# Micro Econometrics

### Sol Yates

### February 21, 2024

### Contents

1	IV		1
	1.1	$Case: length(z) = length(x) \dots \dots$	1
	1.2	2SLS/GIVE	2
	1.3	Properties	3
	1.4	Group Mean Estimator	3

### Lecture 4: IV p2

Tue 20 Feb 16:07

## 1 IV

Model  $y = x\beta + u$  with x a vector of k exogenous and endogenous regressors and z a vector of m IV's (including the exogenous variable)

- 1. m = k: the model is just identified, we have an instrument for each endogenous variable  $\Rightarrow$  use IV
- 2. m < k: the model is not identified, we do not have enough IVs
- 3. m > k: the model is over-identified  $\rightarrow$  we have too many IVs. Use GIVE / 2SLS

## 1.1 Case: length(z) = length(x)

Model:  $y = x\beta + u$ ,  $x = (q, x_2, ..., x_k)$  and  $z = (1, x_2, ..., x_{k-1}, z_1)$ . We know  $Cov(x_j, u) = 0$  for j = 2, ..., k-1 and  $Cov(x_k, u \neq 0)$ 

We have an instrument for  $x_k$ :

- Exogenous  $Cov(z_1, u) = 0$
- Partial Correlation :  $\theta_1 \neq 0$  in  $x_k = \delta_1 + \delta_2 x_2 + \ldots + \delta_{k-1} x_{k-1} + \theta_1 z_1 + r_k$

Where the moment conditions imply:

$$E[z'u] = E[z'(y - x\beta)] = 0$$

We have one instrument at our disposal for this endogenous regressor, we include the constant and all the exogenous regressors because they can be used for instruments for themselves.

Partial correlation best seen by regressing endogenous regressor  $x_k$  on all exogenous variables plus the instrument for  $x_k$  and we need the parameter on the instrument  $\theta_1$  not to be 0, thus partial correlation,

the correlation cannot be 0 after the other effects have been 'netted' out. Different to simple case where sufficient to have covariance between endogenous regressor and instrument  $\neq 0$ 

Exogeneity leads to above expression, plugging in expression for u.

Multiplying the model through with z', taking expectation and using the moment condition:

$$E[z'y] = E[z'x]\beta$$
 if rank  $E[z'x] = k$  
$$\beta = [E[z'x]]^{-1}E[z'y]$$

There is a unique solution only under full rank and it can be shown that if we rule out perfect collinearity in z, full rank holds iff  $\theta_1 \neq 0$ 

Given a random sample, we can estimate consistently:

$$\hat{\beta}^{IV} = \left(\frac{1}{n} \sum_{i=1}^{n} z_i' x_i\right)^{-1} \left(\frac{1}{n} \sum_{i=1}^{n} z_i' y_i\right) = (Z'X)^{-1} Z'Y$$

Where Z and X are  $n \times k$  data matrices and y is  $N \times 1$ Given the assumptions, this estimator is consistent

#### 1.2 2SLS/GIVE

Case: length(z) > length(x):

the idea is to used the fitted values from the first stage regression of the endogenous regressor on all the exogenous variables (including the instruments) and use them as "instruments" in the IV estimator

$$z = (1, x_1, \dots, x_{k-1}, z_1, \dots, z_l) - m = k + I$$
 vector for  $x_k$ 

1. Fitted values from the first stage  $\hat{x}_i = (1, x_1, \dots, x_{k-1}, \hat{x}_k)$ 

$$\hat{x}_{ik} = \hat{\delta}_0 + \hat{\delta}_1 x_{i1} + \dots + \hat{\delta}_{k-1} x_{i,k-1} + \hat{\theta}_1 z_{i1} + \dots + \hat{\theta}_I z_{il}$$

$$\hat{x}_i = z_i \left( \sum z_i' z_i \right)^{-1} z_i' x_i$$

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

For this endogenous regressor you have several potential instruments at your disposal, you would then regress on exogenous variables from initial model and instruments. That gives you a vector of instruments that is equal to (including all potential instruments).

Then you start with obtaining fitted values from First Stage (FS) regressing  $x_k$  on exogenous regressors  $\delta$  and instruments z

Using the fitted values as instruments :

$$\hat{\beta}^{IV} = \left(\frac{1}{n} \sum_{i=1}^{n} \hat{x}_i' x_i\right)^{-1} \left(\frac{1}{n} \sum_{i=1}^{n} \hat{x}_i y_i\right) = (\hat{X}' X)^{-1} \hat{X}' Y$$

Then, using calculus, we can show that  $\hat{X}'X = \hat{X}'\hat{X}$  and hence

$$\hat{\beta}^{IV} = (\hat{X}'\hat{X})^{-1}\hat{X}'Y$$

Which is the GIVE / 2SLS estimator since it equals the OLS estimator on the fitted values from the first stage

Where we have simply replaced with fitted values, and then replaced in matrix form. We can essentially show this is the OLS estimator on the FS using the fitted values. Then we use this to plug in the IV estimator to obtain the OLS estimator on the fitted values

#### To obtain the $\beta$

- 1. First Stage: Obtain the fitted values  $\hat{x}_k$  from the regression  $x_k$  on  $1, x_1, \ldots, x_{k-1}, z_1, \ldots, z_l$
- 2. Second Stage: Run the OLS regression: y on  $1 + x_1, \dots, x_{k-1} + \hat{x_k}$
- However, omitting the exogenous regressors in the first stage is easily done and will lead to inconsistency
- And, SE obtained from the second step are incorrect

Testing for rank condition:  $H_0: \theta_1 = \ldots = \theta_l = 0$  vs at least one  $\theta_s$  for  $s = 1, \ldots, l$  is non zero.

### 1.3 Properties

Here x is 1 x k and generally includes unity, several elements of x may be endogenous, while z includes any exogenous variable

#### Assumption 1:

#### 1.4 Group Mean Estimator

In some situations, have instruments that can be changed into 2 groups, water (of birth/ financial year). 'Chop instrument into groups' like Moulton problem/structure.

It can be shown that group mean estimator is IV, a weighted least squares regression, where it is sufficient to know size of groups and means, do regression and obtain estimator that is equivalent to an IV estimator, that is consistent despite the fact we have an endogenous variable.

- where  $x_{ig}$  is endogenous
- we have g moment conditions, if this is IV, we know exogeneity must hold, whether group 1 or 2, the error term conditional on this group needs to be equal to 0, this must hold for all groups. Essentially, we have g different groups this is really  $E[y_{iq}|z_q=I]=\beta_0+\beta_1 E[x_{iq}|i:z_q=I]$
- To estimate an expectation, we replace with an average, since this is conditional, to estimate the expectation for the first group (born in the 1-st quarter), we take the average for the first group (conditional average by restricting to the first group) taking the means of all the groups
- ullet we do the same for  $\overline{x}$  and intuitively obtain

$$\overline{y_g} = \beta_0 + \beta_1 \overline{x_g} + \overline{u}_g$$

Doing this is the same as using dummy variables for quarter of birth in 2SLS regression, **thus** group means are consistent.

**Exercise 1.** Group mean estimator - what happens if the number of groups = 2 If we have dummy variable, we obtain the Wald estimator (last week). Our instrument, we obtain the same expression In order to derive,

- 1. regress x on dummies, x can only belong to 1, so fitted values are sample means of dependent variable  $x_{ig}$ . fitted values  $\hat{x_{ig}}$  are means  $\overline{x}_g$
- 2. Apply OLS on this after we have found fitted values, the predicted values are our sample means

#### Angrist and Krueger

- Does compulsory school attendance affect schooling and earnings
- using quarter of birth as an instrument
  - 1. Exclusion: Season of birth is a natural experiment and hence unrelated to innate ability, motivation or family connections
  - 2. Relevance: In the US, children were allowed to drop out at 16. Since the age of starting school differs, children have different lengths of schooling when they turn 16
- Potentially weak instrument and potential reasons why quarter of birth might be somewhat correlated with the error

	(1) Born in 1st Quarter of Year	(2) Born in 4th Quarter of Year	(3) Difference (Std. Error (1) – (2)
ln (weekly wage)	5.892	5.905	0135 (.0034)
Years of education	12.688	12.839	151 (.016)
Wald estimate of return to education			.089 (.021)
OLS estimate of return to education			.070 (.0005)

Figure 1: Wald estimates of IV - weak/exclusion violated?

can't test by how much IV exclusion is violated, it might be best to use OLS, but in the same sense it may be incorrect - can we search for better instrument? Or,

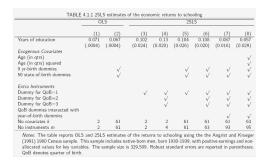


Figure 2: 2SLS estimates of economic returns to schooling

Inflated standard errors :  $0.\frac{0.021}{0.0005} = 42$ , even though the estimate is significant due to a large sample size, the 9% CI is large.

Problem also of a small  $R_{x,z}^2$ : the *instruments are weak*, it might be better to use OLS instead of IV including more instruments and covariates

- Reduces SE but comes to the cost of potentially having a weak instrument
- col 3: just identified 1 instrument
- col 4 : over identified (3 QoB instruments)
- col 5/6: + 59 covariates to 3/4: m k = 0 (or 2 resp)
- col 7 : +30 Ifs (m-k =32)
- col 8 : + age and  $age^2$  to x and z ( m k = 32)

But there is potential to test for over identifying restrictions, using the Sargon test

Why  $\beta$  larger? Asymptotic variance depends on error variance, depends on the  $R_{x,z}^2$  from first stage regression (x on instruments, will be higher if instruments highly relevant and vice versa).

### **Exercise 2.** Consequences of weak instruments

- 1. High SE
- 2. slight violation of exclusion restrictions leads to large bias

**Exercise 3.** Discuss intuition behind Hausman Test Exogenous regressor then both (testing for whether x endogenous) OLS and IV consistent (if we have valid instrument) If we have an endogenous regressor and OLS is not consistent, difference does not converge to 0 any more, test start follows chi-squared distribution with k degrees of freedom, to test for endogeneity 2 main assumptions, 1 test for, 1 assume

• relies on valid instrument, otherwise  $\hat{\beta}$  would not converge at all, this is almost critical assumption in Hausman test

**Tutorial 1.** IV and simultaneity bias we have expression for  $y_1$  and  $y_2$ , we are going to replace this equation, since our asymptotic bias will sum to .? asymptotic bias  $=\frac{Cov(u_1,y_2)}{V(y_2)}$  Then we plug long

covariance into  $Cov(u_1, y_2)$  We know  $cov(u_1, u_2) = 0$  and  $cov(u_1, z_2) = 0$ , but the problem is the variance is typically positive but depending on assumptions we can determine direction of bias based upon  $\alpha_2$  We show that in the formula we replace by  $y_1$ 

**Tutorial 2.** IV is asymptotically unbiased that is *plim*  $\tilde{\alpha}_1 = \alpha_1$  (the IV estimator)

**Tutorial 3.** why can  $z_2$  not be used as an instrument for  $y_1$  to estimate  $\alpha_2$  the slope of the supply curve It is under-identified, we don't have an endogenous variable at our disposal, we don't have a shock for  $y_1 \to \text{we}$  don't have an instrument, two endogenous variables require 2 instruments

**Tutorial 4.** Exercise 1 (tut3) Regression of log wage on education, estimate using OLS, internet, do you expect OLS to be trustworthy? Education on wage includes ability and motivation etc explaining the wages, that are correlate with education  $\rightarrow$  OVB. We expect a positive omitted variable bias, since ability is likely correlated with log wages Testing relevance condition by FS regression: running education on number of siblings, this is significantly different from 0 and f-stat >10 Then running IV regression, we find the instrument has strong enough F-stat, but it could be that the exclusion restriction is violated, before we found coefficient of 0.059, with IV we find 0.122 (12%), which is higher than OLS, revealing inconsistency already, perhaps our assumption about *exogeneity* is not fulfilled.

**Tutorial 5.** Exercise 2 using sibs as iv is not same as plugging sibs into education (as in proxy), we find very different result from our IV estimator, that is big diff from 0 controlling for. Education and birth quarter negatively correlated? b) c) again, we get an increase than the OLS estimator, and larger than when we used siblings as IV. But do we have similar concerns now using birth order Is birth order endogenous? Like the number of siblings? The decision to have children might be related to budget constraints etc. d) identification assumption  $\log(wage) = \beta_0 + \beta_1$  test whether  $\pi_2$  is significantly different from 0, if we estimate our IV, we need to include all exogenous variables as instruments, we estimate a different coefficient.