

Instrumental Variables

Department of Economics
University of Wisconsin-Madison

September 27, 2016

Treatment Effects

Throughout the course we will focus on the “Treatment Effect Model”

For now take that to be

$$Y_i = \alpha T_i + X_i' \beta + u_i$$

α measures the causal effect of T_i on Y_i .

We don't want to just run OLS because we are worried that T_i is not randomly assigned, that is that T_i and u_i are correlated.

There are a number of different reasons that might be true-I think the main thing that we are worried about in the treatment effects literature is omitted variables.

In this course we are going to think about a lot of different ways of dealing with this potential problem and estimating α .

Instrumental Variables

Lets start with instrumental variables

I want to think about three completely different approaches for estimating α

The first is the GMM approach and the second two will come from a Simultaneous Equations framework

To justify OLS we would need

$$E(T_i u_i) = 0$$

$$E(X_i u_i) = 0$$

The focus of IV is to try to relax the first assumption

(There is much less concern about the second)

Lets suppose that we have an instrument Z_i for which

$$E(Z_i u_i) = 0$$

and we continue to assume that

$$E(X_i u_i) = 0$$

We also will stick with the exactly identified case (1 dimensional Z_i)

Define

$$Z_i^* = \begin{bmatrix} Z_i \\ X_i \end{bmatrix}, X_i^* = \begin{bmatrix} T_i \\ X_i \end{bmatrix}, B = \begin{bmatrix} \alpha \\ \beta \end{bmatrix}$$

Then we know

$$E \left[Z_i^* \left(Y_i - X_i^{*'} B \right) \right] = 0$$

Turning this into a GMM estimator we get

$$\begin{aligned} 0 &= \frac{1}{N} \sum_{i=1}^N Z_i^* \left(Y_i - X_i^{*'} \hat{B}_{IV} \right) \\ &= \frac{1}{N} \sum_{i=1}^N Z_i^* Y_i - \left(\frac{1}{N} \sum_{i=1}^N Z_i^* X_i^{*'} \right) \hat{B}_{IV} \end{aligned}$$

which we can solve as

$$\hat{B}_{IV} \equiv \left(\frac{1}{N} \sum_{i=1}^N Z_i^* X_i^{*'} \right)^{-1} \frac{1}{N} \sum_{i=1}^N Z_i^* Y_i$$

And that is IV.

Consistency Notation

I will use \approx to mean asymptotically equivalent. Typically the phrase

$$A_n \approx B_n$$

means that

$$(A_n - B_n) \xrightarrow{p} 0$$

So I want to be formal in this sense

However I will not be completely formal as this could also mean almost sure convergence, convergence in distribution or some sort of uniform convergence. The differences between these are not relevant for anything we will do in class.

Consistency of IV

$$\begin{aligned}\hat{B}_{IV} &= \left(\frac{1}{N} \sum_{i=1}^N Z_i^* X_i^{*'} \right)^{-1} \frac{1}{N} \sum_{i=1}^N Z_i^* Y_i \\ &= \left(\frac{1}{N} \sum_{i=1}^N Z_i^* X_i^{*'} \right)^{-1} \frac{1}{N} \sum_{i=1}^N (Z_i^* X_i^{*'} B + Z_i^* u_i) \\ &= B + \left(\frac{1}{N} \sum_{i=1}^N Z_i^* X_i^{*'} \right)^{-1} \frac{1}{N} \sum_{i=1}^N Z_i^* u_i \\ &\approx B\end{aligned}$$

since

$$\frac{1}{N} \sum_{i=1}^N Z_i^* u_i \approx 0$$

So (assuming iid sampling) this only took two assumptions.
The moment conditions and the fact that you can invert

$$E \left(Z_i^* X_i^{*'} \right)$$

As we will discuss this assumption is typically a bigger deal
than in OLS

Furthermore we get (Huber-White) standard errors from

$$\sqrt{N}(\hat{B}_{IV} - B) = \left(\frac{1}{N} \sum_{i=1}^N Z_i^* X_i^{*'} \right)^{-1} \left[\frac{1}{\sqrt{N}} \sum_{i=1}^N Z_i^* u_i \right]$$

Under standard conditions:

$$\frac{1}{\sqrt{N}} \sum_{i=1}^N Z_i^* u_i \approx N(0, E(u_i^2 Z_i^* Z_i^{*'}))$$

so

$$\sqrt{N}(\hat{B}_{IV} - B) \approx N\left(0, E(Z_i^* X_i^{*'})^{-1} E(u_i^2 Z_i^* Z_i^{*'}) E(X_i^* Z_i^{*'})^{-1}\right)$$

We approximate this by

$$E \left(Z_i^* X_i^{*'} \right) \approx \frac{1}{N} \sum_{i=1}^N Z_i^* X_i^{*'}$$

$$E \left(u_i^2 Z_i^* Z_i^{*' } \right) \approx \sum_{i=1}^N \hat{u}_i^2 Z_i^* Z_i^{*'}$$

Biasedness of IV

It turns out it will be easiest to think of IV in terms of consistency, it is generally biased.

I will first discuss why IV is biased, but then we will go to the more important case of calculating the asymptotic bias of IV or OLS

First lets think about why OLS is unbiased.

Assume that

$$E(u_i | X_i) = 0$$

Then

$$\begin{aligned} E(\widehat{B}_{OLS}) &= B + E\left(\left(\frac{1}{N} \sum_{i=1}^N X_i X_i'\right)^{-1} \frac{1}{N} \sum_{i=1}^N X_i u_i\right) \\ &= B + E\left(E\left[\left(\frac{1}{N} \sum_{i=1}^N X_i X_i'\right)^{-1} \frac{1}{N} \sum_{i=1}^N X_i u_i \mid X_1, \dots, X_N\right]\right) \\ &= B + E\left(\left(\frac{1}{N} \sum_{i=1}^N X_i X_i'\right)^{-1} \frac{1}{N} \sum_{i=1}^N X_i E[u_i \mid X_1, \dots, X_N]\right) \\ &= B \end{aligned}$$

The same trick does not work for IV because in general

$$E\left(\frac{1}{\bar{X}}\right) \neq \frac{1}{E(X)}$$

even though

$$\frac{1}{\bar{X}} \approx \frac{1}{E(X)}$$

This is neither deep nor important-just a fact you should probably be aware of

So now we write the assumption as

$$E(u_i | Z_i^*) = 0$$

It turns out that an analogous conditional expectation argument does not work. What do we condition on Z^* or (Z^*, X^*) ?

First consider Z^*

$$\begin{aligned} E(\widehat{B}_{IV}) &= B + E \left(\left(\frac{1}{N} \sum_{i=1}^N X_i^* Z_i^{*'} \right)^{-1} \frac{1}{N} \sum_{i=1}^N Z_i^* u_i \right) \\ &= B + E \left(E \left[\left(\frac{1}{N} \sum_{i=1}^N X_i^* Z_i^{*'} \right)^{-1} \frac{1}{N} \sum_{i=1}^N Z_i^* u_i \middle| Z_1^*, \dots, Z_N^* \right] \right) \end{aligned}$$

Can't bring the expectation inside so we are stuck here

and conditioning on (Z^*, X^*)

$$\begin{aligned} E(\widehat{B}_{IV}) &= B + E \left(\left(\frac{1}{N} \sum_{i=1}^N X_i^* Z_i^{*'} \right)^{-1} \frac{1}{N} \sum_{i=1}^N Z_i^* u_i \right) \\ &= B + E \left(E \left[\left(\frac{1}{N} \sum_{i=1}^N X_i^* Z_i^{*'} \right)^{-1} \frac{1}{N} \sum_{i=1}^N Z_i^* u_i \middle| Z^*, X^* \right] \right) \\ &= B + E \left(\left(\frac{1}{N} \sum_{i=1}^N X_i^* Z_i^{*'} \right)^{-1} \frac{1}{N} \sum_{i=1}^N Z_i^* E[u_i | Z^*, X^*] \right) \end{aligned}$$

but $E[u_i | Z^*, X^*] \neq 0$ (in general)

So that didn't work so well, we will need to think more about the asymptotic bias

Before that I want to take a detour and discuss partitioned regression which will turn out to be really useful for this and we will use it many times in this course

Partitioned Regression

Think about the standard regression model (in large matrix notation)

$$Y = X_1\beta_1 + X_2\beta_2 + u$$

We will define

$$M_2 \equiv I - X_2 (X_2' X_2)^{-1} X_2'$$

Two facts about M_2

Fact 1: M_2 is idempotent

$$\begin{aligned}M_2 M_2 &= \left(I - X_2 (X_2' X_2)^{-1} X_2' \right) \left(I - X_2 (X_2' X_2)^{-1} X_2' \right) \\&= I - 2X_2 (X_2' X_2)^{-1} X_2' + X_2 (X_2' X_2)^{-1} X_2' X_2 (X_2' X_2)^{-1} X_2' \\&= I - X_2 (X_2' X_2)^{-1} X_2' \\&= M_2\end{aligned}$$

Fact 2: $M_2 Y$ is the Residuals from Regression

For any potential dependent variable (say Y), $M_2 Y$ is the residuals I would get if I regressed Y on X_2

To see that let the regression coefficients be \hat{g} and generically let \tilde{Y} be residuals from a regression of Y on X_2 so that

$$\begin{aligned}\tilde{Y} &\equiv Y - X_2 \hat{g} \\ &= Y - X_2 (X_2' X_2)^{-1} X_2' Y \\ &= \left[I - X_2 (X_2' X_2)^{-1} X_2' \right] Y \\ &= M_2 Y.\end{aligned}$$

An important special case of this is that if I regress something on itself, the residuals are all zero

That is

$$\begin{aligned}M_2 X_2 &= X_2 - X_2(X_2' X_2)^{-1} X_2' X_2 \\ &= 0\end{aligned}$$

If I think of the GMM moment equations for least squares I get the “two equations”

$$0 = X_1' (Y - X_1 \hat{\beta}_1 - X_2 \hat{\beta}_2)$$

$$0 = X_2' (Y - X_1 \hat{\beta}_1 - X_2 \hat{\beta}_2)$$

We can solve for $\hat{\beta}_2$ in the second as

$$\hat{\beta}_2 = (X_2' X_2)^{-1} X_2' (Y - X_1 \hat{\beta}_1)$$

Now plug this into the first

$$\begin{aligned} 0 &= X_1' (Y - X_1 \hat{\beta}_1 - X_2 \hat{\beta}_2) \\ &= X_1' (Y - X_1 \hat{\beta}_1 - X_2 (X_2' X_2)^{-1} X_2' (Y - X_1 \hat{\beta}_1)) \\ &= X_1' M_2 Y - X_1' M_2 X_1 \hat{\beta}_1 \end{aligned}$$

Or

$$\begin{aligned}\hat{\beta}_1 &= (X_1' M_2 X_1)^{-1} X_1' M_2 Y \\ &= (\tilde{X}_1' \tilde{X}_1)^{-1} \tilde{X}_1' \tilde{Y}\end{aligned}$$

That is if I

- ① Run a regression of X_1 on X_2 and form its residuals \tilde{X}_1
- ② Run a regression of Y on X_2 and form its residuals \tilde{Y}
- ③ Run a regression of \tilde{Y} on \tilde{X}_1

Since I derived this from the sample analogue of the GMM moment equations (or the normal equations), this gives me exactly the same result as if I had run the full regression of Y on X_1 and X_2

It turns out the same idea works for IV.

Put everything we had before into large Matrix notation and we can write GMM as:

$$\begin{aligned}0 &= Z' (Y - T\hat{\alpha}_{IV} - X\hat{\beta}_{IV}) \\0 &= X' (Y - T\hat{\alpha}_{IV} - X\hat{\beta}_{IV})\end{aligned}$$

The second can be solved as

$$\hat{\beta}_{IV} = (X'X)^{-1} X' (Y - T\hat{\alpha}_{IV})$$

Now plug this into the first

$$\begin{aligned}0 &= Z' (Y - T\hat{\alpha}_{IV} - X\hat{\beta}_{IV}) \\&= Z' (Y - T\hat{\alpha}_{IV} - X(X'X)^{-1} X' (Y - T\hat{\alpha}_{IV})) \\&= Z'M_X Y - Z'M_X T\hat{\alpha}_{IV}\end{aligned}$$

so

$$\begin{aligned}\hat{\alpha}_{IV} &= (Z' M_X T)^{-1} Z' M_X Y \\ &= \frac{\tilde{Z}' \tilde{Y}}{\tilde{Z}' \tilde{T}} \\ &\approx \frac{\text{cov}(\tilde{Z}_i, \tilde{Y}_i)}{\text{cov}(\tilde{Z}_i, \tilde{T}_i)}\end{aligned}$$

(This last expression assumes that there is an intercept in the model. If not it would be expected values rather than covariances, but covariances make things easier to interpret-at least to me)

To see consistency from this perspective note that

$$\begin{aligned}\tilde{Y} &= M_X Y \\ &= \alpha M_X T + M_X X \beta + M_X u \\ &= \alpha \tilde{T} + \tilde{u}\end{aligned}$$

so

$$\begin{aligned}\hat{\alpha}_{IV} &\approx \frac{\text{cov}(\tilde{Z}_i, \tilde{Y}_i)}{\text{cov}(\tilde{Z}_i, \tilde{T}_i)} \\ &\approx \frac{\text{cov}(\tilde{Z}_i, \alpha \tilde{T}_i + \tilde{u}_i)}{\text{cov}(\tilde{Z}_i, \tilde{T}_i)} \\ &= \alpha + \frac{\text{cov}(\tilde{Z}_i, \tilde{u}_i)}{\text{cov}(\tilde{Z}_i, \tilde{T}_i)}\end{aligned}$$

This formula is helpful. In order for the model to be consistent you need

- $cov(\tilde{Z}_i, \tilde{u}_i) = 0$
- $cov(\tilde{Z}_i, \tilde{T}_i) \neq 0$

But more generally for the asymptotic bias to be small you want

- $cov(\tilde{Z}_i, \tilde{u}_i)$ to be small
- $|cov(\tilde{Z}_i, \tilde{T}_i)|$ to be large

This means that in practice there is some tradeoff between them.

If your instrument is not very powerful, a little bit of correlation in $cov(\tilde{Z}_i, \tilde{u}_i)$ could lead to a large asymptotic bias.

As a concrete example lets compare IV to OLS.

OLS is really just a special case of IV with $Z_i = T_i$

Then we get

$$\hat{\alpha}_{IV} \approx \alpha + \frac{\text{cov}(\tilde{Z}_i, \tilde{u}_i)}{\text{cov}(\tilde{Z}_i, \tilde{T}_i)}$$
$$\hat{\alpha}_{OLS} \approx \alpha + \frac{\text{cov}(\tilde{T}_i, \tilde{u}_i)}{\text{cov}(\tilde{T}_i, \tilde{T}_i)}$$

If $\text{cov}(\tilde{Z}_i, \tilde{u}_i) = 0$ and $\text{cov}(\tilde{T}_i, \tilde{u}_i) \neq 0$ then IV is consistent and OLS is not

However, $\text{cov}(\tilde{Z}_i, \tilde{u}_i) < \text{cov}(\tilde{T}_i, \tilde{u}_i)$ does not guarantee less bias because it also depends on $\text{cov}(\tilde{Z}_i, \tilde{T}_i) = 0$ and $\text{cov}(\tilde{T}_i, \tilde{T}_i) \neq 0$

Assumptions

Lets think about the assumptions

Lets ignore the first stage thing and assume that isn't a problem

we looked at two different things

- GMM

$$E(Z_i u_i) = 0, E(X_i u_i) = 0$$

- My derivation

$$\text{cov}(\tilde{Z}_i, \tilde{u}_i) = 0$$

Are these the same?

Well ...

- We know that GMM gives consistent estimates and if $cov(\tilde{Z}_i, \tilde{u}_i) \neq 0$ then $\hat{\alpha}$ is consistent. So the first assumptions better imply the second
- However, it doesn't have to be the other way around. We need $cov(\tilde{Z}_i, \tilde{u}_i) = 0$ to get consistent estimates of α , but we haven't required consistent estimates of β .

Lets explore

First lets show that $E(Z_i u_i) = 0$, $E(X_i u_i) = 0$ implies
 $cov(\tilde{Z}_i, \tilde{u}_i) = 0$

$$\begin{aligned}\tilde{Z}_i &= Z_i - X_i' \Gamma_{ZX} \\ \tilde{u}_i &= u_i - X_i' \Gamma_{uX}\end{aligned}$$

$E(X_i u_i) = 0$ implies that $\Gamma_{uX} \approx 0$ so $\tilde{u}_i \approx u_i$

but then

$$\begin{aligned}cov(\tilde{Z}_i, \tilde{u}_i) &= cov(Z_i - X_i' \Gamma_{ZX}, u_i) \\ &= cov(Z_i, u_i) - cov(X_i' \Gamma_{ZX}, u_i) \\ &= 0\end{aligned}$$

To see it doesn't go the other way, suppose that I can write

$$u_i = X_i' \delta + \varepsilon_i$$

with $E(X_i, \varepsilon_i) = 0$, $E(Z_i, \varepsilon_i) = 0$

Then $\tilde{u}_i \approx \varepsilon_i$ so

$$\begin{aligned} \text{cov}(\tilde{Z}_i, \tilde{u}_i) &= \text{cov}(Z_i - X_i' \Gamma_{zx}, \varepsilon_i) \\ &= \text{cov}(Z_i, \varepsilon_i) - \text{cov}(X_i' \Gamma_{zx}, \varepsilon_i) \\ &= 0 \end{aligned}$$

We can also see that we can write

$$Y_i = \alpha T_i + X_i' \beta + u_i$$
$$\alpha T_i + X_i' (\beta + \delta) + \varepsilon_i$$

So IV gives a consistent estimate of α but not β

Simultaneous equations

The second and third way to see IV comes from the simultaneous equations framework

$$Y_i = \alpha T_i + X_i' \beta + u_i$$
$$T_i = \rho Y_i + X_i' \gamma + Z_i \delta + \nu_i$$

These are called the “structural equations”

Note the difference between X_i and Z_i in that we restrict what can affect what.

We could also have stuff that affects Y_i but not T_i but let's not worry about that (we are still allowing this as a possibility as some of the γ coefficients could be zero)

The model with $\rho = 0$ simplifies things, but let's focus on what happens when it isn't

We assume that

$$E(u_i | X_i, Z_i) = 0$$

$$E(v_i | X_i, Z_i) = 0$$

but notice that if $\rho \neq 0$, then almost for sure T_i is correlated with u_i because u_i influences T_i through Y_i

It is useful to calculate the “reduced form” for T_i , namely

$$\begin{aligned}T_i &= \rho Y_i + X_i' \gamma + Z_i \delta + \nu_i \\&= \rho [\alpha T_i + X_i' \beta + u_i] + X_i' \gamma + Z_i \delta + \nu_i \\&= \rho \alpha T_i + X_i' [\rho \beta + \gamma] + Z_i' \delta + (\rho u_i + \nu_i) \\&= X_i' \frac{\rho \beta + \gamma}{1 - \rho \alpha} + Z_i' \frac{\delta}{1 - \rho \alpha} + \frac{\rho u_i + \nu_i}{1 - \rho \alpha} \\&= X_i' \beta_2^* + Z_i' \delta_2^* + \nu_i^*\end{aligned}$$

where

$$\begin{aligned}\beta_2^* &\equiv \frac{\rho \beta + \gamma}{1 - \rho \alpha} \\ \delta_2^* &\equiv \frac{\delta}{1 - \rho \alpha} \\ \nu_i^* &\equiv \frac{\rho u_i + \nu_i}{1 - \rho \alpha}\end{aligned}$$

Note that $E(\nu_i^* | X_i, Z_i) = 0$, so one can obtain a consistent estimate of β_2^* and δ_2^* by regressing T_i on X_i and Z_i .

This is called the “reduced form” equation for T_i

Note that the parameters here are not the fundamental structural parameters themselves, but they are a known function of these parameters

To me this is the classic definition of reduced form (you need to have a structural model)

We can obtain a consistent estimate of α as long as we have an exclusion restriction

That is we need some Z_i that affects T_i but not Y_i directly

I want to show this in two different ways

Method 1

We can also solve for the reduced form for Y_i

$$\begin{aligned} Y_i &= \alpha T_i + X_i' \beta + u_i \\ &= X_i \frac{\alpha\gamma + \beta}{1 - \alpha\rho} + Z_i \frac{\alpha\delta}{1 - \alpha\rho} + \frac{\alpha\nu_i + u_i}{1 - \alpha\rho} \\ &= X_i \beta_1^* + Z_i \delta_1^* + u_i^* \end{aligned}$$

with

$$\begin{aligned} \beta_1^* &\equiv \frac{\alpha\gamma + \beta}{1 - \alpha\rho} \\ \delta_1^* &\equiv \frac{\alpha\delta}{1 - \alpha\rho} \\ u_i^* &\equiv \frac{\alpha\nu_i + u_i}{1 - \alpha\rho} \end{aligned}$$

Like the other reduced form, we can get a consistent estimate of β_1^* and δ_1^* by regressing Y_i on X_i and Z_i .

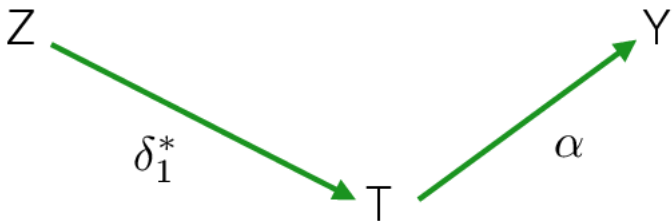
Notice then that

$$\frac{\delta_1^*}{\delta_2^*} = \alpha$$

So we can get a consistent estimate of α simply by taking the ratio of the reduced form coefficients

It also gives another interpretation of IV:

- δ_2^* is the causal effect of Z_i on T_i
- δ_1^* is the causal effect of Z_i on Y_i -it only operates through T_i

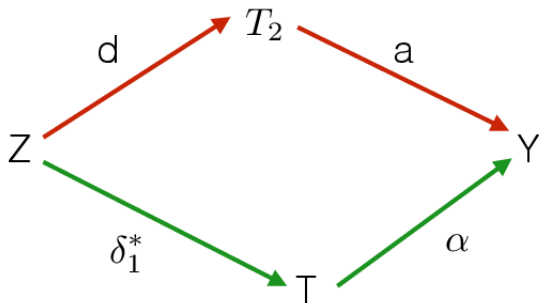


If we increase Z_i by one unit this leads T_i to increase by δ_2^* units which causes Y_i to increase by $\delta_2^*\alpha$ units

Thus the causal effect of Z_i on Y_i is $\delta_2^*\alpha$ units

This illustrates another important way to think about an instrument: the key assumption is that T_i is the only channel through which Z_i influences Y_i

Suppose there was another



Then the causal effect of Z_i on Y_i would be $\delta_2^* \alpha + da$ and IV would be

$$\frac{\delta_2^* \alpha + da}{\delta_2^*}$$

In the exactly identified case (i.e. one Z_i), this is numerically identical to IV.

To see why, note that in a regression of Y_i and T_i on (X_i, Z_i) yields

$$\hat{\delta}_1 = \frac{\tilde{Z}'_i \tilde{Y}_i}{\tilde{Z}'_i \tilde{Z}_i}$$

$$\hat{\delta}_2 = \frac{\tilde{Z}'_i \tilde{T}_i}{\tilde{Z}'_i \tilde{Z}_i}$$

so

$$\begin{aligned} \frac{\hat{\delta}_1}{\hat{\delta}_2} &= \frac{\tilde{Z}'_i \tilde{Y}_i}{\tilde{Z}'_i \tilde{T}_i} \\ &= \hat{\alpha}_{IV} \end{aligned}$$

This is just math-it does not require that the "Structural equation" determining T_i be correct

Method 2

Define

$$T_i^f \equiv X_i' \beta_2^* + Z_i' \delta_2^*$$

and suppose that T_i^f were known to the econometrician

Now notice that

$$\begin{aligned} Y_i &= \alpha T_i + X_i' \beta + u_i \\ &= \alpha [T_i^f + \nu_i^*] + X_i' \beta + u_i \\ &= \alpha T_i^f + X_i' \beta_2 + (\alpha \nu_i^* + u_i) \end{aligned}$$

One could get a consistent estimate of α by regressing Y_i on X_i and T_i^f .

Two Stage Least Squares

In practice we don't know T_i^f but can get a consistent estimate of it from the fitted values of a reduced form regression call this \hat{T}_i (it is crucial that the reduced form gives us consistent estimates of β_2^* and δ_2^*)

That is:

- 1 Regress T_i on X_i and Z_i , form the predicted value \hat{T}_i
- 2 Regress Y on X_i and \hat{T}_i

To run the second regression one needs to be able to vary \hat{T}_i separately from X_i which can only be done if there is a Z_i

It turns out that 2SLS is also numerically identical to IV (with 1 instrument)

Note that

$$\hat{T} = Z^* (Z^{*'} Z^*)^{-1} Z^{*'} T$$

so

$$\hat{B}_{2SLS} = \left(\left[\begin{array}{cc} \hat{T} & X \end{array} \right]' \left[\begin{array}{cc} \hat{T} & X \end{array} \right] \right)^{-1} \left[\begin{array}{cc} \hat{T} & X \end{array} \right]' Y$$

However, note that we can write

$$X = Z^* (Z^{*'} Z^*)^{-1} Z^{*'} X$$

That is projecting X on (X, Z) and using it to predict X will be a perfect fit.

That means that(using notation from earlier) that

$$\begin{bmatrix} \hat{T} & X \end{bmatrix} = Z^* (Z^{*'} Z^*)^{-1} Z^{*'} X^*$$

Then (in the exactly identified case) we can write

$$\begin{aligned} \hat{B}_{2SLS} &= \left(X^{*'} Z^* (Z^{*'} Z^*)^{-1} Z^{*'} Z^* (Z^{*'} Z^*)^{-1} Z^{*'} X^* \right)^{-1} \times \\ &\quad X^{*'} Z^* (Z^{*'} Z^*)^{-1} Z^{*'} Y \\ &= \left(X^{*'} Z^* (Z^{*'} Z^*)^{-1} Z^{*'} X^* \right)^{-1} X^{*'} Z^* (Z^{*'} Z^*)^{-1} Z^{*'} Y \\ &= (Z^{*'} X^*)^{-1} (Z^{*'} Z^*) (X^{*'} Z^*)^{-1} X^{*'} Z^* (Z^{*'} Z^*)^{-1} Z^{*'} Y \\ &= (Z^{*'} X^*)^{-1} (Z^{*'} Z^*) (Z^{*'} Z^*)^{-1} Z^{*'} Y \\ &= (Z^{*'} X^*)^{-1} Z^{*'} Y \\ &= \hat{\beta}_{IV} \end{aligned}$$

3 Interpretations

Thus with 1 instrument we have 3 equivalent ways to derive IV:

- ① GMM estimator or $(Z'T)^{-1}Z'Y$
- ② Ratio of reduced form estimates-rescaling the reduced form
- ③ 2SLS-direct effect of fitted model

With more than one instrument only one of these procedures works-we'll worry about that later

Examples

There are three main reasons people use IV

- ① Simultaneity bias: $\rho \neq 0$
- ② Omitted Variable bias : There are unobservables that contribute to u_i that are correlated with T_i
- ③ Measurement Error: We do not observe T_i perfectly

While the first is the original reason for IV, in practice omitted variable bias is typically the biggest concern

A classic (perhaps the classic) example is the returns to schooling.

Returns to Schooling

This comes from the Card's chapter in the 1999 Handbook of Labor Economics

Lets assume that

$$\log(W_i) = \alpha S_i + X_i' \beta + \varepsilon_i$$

where W_i is wages, S_i is schooling, and X_i is other stuff

The biggest problem is unobserved ability

We are worried about ability bias we want to use instrumental variables

A good instrument should have two qualities:

- It should be correlated with schooling
- It should be uncorrelated with unobserved ability (and other unobservables)

Many different things have been tried. Lets go through some of them

Family Background

If my parents earn quite a bit of money it should be easier for me to borrow for college

Also they might put more value on education

This should make me more likely to go

This has no direct effect on my income-Wisconsin did not ask how much education my Father had when they made my offer

But is family background likely to be uncorrelated with unobserved ability?

Closeness of College

If I have a college in my town it should be much easier to attend college

- I can live at home
- If I live on campus
 - I can travel to college easily
 - I can come home for meals and to get my clothes washed
- I can hang out with my friends from High school

But is this uncorrelated with unobserved ability?

Quarter of Birth

This is the most creative

Consider the following two aspects of the U.S. education system (this actually varies from state to state and across time but ignore that for now),

- People begin Kindergarten in the calendar year in which they turn 5
- You must stay in school until you are 16

Now consider kids who:

- Can't stand school and will leave as soon as possible
- Obey truancy law and school age starting law
- Are born on either December 31, 1972 or January 1, 1973

Those born on December 31 will

- turn 5 in the calendar year 1977 and will start school then (at age 4)
- will stop school on their 16th birthday which will be on Dec. 31, 1988
- thus they will stop school during the winter break of 11th grade

Those born on January 1 will

- turn 5 in the calendar year 1978 and will start school then (at age 5)
- will stop school on their 16th birthday which will be on Jan. 1, 1989
- thus they will stop school during the winter break of 10th grade

The instrument is a dummy variable for whether you are born on Dec. 31 or Jan 1

This is pretty cool:

- For reasons above it will be correlated with education
- No reason at all to believe that it is correlated with unobserved ability

The Fact that not everyone obeys perfectly is not problematic:

An instrument just needs to be correlated with schooling, it does not have to be perfectly correlated

In practice we can't just use the day as an instrument, use "quarter of birth" instead

Policy Changes

Another possibility is to use institutional features that affect schooling

Here often institutional features affect one group or one cohort rather than others

TABLE II
OLS AND IV ESTIMATES OF THE RETURN TO EDUCATION WITH INSTRUMENTS BASED ON FEATURES OF THE SCHOOL SYSTEM

Author	Sample and Instrument		Schooling Coefficients	
			OLS	IV
1. Angrist and Krueger (1991)	1970 and 1980 Census Data, Men. Instruments are quarter of birth interacted with year of birth. Controls include quadratic in age and indicators for race, marital status, urban residence.	1920–29 cohort in 1970	0.070 (0.000)	0.101 (0.033)
		1930–39 cohort in 1980	0.063 (0.000)	0.060 (0.030)
		1940–49 cohort in 1980	0.052 (0.000)	0.078 (0.030)
2. Staiger and Stock (1997)	1980 Census, Men. Instruments are quarter of birth interacted with state and year of birth. Controls are same as in Angrist and Krueger, plus indicators for state of birth. LIML estimates.	1930–39 cohort in 1980	0.063 (0.000)	0.098 (0.015)
		1940–49 cohort in 1980	0.052 (0.000)	0.088 (0.018)
3. Kane and Rouse (1993)	NLS Class of 1972, Women. Instruments are tuition at 2 and 4-year state colleges and distance to nearest college. Controls include race, part-time status, experience. Note: Schooling measured in units of college credit equivalents.	Models without test score or parental education	0.080 (0.005)	0.091 (0.033)
		Models with test scores and parental education	0.063 (0.005)	0.094 (0.042)
4. Card (1995b)	NLS Young Men (1966 Cohort) Instrument is an indicator for a nearby 4-year college in 1966, or the interaction of this with parental education. Controls include race, experience (treated as endogenous), region, and parental education	Models that use college proximity as instrument (1976 earnings)	0.073 (0.004)	0.132 (0.049)
		Models that use college proximity \times family background as instrument	—	0.097 (0.048)

5. Conneely and Uusitalo (1997)	Finnish men who served in the army in 1982, and were working full time in civilian jobs in 1994. Administrative earnings and education data. Instrument is living in university town in 1980. Controls include quadratic in experience and parental education and earnings.	Models that exclude parental education and earnings	0.085 (0.001)	0.110 (0.024)
		Models that include parental education and earnings	0.083 (0.001)	0.098 (0.035)
6. Harmon and Walker (1995)	British Family Expenditure Survey 1978–86 (men). Instruments are indicators for changes in the minimum school leaving age in 1947 and 1973. Controls include quadratic in age, survey year, and region.		0.061 (0.001)	0.153 (0.015)
7. Ichino and Winter-Ebmer (1998)	Austria: 1983 Census, men born before 1946. Germany: 1986 GSOEP for adult men. Instrument is indicator for 1930–35 cohort. (Second German IV also uses dummy for father's veteran status). Controls include age, unemployment rate at age 14, and father's education (Germany only). Education measure is dummy for high school or more.	Austrian Men	0.518 (0.015)	0.947 (0.343)
		German Men	0.289 (0.031)	0.590/0.708 (0.844) (0.279)
8. Lemieux and Card (1998)	Canadian Census, 1971 and 1981: French-speaking men in Quebec and English-speaking in Ontario. Instrument is dummy for Ontario men age 19–22 in 1946. Controls include full set of experience dummies and Quebec-specific cubic experience profile.	1971 Census:	0.070 (0.002)	0.164 (0.053)
		1981 Census:	0.062 (0.001)	0.076 (0.022)
9. Meghir and Palme (1999)	Swedish Level of Living Survey (SLLS) data for men born 1945–55, with earnings in 1991, and Individual Statistics (IS) sample of men born in 1948 and 1953, with earnings in 1993. Instrument is dummy for attending “reformed” school system at age 13. Other controls include cohort, father's education, and county dummies. Models for IS data also include test scores at age 13.	SLLS Data (Years of education)	0.028 (0.007)	0.036 (0.021)
		IS Data (Dummy for 1–2 years of college relative to minimum schooling)	0.222 (0.020)	0.245 (0.082)

TABLE II—Continued

Author	Sample and Instrument		Schooling Coefficients	
			OLS	IV
10. Maluccio (1997)	Bicol Multipurpose Survey (rural Philippines): male and female wage earners age 20–44 in 1994, whose families were interviewed in 1978. Instruments are distance to nearest high school and indicator for local private high school. Controls include quadratic in age and indicators for gender and residence in a rural community.	Models that do not control for selection of employment status or location	0.073 (0.011)	0.145 (0.041)
		Models with selection correction for location and employment status	0.063 (0.006)	0.113 (0.033)
11. Duflo (1999)	1995 Intercensal Survey of Indonesia: men born 1950–72. Instruments are interactions of birth year and targeted level of school building activity in region of birth. Other controls are dummies for year and region of birth and interactions of year of birth and child population in region of birth. Second IV adds controls for year of birth interacted with regional enrollment rate and presence of water and sanitation programs in region.	Model for hourly wage	0.078 (0.001)	0.064/0.091 (0.025) (0.023)
		Model for monthly wage with imputation for self-employed.	0.057 (0.003)	0.064/0.049 (0.017) (0.013)

Notes: See text for sources and more information on individual studies.

Consistently IV estimates are higher than OLS

Why?

- Bad Instruments
- Ability Bias
- Measurement Error
- Publication Bias
- Discount Rate Bias

Measurement Error

Another way people use instruments is for measurement error

In the classic model lets get rid of X 's so we want to measure the effect of T on Y .

$$Y_i = \beta_0 + \alpha T_i + u_i$$

and lets not worry about other issues so assume that $cov(T_i, u_i) = 0$.

The complication is that I don't get to observe T_i , I only get to observe a noisy version of it:

$$\tau_{1i} = T_i + \xi_i$$

where ξ_i is i.i.d measurement error with variance σ_ξ^2

What happens if I run the regression on τ_{1i} instead of T_i ?

$$\begin{aligned}\hat{\alpha} &\approx \frac{\text{Cov}(\tau_{1i}, Y_i)}{\text{Var}(\tau_{1i})} \\ &= \frac{\text{Cov}(T_i + \xi_i, \beta_0 + \alpha T_i + u_i)}{\text{Var}(T_i + \xi_i)} \\ &= \alpha \frac{\text{Var}(T_i)}{\text{Var}(T_i) + \sigma_\xi^2}\end{aligned}$$

Why is OLS biased?

Lets rewrite the model as

$$\begin{aligned} Y_i &= \beta_0 + \alpha T_i + u_i \\ &= \beta_0 + \alpha T_i + \alpha \xi_i + u_i - \alpha \xi_i \\ &= \beta_0 + \alpha \tau_{1i} + (u_i - \alpha \xi_i). \end{aligned}$$

You can see the problem with OLS: $\tau_{1i} = T_i + \xi_i$ is correlated with $(u_i - \alpha \xi_i)$

Now suppose we have another measure of T_i ,

$$\tau_{2i} = T_i + \eta_i$$

where η_i is uncorrelated with everything else in the model.

Using this as an instrument gives us a solution.

τ_{2i} is correlated with τ_{1i} (through T_i), but uncorrelated with $(u_i - \alpha\xi_i)$ so we can use τ_{2i} as an instrument for τ_{1i} .

Twins

(Here we will think about both measurement error and fixed effect approaches)

$$\log(w_{if}) = \theta_f + \alpha S_{if} + u_{if}$$

The problem is that θ_f is correlated with S_{if}

We can solve by differencing

$$\Delta \log(w_f) = \alpha \Delta S_f + \Delta u_f$$

if ΔS_f is uncorrelated with Δu_f , then we can use this to get consistent estimates of α

The problem here is that a little measurement error can screw up things quite a bit because the variance of ΔS_f is small.

A solution of this is to get two measures on schooling

- Ask me about my schooling
- Also ask my brother about my schooling
- do the same think for my brother's schooling

This gives us two different measure of ΔS_{if} .

Use one as an instrument for the other

Table 6

Cross-sectional and within-family differenced estimates of the return to education for twins^a

Author	Sample and specification		Cross-sectional OLS	Differenced	
				OLS	IV
1. Ashenfelter and Rouse (1998)	1991–1993 Princeton Twins Survey. Identical male and female twins. Controls include quadratic in age, gender and race. Added controls include tenure, marital status and union status.	Basic	0.110 (0.010)	0.070 (0.019)	0.088 (0.025)
		Basic + added controls	0.113 (0.010)	0.078 (0.018)	0.100 (0.023)
2. Rouse (1997)	1991–1995 Princeton Twins Survey. Identical male and female twins. Basic controls as above.		0.105 (0.008)	0.075 (0.017)	0.110 (0.023)
3. Miller et al. (1995)	Australian Twins Register. Identical and fraternal twins. Controls include quadratic in age, gender, marital status. Incomes imputed from occupation	Identical twins	0.064 (0.002)	0.025 (0.005)	0.048 (0.010)
		Fraternal twins	0.066 (0.002)	0.045 (0.005)	0.074 (0.008)
4. Behrman et al. (1994)	NAS-NRC white male twins born 1917–1927, plus male twins born 1936–1955 from Minnesota Twins Registry. Controls include quadratic in age ^b	Identical twins	0.071 (0.002)	0.035 (0.005)	0.056 –
		Fraternal twins	0.073 (0.003)	0.057 (0.005)	0.071 –
5. Isacson (1997)	Swedish same-sex twins with both administrative and survey measures of schooling. Controls include sex, marital status, quadratic in age, and residence in a large city ^c	Identical twins	0.049 (0.002)	0.023 (0.004)	0.024 (0.008)
		Fraternal twins	0.051 (0.002)	0.040 (0.003)	0.054 (0.006)

Overidentification

What happens when we have more than one instrument?

Lets think about a general case in which Z_i is multidimensional

- Let K_Z be the dimension of Z_i^*
- Let K_X denote the dimension of X_i^*

Now we have more equations then parameters so we can no longer solve for \hat{B} using

$$0 = Z^{*'} (Y - X^* \hat{B})$$

because this gives us K_Z equations in K_X unknowns.

A simple solution is follow GMM and weight the moments by some $K_Z \times K_Z$ weighting matrix Ω and then minimize

$$\left[Z^{*'} (Y - X^* B) \right]' \Omega \left[Z^{*'} (Y - X^* B) \right]$$

which gives

$$-2X^{*'} Z^* \Omega Z^{*'} (Y - X^* \hat{B}) = 0$$

(notice that in the exactly identified case $X^{*'} Z^* \Omega$ drops out)

We can solve directly for our estimator

$$\hat{B}_{GMM} = \left(X^{*'} Z^* \Omega Z^{*'} X^* \right)^{-1} X^{*'} Z^* \Omega Z^{*'} Y$$

Two staged least squares is a special case of this:

$$\hat{B}_{2SLS} = \left(X^{*'} Z^* (Z^{*'} Z^*)^{-1} Z^{*'} X^* \right)^{-1} X^{*'} Z^* (Z^{*'} Z^*)^{-1} Z^{*'} Y$$

Notice that this is the same as \hat{B}_{GMM} when

$$\Omega = (Z^{*'} Z^*)^{-1}$$

Consistency

$$\begin{aligned}\hat{B}_{GMM} &= \left(\frac{1}{N} X^{*'} Z^* \Omega \frac{1}{N} Z^{*'} X^* \right)^{-1} \frac{1}{N} X^{*'} Z^* \Omega \frac{1}{N} Z^{*'} (X^* B + U) \\ &= B + \left(\frac{1}{N} X^{*'} Z^* \Omega \frac{1}{N} Z^{*'} X^* \right)^{-1} \frac{1}{N} X^{*'} Z^* \Omega \frac{1}{N} Z^{*'} U \\ &\approx B + \left(E \left(X_i^* Z_i^{*'} \right) \Omega E \left(Z_i^* X_i^{*'} \right) \right)^{-1} E \left(X_i^* Z_i^{*'} \right) \Omega E \left(Z_i^* u_i \right) \\ &= B\end{aligned}$$

Inference

$$\sqrt{N}(\hat{B} - B) = \left(\frac{1}{N} X^{*'} Z^* \Omega \frac{1}{N} Z^{*'} X^* \right)^{-1} \frac{1}{N} X^{*'} Z^* \Omega \frac{1}{\sqrt{N}} Z^{*'} U$$

Using a standard central limit theorem with i.i.d. data

$$\begin{aligned} \frac{1}{\sqrt{N}} Z^{*'} U &= \frac{1}{\sqrt{N}} \sum_{i=1}^N Z_i^* u_i \\ &\approx N\left(0, E\left(u_i^2 Z_i^* Z_i^{*'}\right)\right) \end{aligned}$$

Thus

$$\sqrt{N}(\hat{B} - B) \approx N(0, A'VA)$$

with

$$V = E\left(X_i^* Z_i^{*'}\right) \Omega E\left(u_i^2 Z_i^* Z_i^{*'}\right) \Omega E\left(Z_i^* X_i^{*'}\right)$$

$$A = \left(E\left(X_i^* Z_i^{*'}\right) \Omega E\left(Z_i^* X_i^{*'}\right) \right)^{-1}$$

From GMM results we know that the efficient weighting matrix is

$$\Omega = E \left(u_i^2 Z_i^* Z_i^{*'} \right)^{-1}$$

in which case the Covariance matrix simplifies to

$$\left(E \left(X_i^* Z_i^{*'} \right) E \left(u_i^2 Z_i^* Z_i^{*'} \right)^{-1} E \left(Z_i^* X_i^{*'} \right) \right)^{-1}$$

This also means that under homoskedasticity two staged least squares is efficient.

Overidentification Tests

Lets think about testing in the following way.

Suppose we have two instruments so that we have three sets of moment conditions

$$0 = Z_1' (Y - T\hat{\alpha} - X\hat{\beta})$$

$$0 = Z_2' (Y - T\hat{\alpha} - X\hat{\beta})$$

$$0 = X' (Y - T\hat{\alpha} - X\hat{\beta})$$

As before we can use partitioned regression to deal with the X's and then write the first two moment equations as

$$0 = \tilde{Z}'_1 (\tilde{Y} - \tilde{T}\hat{\alpha})$$
$$0 = \tilde{Z}'_2 (\tilde{Y} - \tilde{T}\hat{\alpha})$$

The way I see the overidentification test is whether we can find an $\hat{\alpha}$ that solves both equations.

That is let

$$\hat{\alpha}_1 = \frac{\tilde{Z}'_1 \tilde{Y}}{\tilde{Z}'_2 \tilde{T}}$$

$$\hat{\alpha}_2 = \frac{\tilde{Z}'_1 \tilde{Y}}{\tilde{Z}'_2 \tilde{T}}$$

If

$$\hat{\alpha}_1 \approx \hat{\alpha}_2$$

then the test will not reject the model, otherwise it will

For this reason I am not a big fan of overidentification tests:

- If you have two crappy instruments with roughly the same bias you will fail to reject
- Why not just estimate $\hat{\alpha}_1$ and $\hat{\alpha}_2$ and look at them? It seems to me that you learn much more from that than a simple F-statistic