

ECMA 31100: Intro to Empirical Analysis II

Anderson Rubin Test; Many Instruments

Joe Hardwick

University of Chicago

Winter 2022

$$y = x' \beta + u$$

$$x = z'\pi + v$$

$\pi \approx 0$ because instruments were 'weak'.

$$\bar{\pi} = 0: \text{ (Everything scalar). } \hat{\beta}_{IV} \xrightarrow{d} \text{plim}(\hat{\beta}_{\text{OLS}}) + \frac{\sigma_{uv}}{\sigma_v^2}$$

(Under Standard assumptions).

Model 'weak' but not completely irrelevant instruments using a first stage in which π is 'local to zero'.

$$x = z'\pi_n + v \quad \pi_n = \frac{\pi}{\sqrt{n}}$$

Under 'weak instrument asymptotics'

$$\hat{\beta}_{IV} - \beta \xrightarrow{d} \frac{A}{\pi_n E(z^2) + B} \quad \begin{aligned} A &\rightarrow \text{Ratio of correlated normals} \\ \left(\begin{array}{c} A \\ B \end{array} \right) &\sim N\left(\begin{array}{c} 0 \\ 0 \end{array}, \begin{array}{cc} \sigma_u^2 & \sigma_{uv} \\ \sigma_{uv} & \sigma_v^2 \end{array} \right) \end{aligned}$$

'Usual asymptotics' $\sqrt{n}(\hat{\beta}_{IV} - \beta) \stackrel{d}{\sim} N(0, V)$

$$\hat{\beta}_{IV} - \beta \stackrel{a}{\sim} \frac{A}{\pi_n E(z^2) + B/\sqrt{n}}$$

π is large ≈ 0

$\pi_n \neq 0$, so weak

mt. asymptotics agree with standard asymptotic

Stock, Wright and Yogo (2002)

- Stock, Wright and Yogo (2002) simulate the finite sample distribution of $\hat{\beta}_{IV}$:

$$y = \beta x + u; \quad \sigma_u = \sigma_v = 1. \quad \sigma_{uv} = 0.99.$$
$$x = \pi z + v;$$

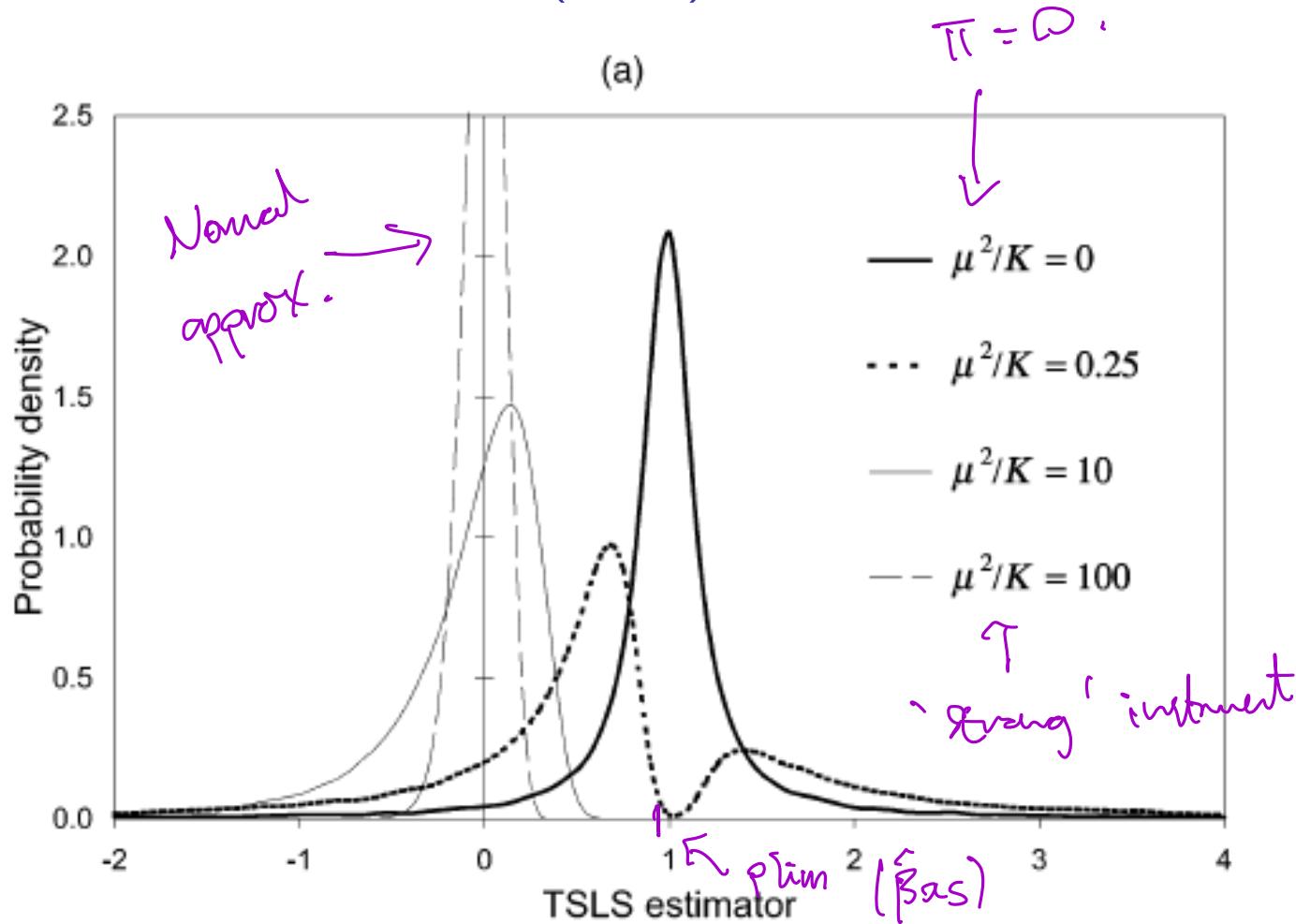
$$\begin{pmatrix} u \\ v \end{pmatrix} \sim \mathcal{N} \left(0, \begin{pmatrix} 1 & 0.99 \\ 0.99 & 1 \end{pmatrix} \right),$$

and $\beta = 0$. Instruments non-random. In this case

$$\hat{\beta}_{OLS} \xrightarrow{P} \frac{\sigma_{uv}}{\sigma_v^2} = 0.99. \quad \mu = \pi \left(\sum_{i=1}^n z_i^2 \right)^{1/2}. \quad K = \dim(z) = 1.$$

\nearrow
Concentration
parameter

Stock, Wright and Yogo (2002)



Conducting inference on β

- Two-step methods that check the first stage F -statistic and then use $\hat{\beta}_{IV}$ to conduct inference on β are often unreliable, and it's generally difficult to come up with a reasonable value for the F -statistic.
- Rule of thumb $F\text{-stat} \geq 10$ reasonable for absolute relative bias of 10% and for the specific case of three or more instruments (where the mean actually exists, rather than being approximated by a Nagar bias). Stock and Yogo (2005) show the critical value lies between 9 and 12 for number of instruments between 3 and 30.
- In joint normal errors/fixed instruments case, $\hat{\beta}_{IV}$ has number of moments equal to the number of excluded instruments - number of endogenous variables (Kinal (1980)).

Conducting inference on β

- Other methods of defining weak instruments exist, for example actual size of t -test should be 'close' to specified significance level.
- With several endogenous variables, the bias of the 2SLS estimate becomes the norm of the vector of biases for the coefficients on each of the endogenous variables.
- Under heteroskedasticity, Montiel Olea and Pflueger (2013) suggest a modification to the F-stat, but it's no longer the case that $\hat{\beta}_{IV}$ is centered at the prob. limit of OLS. See Pflueger and Wang (2015) for a Stata implementation.

Conducting inference on β

- We use a test that is robust to weak instruments, called Anderson-Rubin test. Suppose $y = x'\beta + u$.
- Idea: $E(zu) = 0$ implies $E(z(y - x'\beta)) = 0$, so under $H_0 : \beta = \beta_0$,
$$E(zu(\beta_0)) = 0, \quad E(z(y - x'\beta_0)) = 0.$$
where $u(\beta_0) = y - x'\beta_0$. We can write this equivalently as

$$E(zz')^{-1} E(zu(\beta_0)) = 0,$$

and this is the vector of coefficients of a regression of $u(\beta_0)$ on z .

Instruments are linearly independent -
 $E(zu) = 0 \Leftrightarrow E(zz')^{-1} E(zu) = 0.$

Conducting inference on β $u(\beta_0) = z' \gamma + \varepsilon$. $E(\varepsilon|z)=0$

Run regression estimate $\hat{\gamma}$, Under H_0 : $\gamma=0$.
use F-test to assess $H_0: \gamma=0$.

- Reject if coefficient estimates are significantly different from zero.
- Substitute $y = x'\beta + u$ to give

$$\begin{aligned} u(\beta_0) &= y - x'\beta_0 \\ &= x'\beta - x'\beta_0 + u. \end{aligned}$$

$$\begin{aligned} E(zz')^{-1} E(zu(\beta_0)) &= E(zz')^{-1} E(zx'(\beta - \beta_0) + zu) \\ &= E(zz')^{-1} E(zx') (\beta - \beta_0). \end{aligned} \quad \text{↑ validity.}$$

which equals zero iff $\beta = \beta_0$.

- Therefore, should expect such a test to have power against any alternative $\beta \neq \beta_0$.

Conducting inference on β

- Suppose $z \in \mathbb{R}^I$ and run the regression

$$y_i - x_i' \beta_0 = z' \gamma + \epsilon; \quad E(z\epsilon) = 0.$$

Under H_0 ,

$$\gamma = E(zz')^{-1} E(z(y_i - x_i' \beta_0)) = 0.$$

Assuming all relevant moments exist:

$$\sqrt{n}(\hat{\gamma} - \gamma) \xrightarrow{d} \mathcal{N}(0, V),$$

and so

$$n(\hat{\gamma} - \gamma)' \hat{V}^{-1} (\hat{\gamma} - \gamma) \xrightarrow{d} \chi_I^2.$$

$$\begin{aligned}\hat{\gamma} &\in \mathbb{R}^L \\ \hat{V}^{-\frac{1}{2}} \sqrt{n}(\hat{\gamma} - \gamma) &\xrightarrow{d} N(0, I_L)\end{aligned}$$

Conducting inference on β

- Under H_0 , we have

$$T_n = n(\hat{\gamma})' \hat{V}^{-1}(\hat{\gamma}) \xrightarrow{d} \chi_I^2,$$

and we reject if $T_n > \chi_{I,1-\alpha}^2$.

- This asymptotic argument does not depend on whether instruments are weak, because under H_0 :

$$\sqrt{n}(\hat{\gamma} - \gamma) = \left(\frac{1}{n} \sum_{i=1}^n z_i z_i' \right)^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n z_i u_i.$$

Conducting inference on β

- Let's separate the included and excluded instruments:
 $z = (z_1, z_2)$ where z_1 are included and z_2 are excluded.
- Now write $x = (x_1, z_1)$, where x_1 denote the endogenous variables. Then

$$y = x'_1 \beta_1 + z'_1 \beta_2 + u$$
$$x_1 = z'_1 \pi_1 + z'_2 \pi_2 + v.$$

- Typically only interested in conducting inference on β_1 . Under $H_0 : \beta_1 = \beta_{1,0}$:

$$y - x'_1 \beta_{1,0} = z'_1 \beta_2 + u,$$

where $E(uz_1) = 0$ and $E(uz_2) = 0$.

$$\text{E}(zu) = E((z_2)u) = 0$$

Aside: Overidentification tests

$$y = x' \beta + u \quad E(zu) = 0.$$

Dont make a hypothesis about β .

- Instead of testing whether $\beta = \beta_0$, we can check whether the validity condition $E(zu) = 0$ holds when $\dim(x_1) < \dim(z_2)$.
- Suppose $\dim(z_2) = l_2 > k_1 = \dim(x_1)$. We say there are $l_2 - k_1$ "overidentifying restrictions".
- Now we don't know β , but under $H_0 : E(zz')^{-1} E(zu) = 0$, any GMM estimator is consistent.
- Run the regression

$$y_i - x_i' \hat{\beta}_{GMM} = z_i' \gamma + \epsilon_i.$$

$$E(zu) = 0.$$

*vector of
coefficients of a
regression of
 u on z .*

Run a regression of \hat{u}_i on z_i -

$$y_i = x_i' \beta + u$$

$$y_i - x_i' \hat{\beta} = x_i' \gamma + \epsilon_i$$

Aside: Overidentification tests

$$\Rightarrow \hat{\gamma} = 0.$$

- Assuming homoskedasticity and using the 2SLS estimator yields that the F statistic is asymptotically distributed as $\chi^2_{I_2 - k_1}$ under H_0 . This is known as Sargan's test.
- Alternatives robust to heteroskedasticity exist (see Wooldridge text p134).
- If H_0 is rejected, either our model specification is wrong, one or more of the instruments is invalid, but it does not tell us which instrument is invalid.
- If we fail to reject, it could be that both instruments are invalid but cause similar bias! Can only reliably reject H_0 when some subset of instruments is known to be valid.

Back to inference on β

- Since $E(zu) = 0$, we have

$$y - x_1' \beta_{1,0} = z_1' \beta_2 + z_2' \gamma + u; \quad E(zu) = 0$$

where $\gamma = 0$. If $\gamma \neq 0$ either the specification is incorrect, z_2 is not valid or H_0 is false.

- Suppose $z_1 \in \mathbb{R}^{l_1}, z_2 \in \mathbb{R}^{l_2}$, where $l_1 + l_2 = l$.
- Run this regression and note that for some V :

$$\sqrt{n}(\hat{\gamma} - \gamma) \xrightarrow{d} \mathcal{N}(0, V),$$

so

$$n(\hat{\gamma} - \gamma)' \hat{V}^{-1} (\hat{\gamma} - \gamma) \xrightarrow{d} \chi^2_{l_2}.$$

↖ $\dim(\gamma) = \dim(z_2)$.

Conducting inference on β

- Recall that

$$\begin{aligned} E(zz')^{-1}E(zu(\beta_0)) &= E(zz')^{-1}E(zx'(\beta - \beta_0) + zu) \\ &= E(zz')^{-1}E(zx')(\beta - \beta_0). \end{aligned}$$

- If the instruments are not relevant $E(zx')$ isn't full rank, so can't distinguish between β and β_0 .
- If the instruments are weak, this can be done but power may be poor. AR test controls null rejection probability but does not correct for the fundamental issue of having little information about β from the instruments.

Conducting inference on β

- Critical value $\chi^2_{l_2, 1-\alpha}$ increasing in l_2 , so may expect finite sample power to be poor if there are many excluded instruments.
- Other alternatives under overidentification discussed in Andrews, Stock and Sun (2019).
- Test has power against $\beta \neq \beta_0$ but also against violations of overidentifying restrictions.
- Testing not just whether $\beta = \beta_0$, but whether in fact there exists a β satisfying model assumptions.

Many Instruments

$$\hat{\pi}_i = \left(\sum_i z_i z_i' \right)^{-1} \sum_i z_i v_i$$

$\text{cov}(u, v) \neq 0.$

- Consider the first stage

$$x = z'\pi + v; \quad E(zv) = 0.$$

- Estimate π by OLS to give fitted values

$$\hat{x}_i = z_i' \hat{\pi}.$$

If we knew π , the fitted values would be exogenous, but $\hat{\pi}$ depends on v , which is correlated with u .

Many Instruments

$$\hat{\beta}_{2SLS} = (X' P_Z X)^{-1} X' P_Z Y.$$

$$P_Z X = \hat{X} \leftarrow \begin{matrix} \text{matrix of} \\ \text{fitted} \\ \text{values} \end{matrix}$$

- The 2SLS estimator is given by

$$\begin{aligned}\hat{\beta}_{2SLS} - \beta &= (X' P_Z X)^{-1} X' P_Z U \\ &= \left(\frac{1}{n} \sum_{i=1}^n \hat{x}_i x_i' \right)^{-1} \frac{1}{n} \sum_{i=1}^n \hat{x}_i u_i\end{aligned}$$

- If $\hat{x} = z'\pi$ then \hat{x} is uncorrelated with the error since $E(zu) = 0$.
$$E(\hat{x}u) = E(\pi'zu) = \pi'E(zu) = 0.$$
- Asymptotically, 2SLS is consistent because $\hat{\pi} \xrightarrow{P} \pi$, but what about when the number of instruments is proportional to the sample size?

Many Instruments

- One correction suggested is Jackknife instrumental variables estimator:

$$\hat{\beta}_{JIVE} - \beta = \left(\frac{1}{n} \sum_{i=1}^n \hat{x}_i x_i' \right)^{-1} \frac{1}{n} \sum_{i=1}^n \hat{x}_i u_i,$$

where this time $\hat{x}_i = z_i' \hat{\pi}_{-i}$. $\hat{\pi}_{-i}$ is the OLS estimator in the first stage computed after dropping observation i .

- Because observations are iid, $z_i \hat{\pi}_{-i}$ is now uncorrelated with u_i .
- Chao et. al (2012) work out the details for JIVE showing consistency and asymptotic normality under heteroskedasticity with many instruments.

$$\xrightarrow{n} \frac{\dim(z)}{n} \xrightarrow{\text{distr}} \mathcal{N}(0, 1).$$

Many Instruments

$$K = \dim(z) \rightarrow \infty \text{ as } n \rightarrow \infty.$$

$$\frac{K}{n} \rightarrow \alpha \in (0, 1).$$

- For now, let's derive the inconsistency in OLS/2SLS assuming $\dim(x) = 1$ and

$$\frac{\dim(z)}{n} := \frac{K}{n} \rightarrow \alpha \in (0, 1).$$

$$y = \beta x + u$$

$$x = z' \pi + v.$$

- Note that

$$E(x^2) = E([z' \pi + v]^2) = E([z' \pi]^2) + E(v^2).$$

We could keep the instrument strength (concentration parameter fixed) by assuming $E([z' \pi]^2) \rightarrow H > 0$. What we actually need:

$$\frac{1}{n} \sum_{i=1}^n \pi' z_i z_i' \pi \xrightarrow{P} H.$$

$$\begin{aligned} \pi' z_i z_i' \pi \\ = (\pi' z_i)^2 \end{aligned}$$

$$y_{in} = \beta x_{in} + u_i \quad y, x, u \in \mathbb{R}.$$

$$x_{in} = z_{in}' \pi_n + v_i \quad z_{in} \in \mathbb{R}^{k_n}, v \in \mathbb{R}. \\ (\text{independent sample})$$

$$\xi_i = \begin{pmatrix} u_i \\ v_i \end{pmatrix}, \quad E(\xi_i | Z) = 0 \quad Z := \bigcup_{n=1}^{\infty} \{z_{in}\}_{i=1}^{k_n}$$

$$\text{Var}(\xi_i | Z) = \begin{pmatrix} \sigma_u^2 & \sigma_{uv} \\ \sigma_{uv} & \sigma_v^2 \end{pmatrix} \{y_{in}, x_{in}, z_{in}\}_{i=1}^{k_n}$$

$$Y = X\beta + U \quad Y, X, U \in \mathbb{R}^n \quad \beta \in \mathbb{R}$$

$$X = Z\pi + V \quad X \in \mathbb{R}^n, \quad Z \in \mathbb{R}^{N \times k_n} \quad \pi \in \mathbb{R}^{k_n} \\ V \in \mathbb{R}^n.$$

$$E(\|\xi_i\|^4 | Z) \leq B < \infty.$$

$$\xi_i \sim \text{iid.} \quad E((z_{in}' \pi_n)^2) \rightarrow H > 0.$$

$$\frac{1}{n} \sum \pi_n' z_{in} z_{in}' \pi_n \rightarrow^P H.$$

① Fixed number of strong instruments.

$$z_{in} = z_i \in \mathbb{R}^k \quad \pi \in \mathbb{R}^k.$$

$$E((z' \pi)^2) = H > 0.$$

$$\frac{1}{n} \sum (z_i' \pi)^2 \rightarrow^P E((z_i' \pi)^2) = H.$$

② "Many weak instruments":

$$E((z_n' \pi_n)^2) = \pi_n' E(z_n z_n') \pi_n$$

Assume each component of z_n has mean zero,
variance 1 and is uncorrelated with all other
components of z .

$$E(z_n z_n') = I_{Kn}$$

$$E((z_n' \pi_n)^2) = \pi_n' I_{Kn} \pi_n = \pi_n' \pi_n.$$

Weak instruments: $\pi_n = \underbrace{\left(\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}, \dots, \frac{1}{\sqrt{n}} \right)}_{Kn \text{ times.}}$

$$\pi_n' \pi_n = \frac{Kn}{n} \rightarrow \alpha \in \{0, 1\}.$$

$$\frac{1}{n} \sum_{i=1}^n (z_{in}' \pi_n)^2 - E((z_n' \pi_n)^2) \xrightarrow{P} 0 ?$$

$$X_n \xrightarrow{P} 0 \quad \text{if} \quad f(X_n^2) \rightarrow 0.$$

$$\xrightarrow{L^2} \Rightarrow \xrightarrow{P}$$

$$E \left(\left[\frac{1}{n} \sum (z_{in}' \pi_n)^2 - E[(z_n' \pi_n)^2] \right]^2 \right)$$

2

$$\begin{aligned}
 &= \frac{1}{n^2} E \left[\left(\sum_{i=1}^n (z_{in}' \pi_n)^2 - E[(z_{in}' \pi_n)^2] \right) \right] \\
 &= \frac{1}{n^2} \text{Var}(\text{Sum}) = \frac{1}{n^2} \sum_{i=1}^n \text{Var} \\
 &= \frac{1}{n^2} \sum_{i=1}^n \text{Var}((z_{in}' \pi_n)^2) \\
 &\stackrel{?}{\rightarrow} 0
 \end{aligned}$$

e.g. Let each component of z_{in} be $N(0, 1)$.

$$\pi_n = \underbrace{\left(\frac{1}{\sqrt{n}}, \dots, \frac{1}{\sqrt{n}} \right)}_{k_n \text{ times}}.$$

$$\begin{aligned}
 z_{in}' \pi_n &= \frac{1}{\sqrt{n}} \sum \text{Components of } z_{in} \\
 &\stackrel{d}{=} \frac{1}{\sqrt{n}} N(0, k_n) \\
 &\stackrel{d}{=} \sqrt{\frac{k_n}{n}} N(0, 1).
 \end{aligned}$$

$$(z_{in}' \pi_n)^2 \stackrel{d}{=} \frac{k_n}{n} \chi_1^2.$$

$$\text{Variance of } \frac{k_n}{n} \chi_1^2 = \frac{k_n^2}{n^2} \cdot \text{Var}(\chi_1^2) = \frac{2k_n^2}{n^2}.$$

$$\frac{1}{n^2} \sum \text{Var}(z_{in}' \pi_n) = \frac{1}{n^2} \cdot n \cdot \frac{2k_n^2}{n^2} = \frac{1}{n} \cdot \left(\frac{2k_n^2}{n^2} \right) \rightarrow 0.$$

$$\hat{\beta}^{OLS} - \beta = \left(\frac{1}{n} \sum x_{in} x_{in}' \right)^{-1} \frac{1}{n} \sum x_{in} u_i$$

$$x_{in} = z_{in}' \pi_n + v_i = \left(\frac{1}{n} \sum \pi_n' z_{in} z_{in}' \pi_n + \frac{1}{n} \sum \pi_n' z_{in} v_i + \frac{1}{n} \sum v_i^2 \right)^{-1} \times \left(\frac{1}{n} \sum \pi_n' z_{in} u_i + \frac{1}{n} \sum v_i u_i \right)$$

$\uparrow P$ $\uparrow P$
 \circ \circ
 σ_{uv}

$$WTS: P \left\{ \left| \frac{1}{n} \sum \pi_n' z_{in} u_i - 0 \right| > \varepsilon \mid Z \right\} \rightarrow 0$$

$$\leq \frac{1}{\varepsilon^2} \cdot E \left\{ \left| \frac{1}{n} \sum \pi_n' z_{in} u_i \right|^2 \mid Z \right\}$$

$$= \frac{1}{\varepsilon^2} \cdot \frac{1}{n^2} \sum_{i=1}^n E \left[(\pi_n' z_{in} u_i)^2 \mid Z \right]$$

$$= \frac{1}{\varepsilon^2} \cdot \frac{1}{n^2} \sum \underbrace{(\pi_n' z_{in})^2}_{= \sigma^2} \underbrace{E(u_i^2 \mid Z)}_{= \sigma^2} .$$

$$= \frac{\sigma^2}{\varepsilon^2} \cdot \frac{1}{n} \cdot \underbrace{\frac{1}{n} \sum_{i=1}^n (\pi_n' z_{in})^2}_{\rightarrow P H} = \left(\frac{\sigma^2}{\varepsilon^2} H + O_p(1) \right) \cdot \frac{1}{n}$$

$$= O_p(1), o(1) = o_p(1)$$

$$P\left(\underbrace{\left|\frac{1}{n} \sum_{i=1}^n z_i u_i - 0\right|}_{A_n} > \varepsilon \mid Z\right) \rightarrow^P 0.$$

$$P(A_n \mid Z) \rightarrow^P 0. \quad \text{Does } P(A_n) \rightarrow 0?$$

$$\text{Call } Y_n \equiv P(A_n \mid Z)$$

$$Y_n \rightarrow^P 0.$$

$$E(Y_n) \rightarrow 0 ?$$

$$P(A_n) = E(P(A_n \mid Z))$$

Yes: Vitali's convergence theorem.

$$\rightarrow 0.$$

$$\frac{1}{n} \sum_{i=1}^n T_{i,n} z_i u_i \rightarrow^P 0.$$

Y_n is a conditional probability and therefore bounded.

$$\text{Same holds for } \frac{1}{n} \sum_{i=1}^n T_{i,n} z_i v_i -$$

$$\hat{\beta}_{OLS} - \beta \rightarrow^P (H + \sigma_v^2)^{-1} \sigma_{uv} -$$

Many Instruments

- Next suppose $E(v|z) = 0$, $E(v^2|z) = \sigma_v^2$ and $E(v^4|z) \leq B$ for some constant B . Assume similar for u , and $E(uv|z) = \sigma_{uv}$. Then

$$\begin{aligned} \text{Var} \left(\frac{1}{n} \sum_{i=1}^n \pi' z_i v_i | Z \right) &= \frac{1}{n^2} \sum_{i=1}^n (\pi' z_i)^2 \text{Var}(\pi' z_i v_i | Z) \\ &= \sigma_v^2 \cdot \frac{1}{n^2} \sum_{i=1}^n \cancel{\text{Var}(\pi' z_i v_i | Z)} \xrightarrow{p} 0. \end{aligned}$$

Since $E(\pi' z_i v_i | Z) = 0$ and convergence in L^2 implies convergence in probability:

$$\frac{1}{n} \sum_{i=1}^n \pi' z_i v_i \xrightarrow{p} 0; \quad \frac{1}{n} \sum_{i=1}^n \pi' z_i u_i \xrightarrow{p} 0$$

Many Instruments

- Moreover:

$$\frac{1}{n} \sum_{i=1}^n x_i x'_i = \frac{1}{n} \sum_{i=1}^n \pi' z_i z'_i \pi + \frac{2}{n} \sum_{i=1}^n \pi' z_i v_i + \frac{1}{n} \sum_{i=1}^n v_i^2$$
$$\xrightarrow{P} H + \sigma_v^2$$

and

$$\frac{1}{n} \sum_{i=1}^n x_i u_i = \frac{1}{n} \sum_{i=1}^n \pi' z_i u_i + \frac{1}{n} \sum_{i=1}^n u_i v_i,$$

so

$$\hat{\beta}_{OLS} - \beta \xrightarrow{P} \frac{\sigma_{uv}}{H + \sigma_v^2}.$$

Many Instruments

Doesn't hold: Dimension is growing

- For 2SLS:

$$\begin{aligned}\hat{\beta}_{2SLS} - \beta &= (X' P_Z X)^{-1} X' P_Z U \\ &= \left(\frac{X' P_Z X}{n} \right)^{-1} \left[\frac{(Z\Pi + V)' P_Z U}{n} \right].\end{aligned}$$

Replace X with $Z + \Pi + V$
Since $Z' P_Z = Z'$, we have

$$\Pi' Z' P_Z U = \Pi' Z U.$$

$$\frac{(Z\Pi + V)' P_Z U}{n} = \frac{\Pi' Z U + V' P_Z U}{n} = \frac{1}{n} \sum_{i=1}^n \pi_i' z_i u_i + \frac{V' P_Z U}{n}.$$

- Usually the last term would converge in probability to zero,
but now the dimension of Z is growing with n .

0

Many Instruments

$$E(x_n) \rightarrow \mu \quad \text{Var}(x_n) \rightarrow 0 \\ \Rightarrow x_n \xrightarrow{P} \mu.$$

- We have

Scalar
↓

$$\begin{aligned} E(V' P_Z U | Z) &= \text{tr}(E(V' P_Z U | Z)) = E[\text{tr}(V' P_Z U) | Z] \\ &= \text{tr}(P_Z E(U V' | Z)) = E[\text{tr}(P_Z U V') | Z] \\ &= \text{tr}(P_Z \cdot \sigma_{uv} I_n) \quad \xrightarrow{\text{tr}} = \text{tr}[E(P_Z U V' | Z)] \\ &= \sigma_{uv} (\text{tr}(P_Z)) \\ &= \sigma_{uv} \left(\text{tr}([Z' Z]^{-1} Z' Z) \right) \\ &= \sigma_{uv} \cdot K_n \quad \xrightarrow{\text{tr}} \text{tr}(Z(Z' Z)^{-1} Z') \\ &= \text{tr}[(Z' Z)^{-1} Z' Z]. \end{aligned}$$

Therefore

$$E\left(\frac{V' P_Z U}{n} \middle| Z\right) = \sigma_{uv} \cdot \frac{K_n}{n} \rightarrow \sigma_{uv} \alpha. \quad = \text{tr}(I K_n)$$

$$\frac{V' P_Z U}{n} \xrightarrow{P} \sigma_{uv} \alpha.$$

↑ Vitali Convergence

Many Instruments

$$P\left(\left|\frac{V' P_Z U}{n} - \sigma_{uv}\alpha\right| > \varepsilon(z)\right) \rightarrow^P 0.$$

- Can also show that

$$Var\left(\frac{V' P_Z U}{n} \middle| Z\right) \xrightarrow{P} 0,$$

P 391 of Hansen

so

$$\frac{V' P_Z U}{n} \xrightarrow{P} \sigma_{uv}\alpha,$$

therefore the 2SLS estimator is inconsistent with many instruments.

How do I know Veratile of $V'P_2U$ is a well-defined object?

$$a = V'P_2'$$

$$V'P_2U = V'P_2'P_2U \quad b = P_2U -$$

Cauchy-Schwarz: $|a'b| \leq \|a\|\cdot\|b\|$

$$|V'P_2U| \stackrel{C-S}{\leq} (V'P_2V)^{\frac{1}{2}} |V'P_2U|^{\frac{1}{2}}$$

$$V'V = V'I_n V = V'(P_2 + M_2)V$$

$$(M_2 = I - P_2) = V'P_2V + V'M_2V.$$

$$= \underbrace{V'P_2'P_2V}_{\text{Sum of squares}} + \underbrace{V'M_2'M_2V}_{\text{Both } \geq 0}.$$

$$V'P_2V \leq V'V.$$

$$E(V'P_2U) \leq E((V'P_2V)^{\frac{1}{2}} (U'P_2U)^{\frac{1}{2}})$$

$$\leq E((V'V)^{\frac{1}{2}} \cdot (U'U)^{\frac{1}{2}})$$

$$\begin{aligned} E(XY) &\leq E(X^2)^{\frac{1}{2}} E(Y^2)^{\frac{1}{2}} \\ &\leq (n \cdot \sigma_V^2)^{\frac{1}{2}} (n \cdot \sigma_U^2)^{\frac{1}{2}} < \infty. \end{aligned}$$

Questions?

$$\begin{aligned}\text{Var}(V' P_2 V) &\leq E((V' P_2 V)^2) \\ &\leq E((V' V) \cdot (U' U)) \\ &= n^2 E(V^2 U^2) \\ &\leq n^2 E(V^4)^{\frac{1}{2}} E(U^4)^{\frac{1}{2}}. \\ &< \infty.\end{aligned}$$

Used convergence in L^2 to show that 2SLS in general is inconsistent with a system containing many weak instruments. Alternative: LIML / Modifications of 2SLS which are consistent in this framework.

Causal parameters from estimators

- Return to the model

$$y_i = \beta_0 + \beta_1 d_i + u_i,$$

where y_i is the observed outcome of individual i and $d_i = 1$ if individual i is treated and 0 otherwise.

- A causal interpretation of this model implies that the *ceteris paribus* effect of the treatment, β_1 , is the same for everybody.
- On the other hand, since d_i is binary we can always write

$$y_i = b_0 + b_1 d_i + \epsilon_i; \quad E(\epsilon_i | d_i) = 0.$$

- The latter model does not necessarily have a causal interpretation, but the parameter b_1 may be of interest.
- The error terms u_i and ϵ_i are interpreted differently: u_i contains unobserved determinants of y_i , ϵ_i is the projection residual. It is not necessarily true that $E(u_i | d_i) = 0$.

Causal parameters from estimators

- We may write

$$y_i = y_{i0} d_i + y_{i1} (1 - d_i) = y_{i0} + d_i (y_{i1} - y_{i0}),$$

$\overbrace{\quad \quad \quad}^{= c}$

where y_{i0}, y_{i1} are the potential outcomes for individual i .

- y_{i0} is the outcome we would observe if the individual is not treated, and y_{i1} the outcome if the individual is treated.
- The difference $y_{i1} - y_{i0}$ is the individual treatment effect.
- The average across all individuals, $E(y_{i1} - y_{i0})$, is the average treatment effect (ATE).

$$y_i = y_{i0} + d_i (y_{i1} - y_{i0})$$

$\beta_i \rightarrow$ Random coefficient .

Causal parameters from estimators

Estimate parameters using
OLS/IV, what do I
end up with?

- So far we have argued that if treatment is randomly assigned then the OLS estimate gives the ATE. Adding the assumption that $y_{i1} - y_{i0} = c$ allowed us to estimate c and to interpret the regression as a causal model.
- One reason to be interested in 2SLS is that while we may ask for a parameter of interest (such as the ATE) and then look at how to estimate it, we can go the other way around: Use the estimation method anyway and then ask what parameter it identifies.
- For 2SLS/IV the answer (under some assumptions) is a Local Average Treatment Effect.

Causal parameters from estimators

$$y_i = y_0(1-\mathbb{D}) + y_1 \mathbb{D} = y_0 + \mathbb{D}(y_1 - y_0)$$

- Suppose we allow for heterogeneous treatment effects, so that $y_{i1} - y_{i0} = \beta_{i1}$, and the baseline outcome $\beta_{i0} = y_{i0}$ may also differ across individuals. Then

$$y_i = \beta_{i0} + \beta_{i1} d_i.$$

- Such a model is called a random coefficients model, because the fact we are sampling individuals at random implies that the actual value of β_{i0}, β_{i1} drawn is random also.

Causal parameters from estimators

- The assumption of random assignment to treatment, $(y_{i0}, y_{i1}) \perp\!\!\!\perp d_i$, implies that in the model

$$y_i = b_0 + b_1 d_i + \epsilon_i, \quad E(\epsilon_i | d_i) = 0,$$

we have

$$\begin{aligned} b_1 &= E(y_i | d_i = 1) - E(y_i | d_i = 0) \\ &= E(y_{i1} | d_i = 1) - E(y_{i0} | d_i = 0) \\ &= E(y_{i1} - y_{i0}). \end{aligned}$$

- Random assignment implies OLS estimate of b_1 is a consistent estimate of the ATE.

Causal parameters from estimators

Want some precision of $y_1 - y_0$ ^(all) eg. $E(y_1 - y_0)$
 $E(y_1 - y_0 | D \approx)$

- If, on the other hand, there is selection into treatment, we obtain merely

$$b_1 = E(y_{i1} | d_i = 1) - E(y_{i0} | d_i = 0).$$

- Want to identify some feature of the distribution of $y_{i1} - y_{i0}$.
- Suppose we have a binary instrument $z \in \{0, 1\}$. For example, y might be worker output, $d = 1$ if the worker elects job training and 0 otherwise, and $z = 1$ if the worker is offered job training and 0 otherwise. z may be independent of (y_0, y_1) even if d isn't.

Binary treatment $D \in \{0, 1\}$ $y = D y_1 + (1-D)y_0$

\hookrightarrow Taken computing course? $\mathbb{Y}(N)$.
 y - observed earnings.

$y_0, y_1 \cancel{\perp\!\!\! \perp} D$

Selection into treatment, so can't identify the ATE using a naive comparison.

If we randomly offer some individuals a subsidy to take the computing course, perhaps shifting their treatment choice, can we recover some kind of ATE?

Binary instrument $Z \in \{0, 1\}$ - Subsidy? $\mathbb{Y}(N)$.

Gives "potential treatments": what treatment would I select with the subsidy (D_1) and without (D_0)

$$D = D_1 Z + D_0(1-Z).$$

\hookrightarrow Observed Treatment.

Exogeneity (from hom. TEs): $(y_0, y_1) \perp\!\!\! \perp Z$.

We need: Strong Exogeneity: $(D_0, D_1, y_0, y_1) \perp\!\!\! \perp Z$.

$$Y = y_0 + D(y_i - y_0)$$

$$= E(y_0) + D(y_i - y_0) + (y_0 - E(y_0))$$

$$= \beta_0 + BD + U$$

$(y_i - y_0)$ is a random coefficient.

$$\hat{\beta}^{IV} = \frac{\widehat{\text{Cov}}(Y, Z)}{\widehat{\text{Cov}}(D, Z)} \xrightarrow{P} \frac{\text{Cov}(Y, Z)}{\text{Cov}(D, Z)}.$$

$$\begin{aligned} \text{Cov}(Y, Z) &= \text{Cov}(y_0 + BD, Z) \\ &= \text{Cov}(BD, Z) \\ &= E(BD(Z - E(Z))) \end{aligned}$$

$$\hat{\beta}^{IV} \xrightarrow{P} E\left(B \frac{D(Z - E(Z))}{\text{Cov}(D, Z)}\right)$$

↑ ↑
Treatment effect weight

If B is constant, this simplifies to B .

Note: "Intention to treat" parameter $E(Y|Z=1) - E(Y|Z=0)$
 is the ATE of offering the subsidy on outcome, but it
 is not the effect of taking the cause.

Who is affected by the subsidy? Likely those who were previously indifferent between taking the course and not.

	Complies		Always Takes
	$D_0 = 0, D_1 = 1$		$D_0 = 1, D_1 = 1$
Never-Takes			Depiers
$D_0 = 0, D_1 = 0$			$D_0 = 1, D_1 = 0$

Consider reduced form: Regression of Y on Z .

$$Y = \beta_0 + \beta_1 Z + \varepsilon. E(\varepsilon) = E(Z\varepsilon) = 0.$$

$\beta_1 = \frac{\text{Cov}(Y, Z)}{\text{Var}(Z)}$ is the slope coefficient.

$$\beta_0 = E(Y | Z=0)$$

$$\beta_1 = E(Y | Z=1) - E(Y | Z=0)$$

$$\begin{aligned} E(Y | Z=1) &= E(y_1 D + y_0(1-D) | Z=1) \\ &= E(y_1 D_1 + y_0(1-D_1) | Z=1) \\ &= E(y_0 + D_1(y_1 - y_0)). \end{aligned}$$

$$E(Y|Z=0) = E(y_0 + D_0(y_1 - y_0)) .$$

$$\Rightarrow \frac{\text{Cov}(Y, Z)}{\text{Var}(Z)} = E([D_i - D_o] B) .$$

First Stage: $\frac{\text{Cov}(D_i, Z)}{\text{Var}(Z)} = E(D|Z=1) - E(D|Z=0)$

$$= E(D_i|Z=1) - E(D_o|Z=0)$$

(Strong Exogeneity) $= E(D_i - D_o)$.

Average effect in population of offering subsidy on participation.

Need experiment to offer subsidies randomly, not just to those who experimenter thinks will take the class anyway.
 Otherwise D is correlated with Z (Maybe even $D=Z$)
 and ... $y_0, y_1 \cancel{\perp\!\!\!\perp} Z$. Would cause us to observe exactly the same individuals in treatment and control as if there were no experiment. Only non-zero if $D_i \neq D_o$

$$E[Y(Z=1) - E[Y(Z=0)] = E[D_i - D_o]B]$$

$$= E(B | D_i=1, D_o=0) P(D_i=1, D_o=0)$$

$$- E(B | D_o=1, D_i=0) P(D_o=1, D_i=0)$$

Weighted average of ATEs for complies and defiers, with a negative weight for defiers. How to interpret this? $\text{Cov}(Y, Z)$ could be negative even if both ATEs are positive! Assume away defiers:

Monotonicity: $P(D_1 \geq D_0) = 1$ (No defiers).

or $P(D_0 \geq D_1) = 1$
 Reasonable with button subsidy, but not if the instant values are not naturally ordered - e.g. August + Imbens (1994). Applicants for social program screened by two officials, who likely have different admission rates. If official A accepts with prob p_A and official B with prob $p_B > p_A$, then official B must accept any candidate admitted by official A.

$$P(D_1 = 1, D_0 = 0) > 0.$$

$$P(D_1 = 0, D_0 = 1) > 0.$$

$$\text{Monotonicity} \Rightarrow P(D_1 = 0, D_0 = 1) = 0.$$

$$\Rightarrow E(Y|Z=1) - E(Y|Z=0) = E(B|D_0=0, D_1=1)P(D_0=0, D_1=1)$$

$$\text{First Stage: } \frac{\text{Cov}(D_1, Z)}{\text{Var}(Z)} = E(D_1 - D_0)$$

$$(\text{Monotonicity}) = 1 \cdot P(D_0 = 0, D_1 = 1).$$

$$\text{Provided } P(\text{complier}) > 0: \beta_{IV} = \frac{\text{Cov}(Y, Z)}{\text{Cov}(D_1, Z)} = E(B|D_0=0, D_1=1)$$

ATE measured depends on the instrument, since different experiments (with different Z) will impact treatment choice differently. If we are considering a larger scale reduction of the cost of buying computer server, then we might care about this LATE, but if we care about the average effect in the population of taking computing vs. not, the local average treatment effect will not measure this (necessarily).

Can determine proportions of complies, always takes and never takes:

By monotonicity: if $Z=1$ but $D=0$, must be a never-taker

$Z=0$ but $D=1$, must be always takes

$$P(\text{nev. take}) = P(D=0 | Z=1) -$$

$$P(\text{alw. take}) = P(D=1 | Z=0) .$$

$$P(\text{complies}) = P(D=1 | Z=1) - P(D=1 | Z=0)$$

$$\stackrel{\dagger}{P}(\{\text{at, cp}\}) - P(\{\text{at}\}) = P(\text{complies})$$

$$= P(D=0 | Z=0) - P(D=0 | Z=1) .$$

Special case: "One-sided non-compliance". — Individuals can only take treatment if offered it: No always-takers. Then we have only compliers and never-takers.

Typically, $D=1$ contains at, cp, defies, but without defies and always takes. $D=1$ corresponds to the set of compliers. So $LATE = ATT$.

If everyone were a complier, $=ATE$, but then the experiment would completely determine treatment status.

Causal parameters from estimators

- Now we not only have potential outcomes but also potential treatments:

$$d_i = d_{i1}z_i + d_{i0}(1 - z_i),$$

where the potential outcome $d_{i1} = 1$ if individual i would take job training when they are offered it and 0 otherwise. $d_{i0} = 1$ if individual i would take job training when they are not offered it, and 0 otherwise.

- Question: “OLS produces $b_1 = E(y_{i1}|d_i = 1) - E(y_{i0}|d_i = 0)$, which is not a feature of $y_{i1} - y_{i0}$, but what if I do 2SLS/IV with my instrument?”

Causal parameters from estimators

- First we will make some relevance and validity assumptions analogous to those in the linear model.
- Validity: $(y_0, y_1, d_0, d_1) \perp\!\!\!\perp z$.
- Relevance: Two assumptions
 - Monotonicity assumption: $P(d_1 \geq d_0) = 1$. This implies

$$\begin{aligned} \text{Cov}(d, z) &= E(dz) - E(d)E(z) \\ &= P(d = 1, z = 1) - P(d = 1)P(z = 1) \\ &= [P(d = 1|z = 1) - P(d = 1|z = 0)]P(z = 1)P(z = 0) \\ &= [P(d_1 = 1) - P(d_0 = 1)]Var(z) \\ &= P(d_1 > d_0)Var(z) \geq 0. \end{aligned}$$

- $P(d_0 \neq d_1) > 0$. This says that the instrument must alter treatment choice for some positive fraction of individuals. Together with monotonicity, this implies $P(d_1 > d_0) > 0$.

Causal parameters from estimators

- Note that the first stage regression produces fitted values

$$\hat{d}_i = \hat{\gamma}_0 + \hat{\gamma}_1 z_i,$$

where

$$\hat{\gamma}_1 = \frac{\hat{Cov}(d, z)}{\hat{Var}(z)} \xrightarrow{p} P(d_1 > d_0) > 0.$$

- The second stage regression runs OLS on

$$y_i = \alpha + \beta \hat{d}_i + v_i.$$

- This gives:

$$\hat{\beta}_{2SLS} = \frac{\hat{Cov}(y, z) / \hat{Var}(z)}{\hat{Cov}(d, z) / \hat{Var}(z)} \xrightarrow{p} \frac{Cov(y, z)}{Cov(d, z)}.$$

Causal parameters from estimators

- A similar argument shows that

$$\begin{aligned}\text{Cov}(y, z) &= [\mathbb{E}(y|z = 1) - \mathbb{E}(y|z = 0)] \text{Var}(z) \\&= \mathbb{E}(y_1 d_1 + y_0 (1 - d_1) | z = 1) \text{Var}(z) \\&\quad - \mathbb{E}(y_1 d_0 + y_0 (1 - d_0) | z = 0) \text{Var}(z) \\&= \mathbb{E}(y_1 d_1 + y_0 (1 - d_1)) \text{Var}(z) \\&\quad - \mathbb{E}(y_1 d_0 + y_0 (1 - d_0)) \text{Var}(z) \\&= \mathbb{E}([y_1 - y_0] [d_1 - d_0]) \text{Var}(z) \\&= \mathbb{E}(y_1 - y_0 | d_1 > d_0) \mathbb{P}(d_1 > d_0) \text{Var}(z).\end{aligned}$$

Causal parameters from estimators

- It follows that

$$\hat{\beta}_{2SLS} \xrightarrow{p} \frac{\text{Cov}(y, z)}{\text{Cov}(d, z)} = \mathbb{E}(y_1 - y_0 | d_1 > d_0).$$

- This parameter is the average treatment effect among those who would be switched from no treatment to treatment if the value of the instrument changed.
- In the job training example, it is the average effect among those who would elect training if offered it and not elect it otherwise.
- This parameter may or may not be of interest, and this interpretation depends on the monotonicity assumption.
- The interpretation changes if the instrument is changed (unless, of course, $y_1 - y_0$ is constant).

Questions?

$$\text{Last time } Y = y_0 + D(y_1 - y_0) \quad D \in \{0, 1\},$$

$$D = D_0 + Z(D_1 - D_0) \quad Z \in \{0, 1\}.$$

Without covariates and under strong exogeneity + Monotonicity:

$$\xrightarrow{\text{Estimand}} \beta^{IV} = E(Y_i - y_0 | D_i=1, D_0=0) \quad \text{ATE for "compliers".}$$

Now: Allow for covariates W which may help in non-experimental settings to justify independence of instrument from potential outcomes / treatment.

Conditional Strong Exogeneity: $(y_0, y_1, D_1, D_0) \perp\!\!\!\perp Z | W$

conditional Monotonicity: $P(D_1 > D_0 | W=w) = 1$.

Repeat previous process for each value $W=w$:

$$\beta^{IV}(w) = \text{LATE}(w) = E(Y_i - y_0 | D_i=1, D_0=0, W=w).$$

$$\xrightarrow{\frac{\text{Cov}(Y, Z | W=w)}{\text{Cov}(D, Z | W=w)}}.$$

Consider the probability limit of $\hat{\beta}^{IV}$ conditional on fixing $W=w$.

Reduced form for $W=w$: $Y = \beta_0 + \beta_1 Z + U$.

$$Y 1_{W=w} = (\beta_0 + \beta_1 Z + U) 1_{W=w}.$$

$$= \beta_0 1_{W=w} + \beta_1 Z 1_{W=w} + U 1_{W=w}.$$

$$\textcircled{1} E(ZU 1_{W=w}) = E(Z 1_{W=w} \cdot U 1_{W=w}) = 0.$$

$$\textcircled{2} E(U 1_{W=w}) = E(1_{W=w} \cdot U 1_{W=w}) = 0.$$

$$\textcircled{2} \text{ Plug in } U: E((Y - \beta_0 - \beta_1 Z) \mathbb{1}_{W=w}) = 0.$$

$$E(Y - \beta_0 - \beta_1 Z | W=w) P(W=w) = 0.$$

$$E(Y - \beta_0 - \beta_1 Z | W=w) = 0.$$

$$\Rightarrow \beta_0 = E(Y - \beta_1 Z | W=w).$$

$$\textcircled{1} \quad E(Z(Y - \beta_0 - \beta_1 Z) | W=w) = 0.$$

Plug in β_0 to \textcircled{1}:

$$E(Z(Y - E(Y|W=w) - \beta_1(Z - E(Z|W=w))) | W=w) = 0$$

$$\beta_1^{(w)} = \frac{\text{Cov}(Y, Z | W=w)}{\text{Var}(Z | W=w)}$$

Do the same in the first stage: First stage coefficient in a regression of D on a constant and Z (conditional on $W=w$) is

$$\pi_1^{(w)} = \frac{\text{Cov}(D, Z | W=w)}{\text{Var}(Z | W=w)} . \quad \beta_{IV}^{(w)} = \frac{\text{Cov}(Y, Z | W=w)}{\text{Cov}(D, Z | W=w)}$$

*Same logic as
last class using
new assumptions.*

Cure of dimensionality makes it difficult to estimate these covariances at each value of W , so covariates are usually included in a linear fashion.

Saturated specification: $\gamma = \sum_{w \in W} \beta_{0,w} I(w=w) + \sum_{w \in W} \beta_{1,w} Z I(w=w) + U$

\downarrow
in W

Reduced form: $D = \sum_{w \in W} \pi_{0,w} I(w=w) + \sum_{w \in W} \pi_{1,w} Z I(w=w) + V$

\uparrow \uparrow
 $\epsilon(D|Z=0, w=w)$ $\epsilon(D, Z=1, w=w)$

$\frac{\beta_{1,w}}{\pi_{1,w}} = \text{LATE}(w)$. Since D is binary, what do we get
if we estimate

2nd stage: $\gamma = \sum_{w \in W} \beta_{0,w} I(w=w) + \beta_1 D + U?$

If we use this 2nd stage and the saturated 1st stage, β_1 turns out to be a positively weighted average of LATE(w) over w .

Tanbens + Augrist (1988) "Sature and weight".

Note that the first stage is in fact a non-parametric regression:

We are finding $E(D|Z, w)$.

Most applications of 2SLS do not include a saturated specification in covariates in 2nd stage or 1st stage, so do we still get a positively weighted average of LATEs as in Tanbens + Augrist (1988)?

No.

$$\gamma = \beta_0 + \beta_1 D + w' \beta_2 + U$$

$$D = \pi_{0,w} + \pi_{1,w} Z + w' \pi_2 + V$$

Can estimate reduced form $\gamma = \gamma_0 + \gamma_1 Z + w' \gamma_2 + \epsilon$

Frisch-Waugh: $Z = \delta_0 + w' \delta_2 + \tilde{Z}$ $E(Z) = E(\tilde{Z}) = 0$.

- Regress γ on \tilde{z} (no constant).
- Do the same in 1st stage (regress $D\tilde{z}$ on \tilde{z} (no constant)).

$$\frac{\gamma_1}{\pi_1} = \beta_{IV} = \frac{E(Y\tilde{z}) / \text{Var}(\tilde{z})}{E(D\tilde{z}) / \text{Var}(\tilde{z})} = \frac{E(Y\tilde{z})}{E(D\tilde{z})}$$

$$\begin{aligned} E(Y\tilde{z}) &= E(E(Y\tilde{z}|w)) \\ &= E(\text{Cov}(Y, \tilde{z}|w) + E(Y|w)E(\tilde{z}|w)) \end{aligned}$$

$\text{Cov}(x,y) = E(XY) - E(X)E(Y)$

$$\begin{aligned} \text{Cov}(Y, \tilde{z}|w) &= \text{Cov}(\gamma_1 z - \delta_0 - \delta_1 w | w) \\ &= \text{Cov}(Y, z|w). \end{aligned}$$

$$E(\text{Cov}(Y, z|w)) = E(\text{LATE}(w) \text{Cov}(D, z|w))$$

$$\beta_{IV} = \frac{E(Y\tilde{z})}{E(D\tilde{z})} = E(\text{LATE}(w) \cdot \underbrace{\frac{\text{Cov}(D, z|w)}{E(D\tilde{z})}}_{\text{Are these weights positive?}}) + E\left(\frac{E(Y|w)E(\tilde{z}|w)}{E(D\tilde{z})}\right)$$

$$\frac{\text{Cov}(D, z|w)}{\text{Var}(z|w)} = P(D_0=0, D_1=1|w) > 0.$$

$E(D\tilde{z})$ is positive if the first stage is specified with a "rich enough" set of covariates.

Following Blandford et al. (2022). Proposition 10.

Monotonicity condition says $P(D_1=1|w=w) > P(D_0=1|w=w)$.

In the specification we are currently considering, the first stage is "monotonicity-correct" if $\pi_1 > 0$ in

$$D = \pi_0 + \pi_1 Z + w^\top \pi_2 + V.$$

This will hold if the covariates w are included in a saturated fashion (so each value of w has its own coefficient).

Consider 2nd term $E(E(Y|w)E(\tilde{Z}|w))$

$$\begin{aligned} E(\tilde{Z}|w) &= E(Z|w) - E(\delta_0 + \delta_1' w | w) \\ &= E(Z|w) - (\delta_0 + \delta_1' w) \leftarrow \text{BLP}(Z|w) \\ &= \underbrace{E(Z|w) - L(w)}_{\substack{\text{best linear predictor of } Z|w. \\ \text{Denote by } L(w).}} \end{aligned}$$

$L(w)$ is the best linear approximation

to $E(Z|w)$, so orthogonal to constant and w .

$$E(Y|w) = E(y_0|w) + E(D(y_i - y_0)|w)$$

$$\begin{aligned} \text{Assume } E(y_0|w) &= \alpha_0 + \alpha_1 w. \text{ Then } E(\tilde{Z}|w) E(y_0|w) \\ &= E(\tilde{Z}|w) (\alpha_0 + \alpha_1 w) \end{aligned}$$

$$\text{So: } E(E(\tilde{Z}|w) E(y_0|w)) = 0. \quad \text{B}$$

$$\text{So: } E(E(\tilde{Z}|w) E(Y|w)) = E(E(\tilde{Z}|w) E(\overline{D(y_i - y_0)}|w))$$

$$\begin{aligned} E(DB|w) &= E([D_0 + Z(D_i - D_0)]B|w) \\ &= E(D_0 B|w) + E(Z(D_i - D_0) B|w) \end{aligned}$$

$$\begin{aligned}
 & \text{ATE per always-takes, conditional on } w. \quad E(B | D_0=1, w) P(D_0=1 | w) \\
 & + E(ZB | w, D_0=0, D_1=1) P(D_0=0, D_1=1 | w) \\
 & = \text{ATE(lat, } w) P(\text{lat} | w) \\
 & + E(ZB | w, D_0=0, D_1=1) P(\text{cp} | w) \\
 & \hookrightarrow E(Z | w) E(B | w, D_0=0, D_1=1) \\
 & \hookrightarrow \text{LATE}(w)
 \end{aligned}$$

$$\begin{aligned}
 E(E(Z | w) E(Y | w)) &= E(E(\tilde{Z} | w) E(\widehat{D(y_1 - y_0)} | w)) \\
 &= E\left(E(\tilde{Z} | w) \cdot [ATE(\text{lat}, w) P(\text{lat} | w) + E(Z | w) LATE(w) P(\text{cp} | w)]\right) \\
 \text{Notation: } \beta^{IV} &= \downarrow + E(LATE(w) \cdot \text{cov}(D_z | w)) \\
 &= E_{LATE(w)} \left[\text{cov}(D_z | w) + E(\tilde{Z} | w) E(Z | w) P(\text{cp} | w) \right] \\
 &\quad + E(ATE(\text{lat}, w) P(\text{lat} | w) E(\tilde{Z} | w))
 \end{aligned}$$

Weighted average of LATEs and ATEs for always-takers. Both the weights in the LATE expectation and in the always-takers expectation can be negative. Coefficients in 1st expectation have the same sign.

as $1 - L(w)$, so can be negative.
weights on $A\tilde{E}(\text{lat}, w)$ will be positive and negative
since $E(\tilde{z}|w)$ is a residual.