.? To aid to developing countries?

# Example: School Competition (1 of 3)

Hoxby, Caroline M. (2000), “Does Competition Among Public Schools Benefit Students and Taxpayers?”  
*American Economic Review* 90, 1209–1238

What is the effect of public school competition on student performance?

$Y$  = 12<sup>th</sup> grade test scores

$X$  = measure of choice among school districts (function of # of districts in metro area)

What is the motivation for using instrumental variables?

Proposed instrument:

$Z$  = # small streams in metro area

# Example: School Competition (2 of 3)

## Data – some details

- cross-section, US, metropolitan area, late 1990s ( $n = 316$ ),
- $Y = 12^{\text{th}}$  grade reading score (other measures too)
- $X =$  index taken from industrial organization literature measuring the amount of competition (“Gini index”) – based on number of “firms” and their “market share”
- $Z =$  measure of small streams – which formed natural geographic boundaries.
- $W =$  lots of control variables

# Example: School Competition (3 of 3)

## ***Assessment of validity:***

- Instrument validity:
  - Relevance?
  - Exogeneity?
- Other threats to internal validity:
  1. OV bias
  2. Functional form
  3. Measurement error
  4. Selection
  5. Simultaneous causality
- External validity to today in U.S.? To aid to developing countries?

# Conclusion (SW Section 12.6)

- A valid instrument lets us isolate a part of  $X$  that is uncorrelated with  $u$ , and that part can be used to estimate the effect of a change in  $X$  on  $Y$
- IV regression hinges on having valid instruments:
  1. *Relevance*: Check via first-stage  $F$
  2. *Exogeneity*: Test overidentifying restrictions via the  $J$ -statistic
- A valid instrument isolates variation in  $X$  that is “as if” randomly assigned.
- The critical requirement of at least  $m$  valid instruments cannot be tested – you must use your head.

# Some IV FAQs (1 of 2)

## 1. When might I want to use IV regression?

Any time that  $X$  is correlated with  $u$  and you have a valid instrument. The primary reasons for correlation between  $X$  and  $u$  could be:

- Omitted variable(s) that lead to OV bias
  - Ex: ability bias in returns to education
- Measurement error
  - Ex: measurement error in years of education
- Selection bias
  - Patients select treatment
- Simultaneous causality bias
  - Ex: supply and demand for butter, cigarettes

# Some IV FAQs (2 of 2)

## 2. What are the threats to the internal validity of an IV regression?

- The main threat to the internal validity of IV is the failure of the assumption of valid instruments. Given a set of control variables  $W$ , instruments are valid if they are relevant and exogenous.
  - Instrument relevance can be assessed by checking if instruments are weak or strong: Is the first-stage  $F$ -statistic  $> 10$ ?
  - Instrument exogeneity can be checked using the  $J$ -statistic – as long as you have  $m$  exogenous instruments to start with! In general, instrument exogeneity must be assessed using expert knowledge of the application.