

Sommario

1. Regressione IV: cosa e perché; minimi quadrati in due stadi
2. Il modello generale di regressione IV
3. Verifica della validità degli strumenti
 - a) Strumenti deboli e forti
 - b) Esogeneità degli strumenti
4. Applicazione: domanda di sigarette
5. Esempi: dove trovare gli strumenti?

Regressione IV: perché?

Tre importanti minacce alla validità interna sono:

- Distorsione da variabili omesse per una variabile correlata con X ma inosservata (perciò non può essere inclusa nella regressione) e per cui vi sono variabili di controllo inadeguate;
- Distorsione da causalità simultanea (X causa Y , Y causa X);
- Distorsione da errori nelle variabili (X è misurata con errore)

Tutti e tre i problemi comportano $E(u|X) \neq 0$.

- La regressione con variabili strumentali può eliminare la distorsione quando $E(u|X) \neq 0$ – usando una *variabile strumentale* (IV), Z .

Lo stimatore IV con un singolo regressore e un singolo strumento (Paragrafo 12.1)

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

- La regressione IV divide X in due parti: una che potrebbe essere correlata con u , e una che non lo è. Isolando la parte che non è correlata con u , è possibile stimare β_1 .
- Per fare questo si utilizza una **variabile strumentale**, Z_i , che è correlata con X_i ma incorrelata con u_i .

Terminologia: endogeneità ed esogeneità

Una variabile **endogena** è una variabile correlata con u

Una variabile **esogena** è una variabile incorrelata con u

Nella regressione IV ci concentriamo sul caso in cui X è endogena ed esiste uno strumento, Z , esogeno.

Digressione sulla terminologia: “endogeno” significa letteralmente “determinato all’interno del sistema”. Se X è congiuntamente determinata con Y , allora una regressione di Y su X è soggetta a distorsione da causalità simultanea. Ma questa definizione di endogeneità è troppo stretta perché sia possibile usare la regressione IV per risolvere i problemi di distorsione da variabili omesse e da errori nelle variabili, quindi usiamo la definizione più ampia fornita sopra.

Due condizioni per avere uno strumento valido

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

Perché una variabile strumentale (uno "**strumento**") Z sia valida, deve soddisfare due condizioni:

1. Rilevanza: $\text{corr}(Z_i, X_i) \neq 0$

2. Esogeneità: $\text{corr}(Z_i, u_i) = 0$

Supponiamo per ora di avere un tale Z_i (vedremo più avanti come trovare variabili strumentali); come possiamo usarlo per stimare β_1 ?

Lo stimatore IV con una X e una Z

Spiegazione 1: minimi quadrati in due stadi (TOLS)

Ci sono due stadi – due regressioni:

(1) Si isola la parte di X che non è correlata con u mediante la regressione di X su Z usando gli OLS:

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

- Poiché Z_i non è correlato con u_i , $\pi_0 + \pi_1 Z_i$ non è correlato con u_i . Non conosciamo π_0 o π_1 ma li abbiamo stimati, perciò...
- Si calcolano i valori predetti di X_i , \hat{X}_i , dove $\hat{X}_i = \hat{\pi}_0 + \hat{\pi}_1 Z_i$, $i = 1, \dots, n$.

Minimi quadrati in due stadi (continua)

(2) Si sostituisce X_i con \hat{X}_i nella regressione di interesse: si esegue la regressione di Y su \hat{X}_i usando gli OLS:

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

- **Poiché \hat{X}_i è incorrelato con u_i , la prima assunzione dei minimi quadrati vale per la regressione (2).** (Ciò richiede che n sia grande in modo che π_0 e π_1 siano stimati con precisione)
- Quindi, in grandi campioni, β_1 può essere stimato con gli OLS usando la regressione (2)
- Lo stimatore risultante è detto *stimatore dei minimi quadrati in due stadi (TSLS)*, $\hat{\beta}_1^{TSLS}$.

Minimi quadrati in due stadi: riepilogo

Supponiamo che Z_i , soddisfi le due condizioni per uno strumento valido:

1. Rilevanza: $\text{corr}(Z_i, X_i) \neq 0$

2. Esogeneità: $\text{corr}(Z_i, u_i) = 0$

Minimi quadrati in due stadi:

Stadio 1: Regressione di X_i su Z_i (inclusa intercetta), ottenendo i valori predetti \hat{X}_i

Stadio 2: Regressione di Y_i su \hat{X}_i (inclusa intercetta); il coefficiente di \hat{X}_i è lo stimatore TSLS, $\hat{\beta}_1^{TSLS}$.

$\hat{\beta}_1^{TSLS}$ è uno stimatore consistente di β_1 .

Lo stimatore IV, una X e una Z (continua)

Spiegazione 2: derivazione algebrica diretta

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

Allora

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

dove $\text{cov}(u_i, Z_i) = 0$ per l'esogeneità dello strumento; quindi

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

Lo stimatore IV, una X e una Z (continua)

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

Lo stimatore IV sostituisce queste covarianze della popolazione con covarianze campionarie:

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

s_{YZ} e s_{XZ} sono covarianze campionarie. Questo è lo stimatore TOLS – con una derivazione diversa!

Lo stimatore IV, una X e una Z (continua)

Spiegazione 3: derivazione dalla "forma ridotta"

La "forma ridotta" mette in relazione Y a Z e X a Z:

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

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

dove w_i è un termine d'errore. Poiché Z è esogena, è incorrelata con v_i e con w_i .

L'idea: una variazione unitaria in Z_i comporta una variazione in X_i di π_1 e una variazione in Y_i di γ_1 . Poiché tale variazione in X_i nasce dalla variazione esogena in Z_i , tale variazione in X_i è esogena. Quindi una variazione esogena in X_i di π_1 unità è associata a una variazione in Y_i di γ_1 unità – perciò l'effetto su Y di una variazione esogena in X è $\beta_1 = \gamma_1 / \pi_1$ unità.

Lo stimatore IV dalla forma ridotta (continua)

I calcoli:

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

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

Risolviamo l'equazione di X in Z :

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

Sostituiamo nell'equazione di Y e raccogliamo i termini:

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

$$= \gamma_0 + \gamma_1 [-\pi_0/\pi_1 + (1/\pi_1)X_i - (1/\pi_1)v_i] + w_i$$

$$= [\gamma_0 - \pi_0\gamma_1/\pi_1] + (\gamma_1/\pi_1)X_i + [w_i - (\gamma_1/\pi_1)v_i]$$

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

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

Lo stimatore IV dalla forma ridotta (continua)

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

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

quindi

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

dove

$$\beta_1 = \gamma_1 / \pi_1$$

Interpretazione: una variazione esogena in X_i di π_1 unità è associata a una variazione in Y_i di γ_1 unità – perciò l'effetto su Y di una variazione unitaria esogena in X è $\beta_1 = \gamma_1 / \pi_1$.

Esempio 1: effetto dello studio sui voti

Qual è l'effetto sui voti di studiare un'ora in più al giorno?

Y = media voti

X = tempo di studio (ore al giorno)

Dati: voti e ore di studio di studenti del primo anno di college.

Vi aspettate che lo stimatore OLS di β_1 (l'effetto sulla media voti di studiare un'ora in più al giorno) sia non distorto? Perché, o perché no?

Effetto dello studio sui voti (continua)

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.

- $n = 210$ studenti del primo anno al Berea College (Kentucky) nel 2001
- Y = media voti del primo semestre
- X = media ore di studi al giorno (sondaggio)
- I compagni di stanza sono stati assegnati a caso
- $Z = 1$ se il compagno di stanza ha portato un videogioco, = 0 altrimenti

Pensate che Z_i (indica se un compagno ha portato un videogioco) sia uno strumento valido?

1. È rilevante (correlato con X)?
2. È esogeno (incorrelato con u)?

Effetto dello studio sui voti (continua)

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

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

$Y = \text{media voti (scala 4 punti)}$

$X = \text{tempo di studio (ore al giorno)}$

$Z = 1 \text{ se il compagno ha portato un videogioco, } = 0 \text{ altrimenti}$

Risultati di Stinebrinckner e Stinebrinckner

$$\hat{\pi}_1 = -0,668$$

$$\hat{\gamma}_1 = -0,241$$

$$\hat{\beta}_1^{IV} = \frac{\hat{\gamma}_1}{\hat{\pi}_1} = \frac{-0,241}{-0,668} = 0,360$$

Quali sono le unità? Queste stime hanno senso nel mondo reale? (Nota: in realtà hanno eseguito le regressioni con regressori aggiuntivi – ci torneremo più avanti)

Consistenza dello stimatore TSLS

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

Le covarianze campionarie sono consistenti:

$s_{YZ} \xrightarrow{p} \text{cov}(Y,Z)$ e $s_{XZ} \xrightarrow{p} \text{cov}(X,Z)$. Quindi

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

- La condizione di rilevanza dello strumento, $\text{cov}(X,Z) \neq 0$, assicura che non si esegua una divisione per zero.

Esempio 2: offerta e domanda di burro

La regressione IV è stata sviluppata in origine per stimare l'elasticità della domanda per beni agricoli, per esempio il burro:

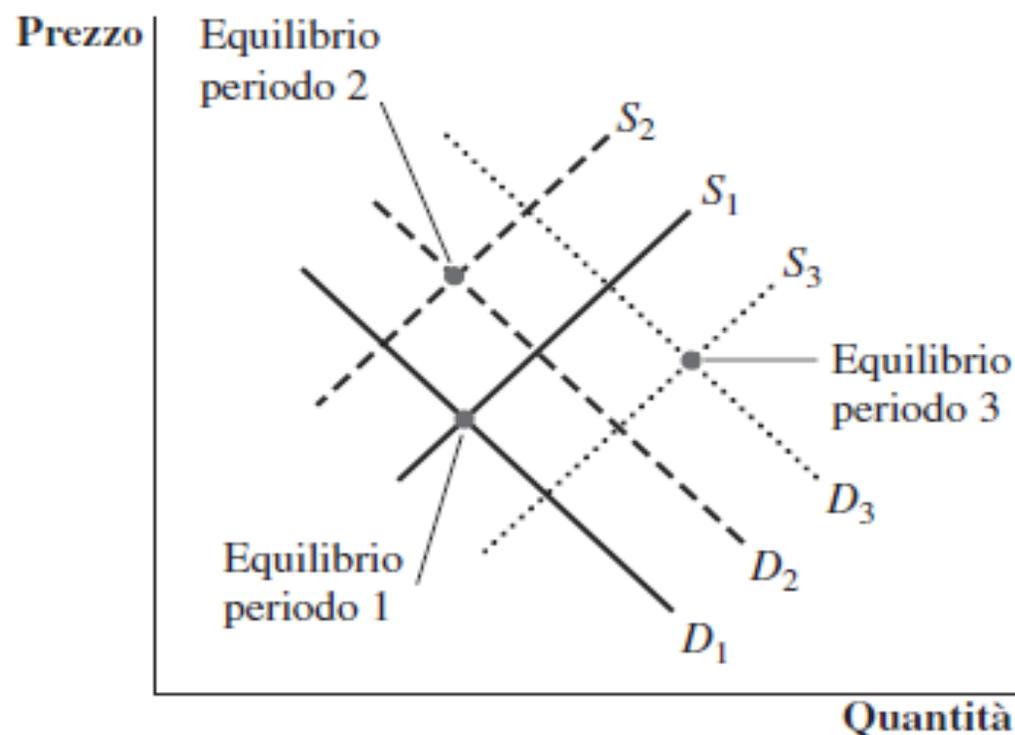
$$\ln(Q_i^{butter}) = \beta_0 + \beta_1 \ln(P_i^{butter}) + u_i$$

• β_1 = elasticità del burro = variazione percentuale in quantità per una variazione dell'1% in prezzo (si ricordi la discussione sulla specifica log-log)

• Dati: osservazioni su prezzo e quantità di burro per diversi anni

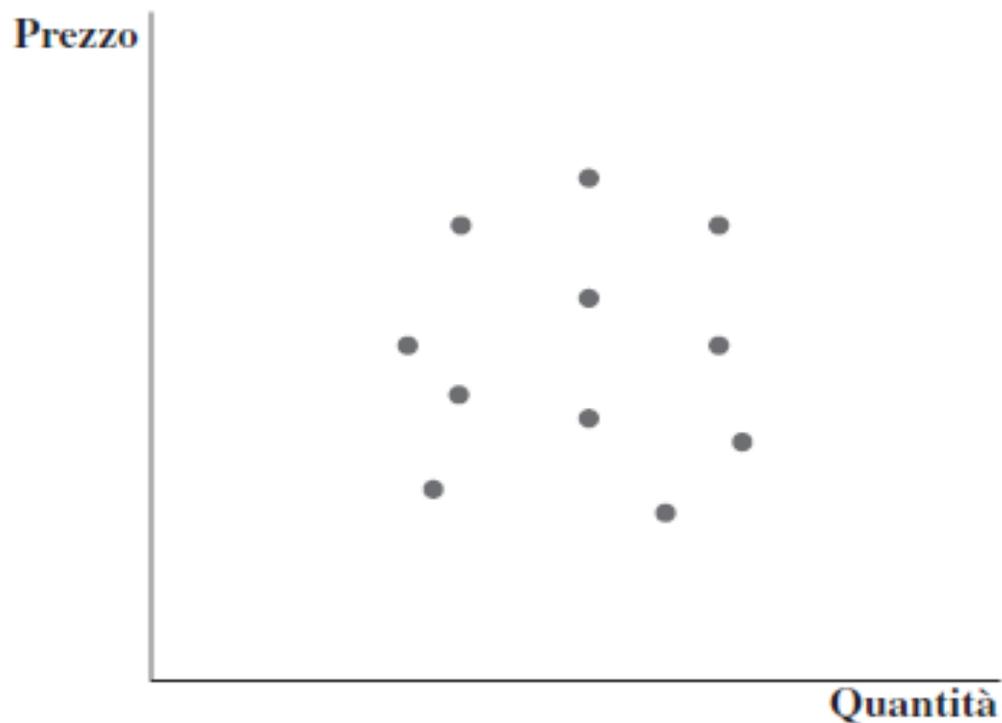
• La regressione OLS di $\ln(Q_i^{butter})$ su $\ln(P_i^{butter})$ soffre di distorsione da causalità simultanea (*perché?*)

La distorsione da causalità simultanea nella regressione OLS di $\ln(\text{Prezzo})$ su $\ln(Q_i^{butter})$ nasce perché prezzo e quantità sono determinati dall'interazione di domanda e offerta:



(a) Domanda e offerta in tre periodi

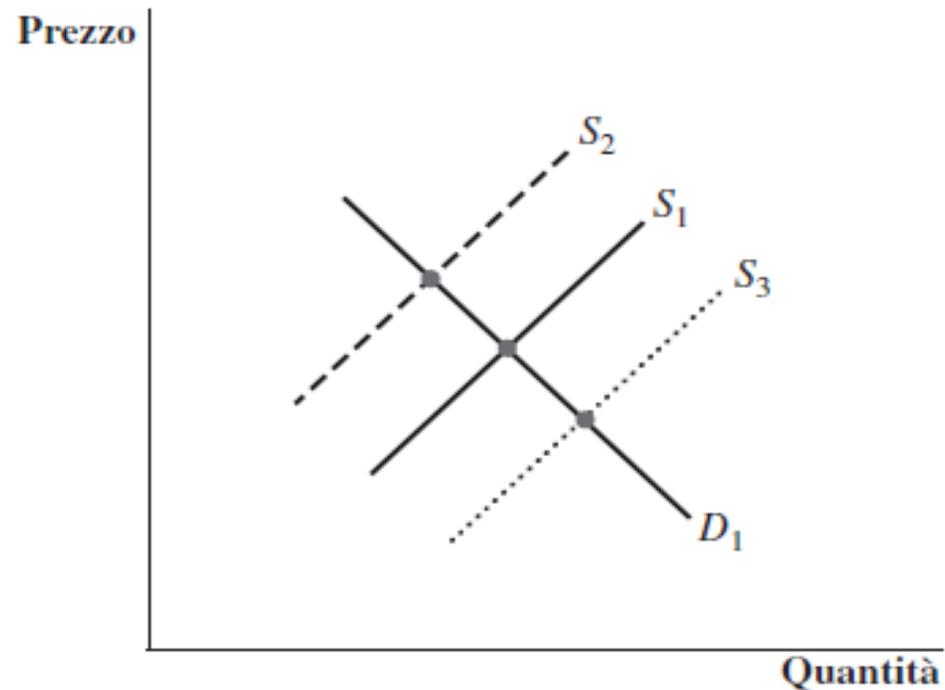
Questa interazione di domanda e offerta produce dati come...



(b) Prezzo e quantità di equilibrio per 11 periodi

Una regressione con questi dati produrrebbe la curva di domanda?

E se si spostasse solo l'offerta?



(c) Prezzo e quantità di equilibrio
quando solo la curva di offerta si sposta

- TSLS stima la curva di domanda isolando gli spostamenti di prezzo e quantità conseguenti a spostamenti dell'offerta.
- Z è una variabile che sposta l'offerta ma non la domanda.

TOLS nell'esempio di domanda e offerta:

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

Sia Z = pioggia nelle aree di produzione lattiera.

Z è uno strumento valido?

(1) Rilevante? $\text{corr}(rain_i, \ln(P_i^{butter})) \neq 0$?

Plausibilmente: pioggia insufficiente significa meno pascolo quindi meno burro e quindi prezzi più alti

(2) Esogeno? $\text{corr}(rain_i, u_i) = 0$?

Plausibilmente: la pioggia nelle aree di produzione lattiera non dovrebbe influenzare la domanda di burro

TSLS nell'esempio di domanda e offerta (continua)

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

$Z_i = rain_i$ = pioggia nelle aree di produzione lattiera.

Passo 1: regressione di $\ln(P_i^{butter})$ su $rain$, dà $\widehat{\ln(P_i^{butter})}$

$\widehat{\ln(P_i^{butter})}$ isola le variazioni nel log del prezzo conseguenti all'offerta (o almeno a parte di essa)

Passo 2: regressione di $\ln(Q_i^{butter})$ su $\widehat{\ln(P_i^{butter})}$

Controparte dell'uso degli spostamenti della curva di offerta per tracciare la curva di domanda.

Esempio 3: punteggi nei test e dimensioni delle classi

- Le regressioni per punteggi nei test/dimensioni delle classi in California potrebbero avere distorsione da variabili omesse (per esempio per interessamento dei genitori).
- In linea di principio questa distorsione può essere eliminata dalla regressione IV (TSLS).
- La regressione IV richiede uno strumento valido, cioè che sia:

1. rilevante: $\text{corr}(Z_i, STR_i) \neq 0$

2. esogeno: $\text{corr}(Z_i, u_i) = 0$

Esempio 3: punteggi nei test e dimensioni delle classi (continua)

Ecco uno strumento ipotetico:

- alcuni distretti, colpiti casualmente da un terremoto, “raddoppiano” le classi:

$$Z_i = Quake_i = 1 \text{ se colpito da terremoto, } = 0 \text{ altrimenti}$$

- *Valgono le due condizioni per la validità dello strumento?*
- Il terremoto crea una situazione *come se* i distretti rientrassero in un esperimento con assegnazione casuale. Quindi, la variazione in *STR* conseguente al terremoto è esogena.
- Il primo stadio del TSLS prevede la regressione di *STR* su *Quake*, isolando così la parte esogena di *STR* (la parte “come se” fosse assegnata casualmente)

Inferenza con TSLS

- In grandi campioni, la distribuzione campionaria dello stimatore TSLS è normale
- L'inferenza (verifiche di ipotesi, intervalli di confidenza) procede nel modo consueto, ovvero $\pm 1,96SE$
- Il concetto alla base della distribuzione normale in grandi campioni dello stimatore TSLS è che – come tutti gli altri stimatori che abbiamo considerato – comporta variabili casuali i.i.d. con media nulla, a cui possiamo applicare il TLC.
- Di seguito riportiamo i calcoli abbozzati (si veda l'Appendice 12.3 per i dettagli)...

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

$$= \frac{\sum_{i=1}^n Y_i (Z_i - \bar{Z})}{\sum_{i=1}^n X_i (Z_i - \bar{Z})}$$

Sostituiamo in $Y_i = \beta_0 + \beta_1 X_i + u_i$ e semplifichiamo:

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

quindi...

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

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

Moltiplicando entrambi i membri per \sqrt{n}

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

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

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

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

quindi: $\hat{\beta}_1^{TOLS}$ ha distribuzione appr. $N(\beta_1, \sigma_{\hat{\beta}_1^{TOLS}}^2)$

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

dove $\text{cov}(X, Z) \neq 0$ perché lo strumento è rilevante

Inferenza con TSLS (continua)

$\hat{\beta}_1^{TSLS}$ ha distribuzione appr. $N(\beta_1, \sigma_{\hat{\beta}_1}^{2,TSLS})$

- L'inferenza statistica procede nel modo consueto.
- La giustificazione è (come di consueto) basata su grandi campioni
- Tutto questo assume che gli strumenti siano validi – vedremo tra breve che cosa accade se non lo sono.
- **Nota importante sugli errori standard:**
 - Gli errori standard OLS dalla regressione del secondo stadio non sono corretti – non tengono conto della stima al primo stadio (è \hat{x}_i stimata).
 - Si utilizza invece un singolo comando apposito che calcola lo stimatore TSLS e gli errori standard corretti.
 - Come di consueto, si usano errori standard robusti all'eteroschedasticità

Esempio 4: domanda di sigarette

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

Perché lo stimatore OLS di β_1 è probabilmente distorto?

- Data set: dati panel sul consumo annuo e i prezzi medi (comprese le imposte) delle sigarette, per stato, per i 48 stati contigui USA, 1985-1995.
 - Variabile strumentale proposta:
 - Z_i = imposta generale sulle vendite al pacchetto nello stato = $SalesTax_i$
 - Pensate che questo strumento sia valido?
1. Rilevante? $\text{corr}(SalesTax_i, \ln(P_i^{\text{cigarettes}})) \neq 0$?
 2. Esogeno? $\text{corr}(SalesTax_i, u_i) = 0$?

Domanda di sigarette (continua)

Per ora usiamo solo i dati del 1995.

Primo stadio regressione OLS:

$$\widehat{\ln(P_i^{cigarettes})} = 4.63 + 0,031SalesTax_i, n = 48$$

Secondo stadio regressione OLS:

$$\widehat{\ln(Q_i^{cigarettes})} = 9,72 - 1,08 \quad \widehat{\ln(P_i^{cigarettes}, n)} = 48$$

Regressione TSLS combinata con errori standard corretti, robusti all'eteroschedasticità:

$$\widehat{\ln(Q_i^{cigarettes})} = 9,72 - 1,08, \quad \widehat{\ln(P_i^{cigarettes}, n)} = 48$$

(1,53) (0,32)

Esempio STATA: domanda di sigarette, primo stadio

Strumento = $Z = rtaxso$ = imposta vendite
(dollari reali/pacchetto)

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]  
-----+-----  
    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;      Ora abbiamo i valori predetti dal primo stadio
```

Secondo stadio

Y *X-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
```

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

- Questi coefficienti sono le stime TSLS
- Gli errori standard sono sbagliati perché ignorano la stima al primo stadio

Tutto in un unico comando:

```

      Y           X           Z
. ivregress 2sls lpackpc (lragvprs = 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				
	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
lpackpc						
lragvprs	-1.083587	.3122035	-3.47	0.001	-1.695494	-.471679
_cons	9.719876	1.496143	6.50	0.000	6.78749	12.65226

Instrumented: lragvprs *This is the endogenous regressor*
Instruments: rtaxso *This is the instrumental variable*

Equazione della domanda di sigarette stimata:

$$\ln(Q_i^{cigarettes}) = 9.72 - 1,08 \ln(P_i^{cigarettes}, \eta) = 48$$

(1,53) (0,31)

Riepilogo della regressione IV con singola X e Z

- Uno strumento valido Z deve soddisfare due condizioni:
 1. *rilevanza*: $\text{corr}(Z_i, X_i) \neq 0$
 2. *esogeneità*: $\text{corr}(Z_i, u_i) = 0$
- TSLS procede eseguendo prima la regressione di X su Z per ottenere \hat{X} , poi di Y su \hat{X}
- Il concetto chiave è che il primo stadio isola la parte della variazione in X che è incorrelata con u
- Se lo strumento è valido, allora la distribuzione in grandi campioni dello stimatore TSLS è normale, perciò l'inferenza procede come di consueto

Il modello generale di regressione IV (Paragrafo 12.2)

- Finora abbiamo considerato la regressione IV con un singolo regressore endogeno (X) e un singolo strumento (Z).
- Dobbiamo estenderla a:
 - più regressori endogeni (X_1, \dots, X_k)
 - più variabili incluse esogene (W_1, \dots, W_r) o variabili di controllo, che devono essere incluse per il consueto motivo delle variabili omesse
 - più variabili strumentali (Z_1, \dots, Z_m). Più strumenti (rilevanti) possono produrre una minore varianza del TSLS: l' R^2 del primo stadio aumenta, perciò si ha maggiore variazione in \hat{X} .
- *Nuovi termini*: identificazione e sovraidentificazione

Identificazione

- In generale si dice che un parametro è ***identificato*** se diversi valori del parametro producono distribuzioni diverse dei dati.
- Nella regressione IV, il fatto che i coefficienti siano identificati dipende dalla relazione tra il numero di strumenti (m) e il numero di regressori endogeni (k)
- Intuitivamente, se ci sono meno strumenti che regressori endogeni, non possiamo stimare β_1, \dots, β_k
 - Per esempio, supponiamo $k = 1$ ma $m = 0$ (nessuno strumento)!

Identificazione (continua)

I coefficienti β_1, \dots, β_k si dicono:

- **esattamente identificati** se $m = k$.

Ci sono esattamente gli strumenti sufficienti per stimare β_1, \dots, β_k .

- **sovraidentificati** se $m > k$.

Ci sono più strumenti di quelli necessari per stimare β_1, \dots, β_k . In questo caso si può verificare se gli strumenti sono validi (test delle "restrizioni sovraidentificanti") – torneremo sul tema in seguito

- **sottoidentificati** se $m < k$.

Ci sono troppo pochi strumenti per stimare β_1, \dots, β_k . In questo caso occorre procurarsi più strumenti!

Il modello generale di regressione IV: riepilogo della terminologia

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

- Y_i è la **variabile dipendente**
- X_{1i}, \dots, X_{ki} sono i **regressori endogeni** (potenzialmente correlati con u_i)
- W_{1i}, \dots, W_{ri} sono i **regressori esogeni inclusi** (incorrelati con u_i) o **variabili di controllo** (inclusi in modo che Z_i sia incorrelata con u_i , una volta inclusi i W)
- $\beta_0, \beta_1, \dots, \beta_{k+r}$ sono i coefficienti di regressione ignoti
- Z_{1i}, \dots, Z_{mi} sono le m **variabili strumentali (variabili esogene escluse)**
- I coefficienti sono **sovraidentificati** se $m > k$; **esattamente identificati** se $m = k$; **sottoidentificati** se $m < k$.

TSLS con un singolo regressore endogeno

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

- m strumenti: Z_{1i}, \dots, Z_{mi}
- Primo stadio
 - Regressione di X_1 su *tutti* i regressori esogeni: regressione di X_1 su $W_{1i}, \dots, W_{ri}, Z_{1i}, \dots, Z_{mi}$, e un'intercetta, usando OLS
 - Calcolo dei valori predetti $\hat{X}_{1i} = 1, \dots, n$
- Secondo stadio
 - Regressione di Y su $\hat{X}_{1i}, W_{1i}, \dots, W_{ri}$, e un'intercetta, usando OLS
 - I coefficienti di questa regressione del secondo stadio sono gli stimatori TSLS, ma gli errori standard sono sbagliati
- Per ottenere errori standard corretti, occorre procedere in un singolo passaggio con il software di regressione

Esempio 4: ancora la domanda di sigarette

Si supponga che il reddito sia esogeno (è plausibile – *perché?*), e di voler anche stimare l'elasticità:

$$\ln(\overbrace{\ln(Q_i^{cigarettes})}) = \beta_0 + \beta_1 \ln(\overbrace{\ln(P_i^{cigarettes})}) + \beta_2 \ln(Income_i) + u_i$$

Abbiamo due strumenti:

Z_{1i} = imposta generale sulle vendite

Z_{2i} = imposta specifica sulle sigarette

- Variabile endogena: $\ln(\overbrace{\ln(P_i^{cigarettes})})$ (“una sola X ”)
- Variabile esogena inclusa: $\ln(Income_i)$ (“una sola W ”)
- Strumenti (variabili endogene escluse): imposta generale vendite, imposta specifica sulle sigarette (“due Z ”)
- β_1 è sotto, sopra o esattamente identificata?

Esempio: domanda di sigarette, un solo strumento

IV: rtaxso = real overall sales tax in state

Y W X Z

```
. ivreg lpackpc lperinc (lragvprs = 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]	
lragvprs	-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: lragvprs

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

Esempio: domanda di sigarette, due strumenti

```

      Y      W      X      Z1      Z2
. ivreg lpackpc lperinc (lragvprs = 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

```

	Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]	
lpackpc						
lragvprs	-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: lragvprs

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

Stime TSLS, $Z =$ imposta vendite ($m = 1$)

$$\widehat{\ln(Q_i^{cigarettes})} = 9,43 - 1,14 \widehat{\ln(P_i^{cigarettes})} + 0,21 \ln(Income_i)$$

(1,26) (0,37) (0,31)

Stime TSLS, $Z =$ imposta vendite e imposta sig. ($m = 2$)

$$\widehat{\ln(Q_i^{cigarettes})} = 9,89 - 1,28 \widehat{\ln(P_i^{cigarettes})} + 0,28 \ln(Income_i)$$

(0,96) (0,25) (0,25)

- **Errori standard minori per $m = 2$.** Con 2 strumenti si hanno più informazioni, più "variazione come se casuale"
- Bassa elasticità al reddito (non è un bene di lusso); elasticità al reddito non significativamente diversa da zero a livello statistico
- Elasticità al prezzo sorprendentemente elevata

Assunzioni generali per la validità di uno strumento

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

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

(2) **Rilevanza**: caso generale, più X

Supponiamo che la regressione del secondo stadio possa essere eseguita usando i valori predetti dalla regressione del primo stadio. Allora non vi è perfetta collinearità in questa (impraticabile) regressione del secondo stadio.

- *Caso speciale di una sola X* : l'assunzione generale è equivalente a (a) almeno uno strumento deve entrare nella controparte della regressione del primo stadio e (b) i W non sono perfettamente collineari.

Assunzioni della regressione IV

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

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

- l'assunzione 1 dice "i regressori esogeni sono esogeni"

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

- l'assunzione 2 non è nuova

3. X, W, Z e Y hanno momenti quarti finiti non nulli

- l'assunzione 3 non è nuova

4. Gli strumenti (Z_{1i}, \dots, Z_{mi}) sono validi.

- Ne abbiamo parlato

- Sotto le assunzioni 1-4, il TSLS e la sua statistica t hanno distribuzione normale

- Il requisito fondamentale è che gli strumenti siano validi

W come variabili di controllo

- In molti casi le W sono incluse allo scopo di controllare per fattori omessi, cosicché, una volta incluse le W , Z è incorrelata con u . In questo caso le W non devono essere esogene, ma devono essere variabili di controllo effettive nel senso discusso nel Capitolo 7 – ora però focalizzandosi sulla produzione di uno strumento esogeno.
- Tecnicamente, la condizione perché le W siano variabili di controllo effettive è che la media condizionata degli u_j non dipenda da Z_j , date W_j :

$$E(u_j | W_j, Z_j) = E(u_j | W_j)$$

W come variabili di controllo (continua)

- Quindi un'alternativa alla prima assunzione della regressione IV è che valga l'indipendenza in media condizionata:

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

Questa è la versione IV dell'assunzione dell'indipendenza in media condizionata del Capitolo 7.

- *Ecco il punto chiave:* in molte applicazioni occorre includere variabili di controllo (W) affinché Z sia verosimilmente esogena (incorrelata con u).
- Per i dettagli si veda l'Appendice 12.6

Esempio 1: effetto dello studio sui voti (continua)

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

Y = media voti del primo semestre

X = media di ore di studio al giorno

$Z = 1$ se il compagno di stanza ha portato un videogioco, = 0 altrimenti

I compagni di stanza sono stati assegnati a caso

Conoscete un motivo per cui Z potrebbe essere correlata con u – anche se è assegnata casualmente? Che cos'altro entra nel termine d'errore, quali sono altri determinanti dei voti, oltre al tempo speso studiando?

Esempio 1: effetto dello studio sui voti (continua)

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

Perché Z potrebbe essere correlata u ?

- Ecco una *ipotetica* possibilità: il genere. Supponiamo:
 - le donne ottengono voti migliori degli uomini, a parità di ore di studio
 - Gli uomini hanno più probabilità di portare un videogioco, rispetto alle donne
 - Allora $\text{corr}(Z_i, u_i) < 0$ (i maschi hanno più probabilità di avere un compagno di stanza [maschio] che porti un videogioco, ma i maschi tendono anche ad avere voti inferiori, a parità di tempo di studio).
- È solo un altro caso di distorsione da variabili omesse. La soluzione sta nel controllare per (o includere) la variabile omessa, in questo caso il genere.

Esempio 1: effetto dello studio sui voti (continua)

- Questa logica porta a includere W = genere come variabile di controllo nella regressione IV:

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

- La stima TSLS qui riportata proviene da una regressione che ha incluso il genere come variabile W – tra altre variabili come la materia di specializzazione.

Verifica della validità degli strumenti (Paragrafo 12.3)

Ricordiamo i due requisiti per strumenti validi:

1. *Rilevanza* (caso speciale di una sola X)

Almeno uno strumento deve entrare nella controparte di popolazione della regressione del primo stadio.

2. *Esogeneità*

Tutti gli strumenti devono essere incorrelati con il termine d'errore:

$$\text{corr}(Z_{1j}, u_j) = 0, \dots, \text{corr}(Z_{mj}, u_j) = 0$$

Che cosa accade se uno di questi requisiti non è soddisfatto? Come si può verificare? Che cosa occorre fare?

Se si hanno più strumenti, quale si deve usare?

Verifica dell'assunzione 1: rilevanza dello strumento

Ci concentreremo su un singolo regressore incluso:

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

Regressione del primo stadio:

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

- Gli strumenti sono rilevanti se almeno uno dei π_1, \dots, π_m è diverso da zero.
- Gli strumenti si dicono **deboli** se tutti i π_1, \dots, π_m sono uguali o vicini a zero.
- **Gli strumenti deboli** dicono molto poco sulla variazione in X , oltre a ciò che dicono le W

Quali sono le conseguenze di strumenti deboli?

Se gli strumenti sono deboli, la distribuzione campionaria del TSLS e della sua statistica t non è normale, anche con n grande.

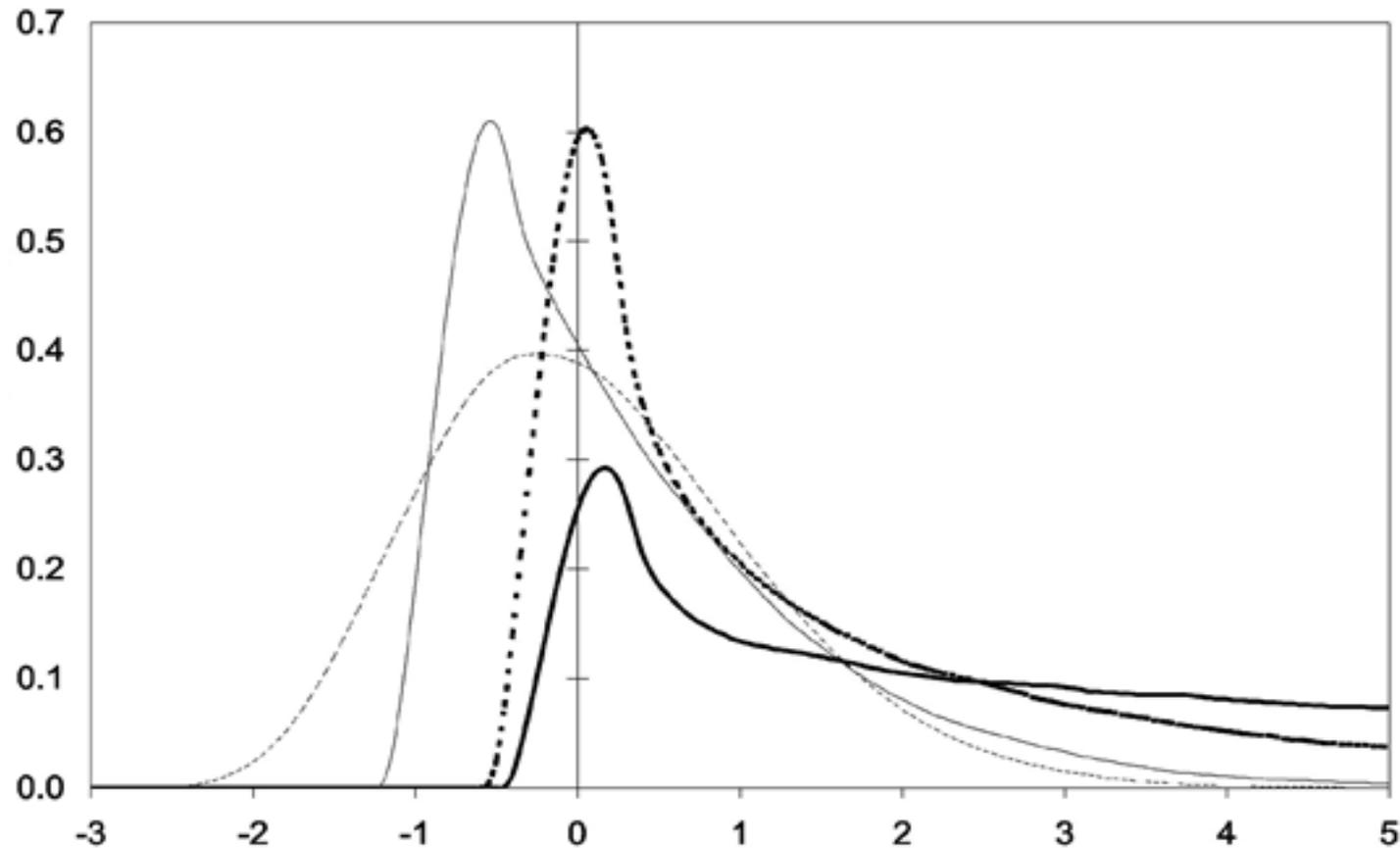
Consideriamo il caso più semplice:

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

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

- Lo stimatore IV è $\hat{\beta}_1^{TSLS} = \frac{s_{YZ}}{s_{XZ}}$
- Se $\text{cov}(X, Z)$ è zero o minore, allora s_{XZ} sarà piccolo: con strumenti deboli, il denominatore è quasi zero.
- In questo caso, la distribuzione campionaria di $\hat{\beta}_1^{TSLS}$ (e la sua statistica t) non è ben approssimata dall'approssimazione normale per n grande...

***Esempio:* la distribuzione campionaria della statistica t del TOLS con strumenti deboli**



Linea scura = strumenti non rilevanti

Linea chiara tratteggiata = strumenti forti

Perché la nostra approssimazione normale ci tradisce?

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

- Se $\text{cov}(X,Z)$ è piccola, piccole variazioni in s_{XZ} (da un campione al successivo) può indurre grandi variazioni in $\hat{\beta}_1^{TSLs}$
- Supponiamo di calcolare in un campione $s_{XZ} = 0,00001\dots$
- Allora l'approssimazione normale per n grande non è una buona approssimazione della distribuzione campionaria di $\hat{\beta}_1^{TSLs}$
- Un'approssimazione migliore è quella di $\hat{\beta}_1^{TSLs}$ come il *rapporto* di due variabili casuali normali correlate (si veda l'Appendice 12.4)
- Se gli strumenti sono deboli, i consueti metodi di inferenza sono inaffidabili – potenzialmente molto inaffidabili.

Misurazione della forza degli strumenti in pratica: la statistica F del primo stadio

- La regressione del primo stadio (una sola X):
- Regressione di X su $Z_1, \dots, Z_m, W_1, \dots, W_k$.
- Strumenti totalmente irrilevanti \leftrightarrow *tutti* i coefficienti di Z_1, \dots, Z_m sono zero.
- La **statistica F del primo stadio** verifica l'ipotesi che Z_1, \dots, Z_m non entrino nella regressione del primo stadio.
- Strumenti deboli implicano un valore basso della statistica F del primo stadio.

Verifica di strumenti deboli con una singola X

- Si calcola la statistica F del primo stadio.

Regola empirica: se la statistica F del primo stadio è minore di 10, allora l'insieme di strumenti è debole.

- In questo caso, lo stimatore TSLS sarà distorto, e le inferenze statistiche (errori standard, verifiche di ipotesi, intervalli di confidenza) possono essere fuorvianti.

Verifica di strumenti deboli con una singola X (continua)

- Perché confrontare la statistica F del primo stadio con 10?
- Non è sufficiente respingere l'ipotesi nulla che i coefficienti delle Z siano zero – serve un contenuto predittivo sostanziale per una buona approssimazione normale.
- Il confronto della statistica F del primo stadio con 10 verifica se la distorsione del TSLS, rispetto all'OLS, è minore del 10%. Se la F è minore di 10, la distorsione relativa è superiore al 10%, cioè il TSLS può avere una distorsione sostanziale (si veda l'Appendice 12.5).

Che cosa fare se si hanno strumenti deboli

- Procurarsi strumenti migliori (più facile a dirsi che a farsi!)
- Se si hanno molti strumenti, alcuni sono probabilmente più deboli di altri ed è una buona idea scartare i più deboli (scartando uno strumento irrilevante si aumenta la statistica F del primo stadio)
- Se si hanno pochi strumenti, e sono tutti deboli, allora occorre eseguire un'analisi IV al di là del TSLS...
 - Separare il problema della stima di β_1 e della costruzione di intervalli di confidenza
 - Sembra strano, ma se il TSLS non ha distribuzione normale, ha senso (davvero?)

Intervalli di confidenza con strumenti deboli

- Con strumenti deboli, gli intervalli di confidenza TSLS non sono validi, ma altri intervalli di confidenza lo sono. Riportiamo due modi per calcolare intervalli di confidenza validi in grandi campioni anche se gli strumenti sono deboli:
 1. L'intervallo di confidenza di Anderson-Rubin
 - L'intervallo di confidenza di Anderson-Rubin si basa sulla statistica test di Anderson-Rubin che verifica $\beta_1 = \beta_{1,0}$:
 - Si calcola $u_i = Y_i - \beta_{1,0}X_i$
 - Si esegue la regressione su $W_{1j}, \dots, W_{ri}, Z_{1j}, \dots, Z_{mi}$
 - Il test AR è la statistica F su Z_{1j}, \dots, Z_{mi}
 - Ora si inverte il test: l'intervallo di confidenza AR al 95% è l'insieme di β_1 non rifiutati al livello del 5% dal test AR.
 - Calcolo: si usa software specialistico.

Intervalli di confidenza con strumenti deboli (continua)

2. Intervallo di confidenza del rapporto di verosimiglianza condizionato di Moreira
 - L'intervallo di confidenza del rapporto di verosimiglianza condizionato è basato sull'inverso del test del rapporto di verosimiglianza condizionato di Moreira. Per calcolare questo test, il suo valore critico e l'intervallo di confidenza del rapporto di verosimiglianza condizionato, è richiesto un software specialistico.
 - Questo intervallo di confidenza tende a essere più ristretto di quello di Anderson-Rubin, soprattutto quando vi sono molti strumenti.
 - Se si dispone di un software che produce questo intervallo, è il caso di usarlo.

Stima con strumenti deboli

Non ci sono stimatori non distorti se gli strumenti sono deboli o irrilevanti. Tuttavia, alcuni stimatori hanno una distribuzione più centrata su β_1 del TSLS.

- Uno di questi stimatori è quello di massima verosimiglianza con informazione limitata (LIML)
- Lo stimatore LIML
 - può essere derivato come stimatore di massima verosimiglianza
 - è il valore di β_1 che minimizza il valore- p del test AR (!)
- Per approfondire stimatori, verifiche e intervalli di confidenza nel caso di strumenti deboli, si veda l'Appendice 12.5

Verifica dell'assunzione 2: esogeneità dello strumento

- Esogeneità dello strumento: **Tutti** gli strumenti sono correlati con il termine d'errore:
 $\text{corr}(Z_{1j}, u_j) = 0, \dots, \text{corr}(Z_{mj}, u_j) = 0$
- Se gli strumenti sono correlati con il termine d'errore, il primo stadio del TSLS non può isolare una componente di X incorrelata con il termine d'errore, perciò \hat{X} è correlata con u e il TSLS è inconsistente.
- Se ci sono più strumenti che regressori endogeni, è possibile verificare – *parzialmente* – l'esogeneità dello strumento.

Verifica di restrizioni di sovraidentificazione

Consideriamo il caso più semplice:

$$Y_i = \beta_0 + \beta_1 X_i + u_{ir}$$

- Supponiamo che vi siano due strumenti validi: Z_{1ir} , Z_{2i}
- Allora potremmo calcolare due stime TSLS separate.
- Intuitivamente, se queste due stime TSLS sono molto diverse tra loro, ci dev'essere qualcosa di sbagliato: uno strumento o l'altro (o entrambi) devono essere non validi.
- Il test J di restrizioni sovraidentificanti esegue questo confronto in un modo statisticamente preciso.
- Si può fare soltanto se il numero di Z è maggiore del numero di X (sovraidentificazione).

Il test J di restrizioni di sovraidentificazione

Supponiamo che il numero di strumenti = $m >$ numero di $X = k$
(sovraidentificazione)

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

Il test J è il test di Anderson-Rubin, usando lo stimatore TSLS al posto del valore ipotizzato $\beta_{1,0}$. Procedura:

1. Prima si stima l'equazione di interesse usando TSLS e tutti gli m strumenti; si calcolano i valori predetti \hat{Y}_i , usando le X effettive (non le \hat{X} usate per stimare il secondo stadio)
2. Si calcolano i residui $\hat{u}_i = Y_i - \hat{Y}_i$
3. Si esegue la regressione rispetto a $Z_{1i}, \dots, Z_{mi}, W_{1i}, \dots, W_{ri}$
4. Si calcola la statistica F che verifica l'ipotesi che i coefficienti di Z_{1i}, \dots, Z_{mi} siano tutti zero;
5. La **statistica J** è $J = mF$

Il test J (continua)

$J = mF$, dove F = la statistica F che verifica i coefficienti di Z_{1i}, \dots, Z_{mi} in una regressione dei residui TOLS rispetto a $Z_{1i}, \dots, Z_{mi}, W_{1i}, \dots, W_{ri}$.

Distribuzione della statistica J

- Sotto l'ipotesi nulla che tutti gli strumenti siano esogeni, J ha una distribuzione chi-quadro con $m-k$ gradi di libertà
- Se $m = k$, $J = 0$ (ha senso?)
- Se alcuni strumenti sono esogeni e altri sono endogeni, la statistica J sarà grande, e l'ipotesi nulla che tutti gli strumenti sono esogeni sarà respinta.

Verifica della validità degli strumenti: riepilogo

Questo riepilogo considera il caso di una singola X . I due requisiti per la validità degli strumenti sono:

1. *Rilevanza*

- Almeno uno strumento deve entrare nella controparte della regressione del primo stadio.
- Se gli strumenti sono deboli, allora lo stimatore TSLS è distorto e la statistica t ha una distribuzione non normale
- Per verificare strumenti deboli con un singolo regressore endogeno incluso, si verifica la statistica F del primo stadio
 - Se $F > 10$, gli strumenti sono forti – si usa il TSLS
 - Se $F < 10$, gli strumenti sono deboli – si fa qualcosa.

2. Esogeneità

- **Tutti** gli strumenti devono essere incorrelati con il termine d'errore: $\text{corr}(Z_{1i}, u_i) = 0, \dots, \text{corr}(Z_{mi}, u_i) = 0$
- Possiamo eseguire una verifica parziale di esogeneità: se $m > 1$, possiamo verificare l'ipotesi nulla che tutti gli strumenti siano esogeni contro l'alternativa che almeno $m-1$ siano endogeni (correlati con u)
- Si usa il test J , realizzato usando i residui TSLS.
- Se il J respinge l'ipotesi, allora almeno alcuni degli strumenti sono endogeni, perciò occorre prendere una decisione difficile e scartare alcuni (o tutti) gli strumenti.

Applicazione alla domanda di sigarette (Paragrafo 12.4)

Perché siamo interessati a conoscere l'elasticità della domanda di sigarette?

- Teoria della tassazione ottimale. L'aliquota d'imposta ottimale è inversamente proporzionale all'elasticità al prezzo: maggiore è l'elasticità, minore la quantità influenzata da una data percentuale d'imposta, perciò minore è la variazione di consumo e perdita secca.
- Esternalità del fumo – ruolo dell'intervento pubblico per scoraggiare il fumo
 - effetti di salute del fumo passivo? (non monetari)
 - esternalità monetarie

Dati panel

- Consumo annuo di sigarette, prezzi medi pagati dal consumatore finale (tasse incluse), reddito personale e percentuali d'imposta (specifiche per le sigarette e generali sulle vendite nello stato)
- 48 stati continentali USA, 1985-1995

Strategia di stima

- Dobbiamo usare metodi di stima IV per gestire la distorsione da causalità simultanea che nasce dall'interazione di offerta e domanda.
- Indicatori binari di stato = variabili W (variabili di controllo) che controllano per caratteristiche inosservate a livello di stato che influiscono sulla domanda di sigarette e la percentuale d'imposta, purché tali caratteristiche non varino nel tempo.

Modello con effetti fissi della domanda di sigarette

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

- $i = 1, \dots, 48$, $t = 1985, 1986, \dots, 1995$
- $\text{corr}(\ln(P_{it}^{cigarettes}), u_{it})$ è verosimilmente diversa da zero a causa di interazioni offerta-domanda
- α_i riflette fattori omessi inosservati che variano tra stati ma non nel tempo, per esempio l'atteggiamento verso il fumo
- Strategia di stima:
 - Usiamo metodi di regressione con dati panel per eliminare α_i
 - Usiamo TSLS per gestire la distorsione da causalità simultanea
 - Usiamo $T = 2$ con variazioni 1985 – 1995 (metodo "prima e dopo") – osserviamo la risposta a lungo termine, non la dinamica di breve termine (elasticità a breve v. lungo termine)

Il metodo "prima e dopo" (quando $T=2$)

- Un modo per modellare gli effetti a lungo termine è quello di considerare variazioni su 10 anni, tra il 1985 e il 1995
- Riscriviamo la regressione in forma "prima e dopo" :

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

- Creiamo variabili di "variazione a 10 anni", per esempio:
- Variazione a 10 anni nel log del prezzo = $\ln(P_{i1995}) - \ln(P_{i1985})$
- Poi stimiamo l'elasticità della domanda mediante TSLS usando variazioni a 10 anni nelle variabili strumentali
- Questo è equivalente a usare i dati originali e includere gli indicatori binari di stato (variabili "W") nella regressione

STATA: domanda di sigarette

Prima si creano variabili di “variazione a 10 anni”

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]);  
. 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];
```

_n-10 è il valore a 10 anni

Usiamo TSLS per stimare l'elasticità della domanda con la specifica "variazione a 10 anni"

```
. ivregress 2sls Y W X Z
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
```

	Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]	
dlpackpc						
dlavgprs	<i>-0.9380143</i>	<i>0.2075022</i>	-4.52	0.000	-1.355945	-.5200834
dlperinc	<i>0.5259693</i>	<i>0.3394942</i>	1.55	0.128	-.1578071	1.209746
_cons	0.2085492	0.1302294	1.60	0.116	-.0537463	.4708446

```
Instrumented:  dlavgprs
Instruments:   dlperinc drtaxso
```

NOTE:

- *Tutte le variabili - Y, X, W e Z - sono in variazioni a 10 anni*
- *Elasticità stimata = -0,94 (SE = 0,21) - sorprendentemente elastica!*
- *Elasticità del reddito piccola, non statisticamente diversa da zero*
- *Occorre verificare se lo strumento è rilevante...*

Verifica della rilevanza dello strumento: si calcola la statistica F del primo stadio

```
. reg dlavgrs 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	

	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
```

```
F( 1, 45) = 46.41  
Prob > F = 0.0000
```

*Non serviva eseguire "test" qui!
Con strumento $m=1$, la stat F
è il quadrato della stat t :
 $6,81*6,81 = 46,41$*

F del primo stadio = 46,5 > 10 perciò lo strumento non è debole

*Possiamo verificare l'esogeneità dello strumento? **No**: $m = k$*

Domanda di sigarette, variazioni a 10 anni – 2 IV

```

      Y           W           X           Z1   Z2
. ivregress 2sls dlpackpc 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

```

	Coef.	Robust Std. Err.	z	P> z	[95% Conf. Interval]	
dlpackpc						
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 = solo imposta generale sull vendite

drtax = solo imposta specifica sulle sigarette

Elasticità stimata = -1,2, anche più elastica rispetto all'uso della sola imposta generale sulle vendite!

Statistica F del primo stadio – entrambi gli strumenti

```

      X      Z1      Z2      W
. reg dlavgprs drtaxso drtax dlperinc ;

```

Source	SS	df	MS	Number of obs =	48
Model	.289359873	3	.096453291	F(3, 44) =	51.36
Residual	.082627329	44	.001877894	Prob > F =	0.0000
Total	.371987202	47	.007914621	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 perciò gli strumenti non sono deboli
```

```
Prob > F = 0.0000
```

Con $m > k$, possiamo verificare le restrizioni di sovraidentificazione...

Verifica delle restrizioni di sovraidentificazione

```
. predict e, resid;           Calcola valori predetti per l'ultima regressione
                               (la precedente regressione TSLS)
. 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           Calcola la statistica J, che è m*F,
( 2) drtax = 0           dove F verifica se i coefficienti
                          degli strumenti sono zero
                          perciò J = 2 × 2.47 = 4.93
F( 2, 44) = 2.47
Prob > F = 0.0966        ** ATTENZIONE - usa la f.d. sbagliata **
```

I gradi di libertà corretti per la statistica J sono $m-k$:

- $J = mF$, dove F = la statistica F che verifica i coefficienti di Z_{1i}, \dots, Z_{mi} in una regressione dei residui TSLS rispetto a $Z_{1i}, \dots, Z_{mi}, W_{1i}, \dots, W_{mi}$.
- Sotto l'ipotesi nulla che tutti gli strumenti siano esogeni, J ha una distribuzione chi-quadro con $m-k$ gradi di libertà
- Qui $J = 4,93$, distribuzione chi-quadro con f.d. = 1; il valore critico al 5% è 3,84, perciò respinge al livello di significatività del 5%.
- 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
```

$$J = 2 \times 2.47 = 4.93$$

valore-p da distribuzione chi-quadro(1)

E ora???

Riepilogo dei risultati in forma di tabella

Tabella 12.1 Stime TOLS del consumo di sigarette su dati panel per 48 stati USA.

Variabile dipendente: $\ln(Q_{i,1995}^{cigarettes}) - \ln(Q_{i,1985}^{cigarettes})$

Regressore	(1)	(2)	(3)
$\ln(P_{i,1995}^{cigarettes}) - \ln(P_{i,1985}^{cigarettes})$	-0,94** (0,21)	-1,34** (0,23)	-1,20** (0,20)
$\ln(Inc_{i,1995}) - \ln(Inc_{i,1985})$	0,53 (0,34)	0,43 (0,30)	0,46 (0,31)
Intercetta	-0,12 (0,07)	-0,02 (0,07)	-0,05 (0,06)
Variabili strumentali	Sales tax	Sigarette-specific tax	Both sales tax and sigarette-specific tax
Statistica F first stage	33,70	107,20	88,60
Restrizioni di sovraidentificazione test J e p -value			4,93 (0,026)

Queste regressioni sono state stimate sui dati per 48 stati USA (48 osservazioni su differenze decennali). I dati sono descritti nell'Appendice 12.1. Il test- J per le restrizioni di sovraidentificazione è descritto nel Concetto chiave 12.6 (il suo valore- p è riportato tra parentesi) e la statistica F per il primo stadio è descritta nel Concetto chiave 12.5. I coefficienti sono statisticamente significativi al livello *5% o **1%.

Come dobbiamo interpretare il rifiuto del test J ?

- Il test J rifiuta l'ipotesi nulla che entrambi gli strumenti siano esogeni
- Questo significa che o $rtaxso$ è endogeno, o $rtax$ è endogeno, o entrambi!
- Il test J non ci dice quale! *Occorre ragionare...*
- Perché $rtax$ (imposta sulle sigarette) potrebbe essere endogeno?
 - Forze politiche: tradizione del fumare o molti fumatori -> pressione politica per basse imposte sulle sigarette
 - In questo caso, l'imposta sulle sole sigarette è endogena
- Questo ragionamento non vale per l'imposta generale sulle vendite
- → usiamo solo uno strumento, l'imposta generale sulle vendite

La domanda di sigarette: riepilogo di risultati empirici

- Usiamo l'elasticità stimata in base al TSLS con l'imposta generale sulle vendite come unico strumento:

Elasticità = $-0,94$, $SE = 0,21$

- Questa elasticità è sorprendentemente elevata (non anelastica) – un incremento dell'1% nei prezzi riduce le vendite di sigarette di quasi l'1%. È un'elasticità molto maggiore di quanto si pensi comunemente nella letteratura sull'economia e la salute.
- È un'elasticità di lungo periodo (variazione a 10 anni). *Che cosa vi aspettate riguardo l'elasticità di breve periodo (variazione a un anno), sarà maggiore o minore?*

Valutazione della validità dello studio

Altre minacce alla validità interna?

1. Distorsione da variabili omesse?

- *Lo stimatore con effetti fissi controlla per fattori inosservati che variano tra stati ma non nel tempo*

2. Errata specificazione della forma funzionale? (*si potrebbe verificare*)

3. Altra distorsione da causalità simultanea?

- *Non se l'imposta generale sulle vendite è uno strumento valido, una volta inclusi gli effetti fissi di stato!*

4. Distorsione da errori nelle variabili?

5. Distorsione da selezione campionaria? (*no, abbiamo tutti gli stati*)

6. Altre minacce alla validità interna degli studi di regressione IV riguardano il dubbio che lo strumento sia (1) rilevante e (2) esogeno. *Quanto sono significative queste minacce nell'applicazione dell'elasticità della domanda di sigarette?*

Valutazione della validità dello studio (continua)

Validità esterna?

- Abbiamo stimato un'elasticità nel lungo periodo, possiamo generalizzarla al breve periodo? Perché, o perché no?
- Supponiamo di voler usare l'elasticità stimata di $-0,94$ per orientare la politica odierna. Riportiamo due variazioni a partire dal periodo coperto dai dati (1985-95) – queste variazioni pongono una minaccia alla validità esterna (generalizzazione dal periodo 1985-95 a oggi)?
 - Oggi si fuma meno che nel 1985-1995
 - L'atteggiamento culturale verso il fumo è cambiato in senso negativo dal 1985-95.

Dove trovare strumenti validi? (Paragrafo 12.5)

Note generali

Nell'analisi IV il difficile è trovare strumenti validi

- Metodo 1: "variabili in un'altra equazione" (per es. fattori di spostamento dell'offerta che non hanno effetto sulla domanda)
- Metodo 2: cercare una variazione esogena (Z) che sia "come se" assegnata casualmente (non influisce direttamente su Y) ma influisca su X .
- Questi sono due modi diversi di pensare agli stessi problemi – vedere il collegamento...
 - La pioggia sposta la curva di offerta del burro ma non la curva di domanda. È "come se" assegnata casualmente
 - L'imposta sulle vendite sposta la curva di offerta per le sigarette ma non la curva di domanda; le imposte sulle vendite sono "come se" assegnate casualmente

Esempio: cateterizzazione cardiaca

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.

La cateterizzazione cardiaca migliora la longevità dei pazienti affetti da attacco cardiaco?

Y_i = sopravvivenza (in giorni) di pazienti colpiti da attacco cardiaco
 $X_i = 1$ se il paziente riceve la cateterizzazione cardiaca,
= 0 altrimenti

- Trial clinici mostrano che *CardCath* influisce su *SurvivalDays*.
- Ma il trattamento è efficace "sul campo"?

Cateterizzazione cardiaca (continua)

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

- L'OLS è non distorto? La decisione di trattare un paziente con cateterizzazione cardiaca è endogena – è (*stata*) presa sul campo dal soccorritore e dipende da u_i (caratteristiche di salute del paziente inosservate)
- Se i pazienti più sani sono cateterizzati, allora l'OLS è affetto da distorsione da causalità simultanea e sovrastima l'effetto del trattamento
- Strumento proposto: distanza del più vicino ospedale dotato di cateterizzazione cardiaca meno distanza del più vicino ospedale "normale"

Cateterizzazione cardiaca (continua)

- Z = distanza differenziale dall'ospedale dotato di CC
 - Rilevante? Se un ospedale dotato di CC è lontano, il paziente non vi sarà portato e non sarà trattato con CC
 - Esogena? Se la distanza dell'ospedale dotato di CC non influisce sulla sopravvivenza, se non per l'effetto su $CardCath_i$, allora $corr(distance, u_i) = 0$ perciò è esogena
 - Se la posizione del paziente è casuale, allora la distanza differenziale è "come se" fosse assegnata casualmente.
 - *Il primo stadio è un modello di probabilità lineare: la distanza influisce sulla probabilità di ricevere il trattamento*
- Risultati:
 - OLS stima un significativo e ampio effetto della CC
 - TSLS stima un effetti piccolo, spesso insignificante

Esempio: crowding out (spiazzamento) della spesa privata in beneficenza

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

La spesa sociale pubblica spiazza la spesa di beneficenza privata (chiesa, croce rossa, ecc.)?

Y = spesa di beneficenza privata (chiese)

X = spesa pubblica

Qual è il motivo per usare variabili strumentali?

Strumento proposto:

Z = forza della delega congressuale

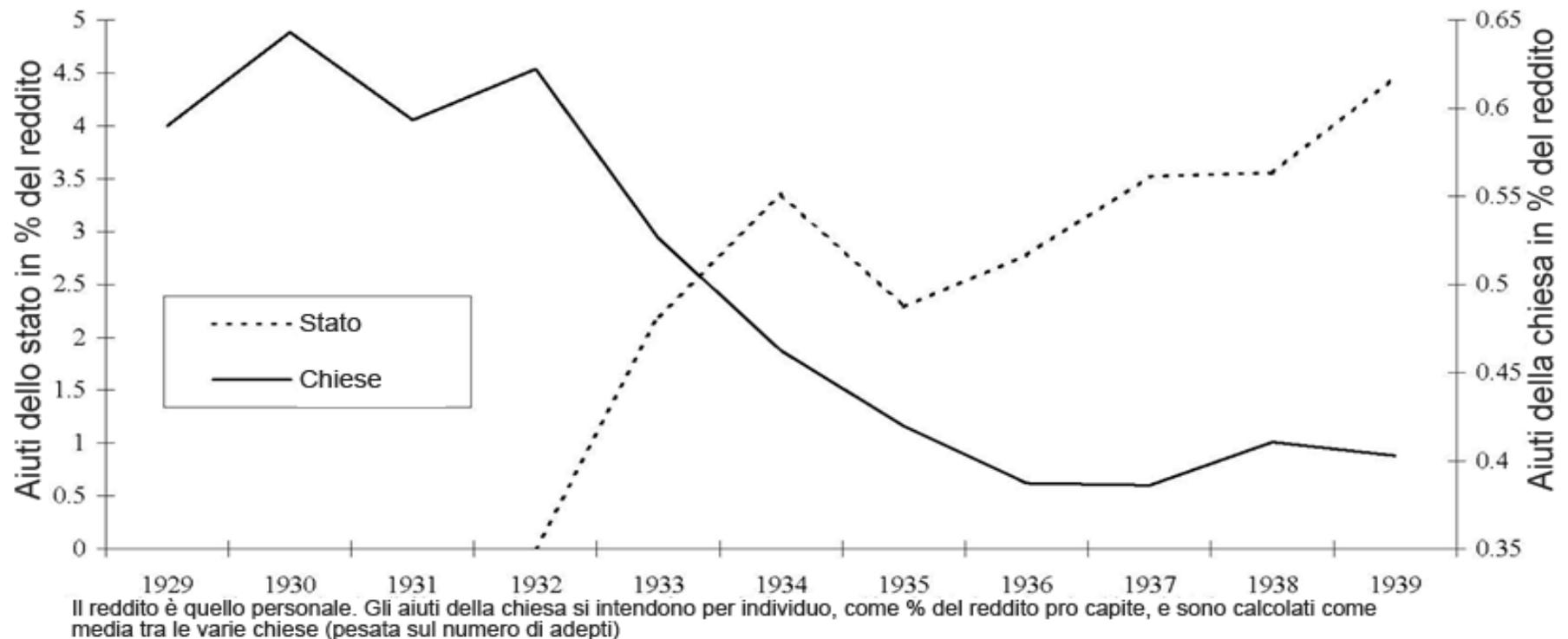
Spesa di beneficenza privata (continua).

I dati – alcuni dettagli

- Dati panel, annui, per stato, 1929-1939, U.S.A.
- Y = totale spesa beneficenza da parte di sei chiese (CCC, Lutheran, Northern Baptist, Presbyterian (2), Southern Baptist); beneficenza = $\frac{1}{4}$ della spesa totale delle chiese.
- X = spesa sociale federale sotto il New Deal (General Relief, Work Relief, Civil Works Administration, Aid to Dependent Children,...)
- Z = durata mandato dei rappresentanti dello stato al House & Senate Appropriations Committees, in mesi
- W = insiemi di effetti fissi

Spesa di beneficenza privata (continua)

Aiuti dello stato e della chiesa durante la Grande Depressione



Spesa di beneficenza privata (continua)

Valutazione di validità:

- Validità strumento:
 - Rilevanza?
 - Esogeneità?
- Altre minacce alla validità interna:
 1. Distorsione da variabili omesse
 2. Forma funzionale
 3. Errore di misura
 4. Selezione del campione
 5. Causalità simultanea
- Validità esterna a oggi negli USA? Per aiuti ai paesi in via di sviluppo?

Esempio: concorrenza tra scuole

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

Qual è l'effetto della concorrenza tra scuole pubbliche sulla performance degli studenti?

Y = punteggi nei test al 12-esimo livello

X = misura di scelta tra distretti scolastici (funzione del numero di distretti nell'area)

Qual è la motivazione per usare variabili strumentali?

Strumento proposto:

Z = numero di piccoli corsi d'acqua nell'area

Concorrenza tra scuole (continua)

Dati – alcuni dettagli

- dati sezionali, USA, area metropolitana, fine anni Novanta ($n = 316$),
- Y = voto al 12-esimo livello scolastico (anche altre misure)
- X = indice preso da una revisione della letteratura di settore misurando il livello di concorrenza (“indice di Gini”) – in base al numero di “imprese” e alla “quota di mercato”
- Z = misura del numero di piccoli corsi d’acqua – che hanno formato confini geografici naturali.
- W = insiemi di variabili di controllo

Concorrenza tra scuole (continua)

Valutazione di validità:

- Validità strumento:
 - Rilevanza?
 - Esogeneità?
- Altre minacce alla validità interna:
 1. Distorsione da variabili omesse
 2. Forma funzionale
 3. Errore di misura
 4. Selezione del campione
 5. Causalità simultanea
- Validità esterna a oggi negli USA?

Conclusioni

(Paragrafo 12.6)

- Uno strumento valido ci consente di isolare una parte di X che è incorrelata con u , e quella parte può essere usata per stimare l'effetto su Y di una variazione in X
- La regressione IV richiede strumenti validi:
 1. *Rilevanza*: verifica tramite statistica F del primo stadio
 2. *Esogeneità*: verifica di restrizioni di sovraidentificazione tramite la statistica J
- Uno strumento valido isola la variazione in X che è "come se" assegnata casualmente.
- Il requisito fondamentale di almeno m strumenti validi non può essere verificato – *occorre usare la testa*.

Domande e risposte sulla regressione IV

1. Quando usare la regressione IV?

Ogni volta che X è correlata con u e si ha uno strumento valido. I motivi principali per la correlazione tra X e u potrebbero essere:

- Variabili omesse che portano a distorsione
 - Esempio: distorsione da talento nel rendimento dell'istruzione
- Errore di misura
 - Esempio: errore di misura negli anni di istruzione
- Distorsione da selezione del campione
 - I pazienti scelgono il trattamento
- Distorsione da causalità simultanea
 - Esempio: offerta e domanda di burro, sigarette

2. Quali sono le minacce alla validità interna di una regressione IV?

- La minaccia principale alla validità interna di una regressione IV è la non validità dell'assunzione di strumenti validi. Dato un insieme di variabili di controllo W , gli strumenti sono validi se sono rilevanti ed esogeni.
 - La rilevanza può essere valutata verificando se gli strumenti sono deboli o forti: la statistica F del primo stadio è > 10 ?
 - L'esogeneità può essere verificata usando la statistica J – purché si abbiano m strumenti esogeni con cui partire! In generale, l'esogeneità deve essere valutata basandosi su una conoscenza approfondita dell'applicazione considerata.