

REGRESSION ANALYSIS AND CAUSALITY WITH R

Instrumental Variables

João Cerejeira¹ Miguel Portela^{1,2,3}

¹NIPE – UMinho

²IZA, Bonn

³Banco de Portugal

October 26, 2021

How to get a consistent estimator of β if X is correlated with ε ?

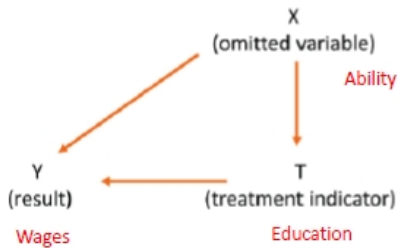
$$Y = X\beta + \varepsilon, \tag{1}$$

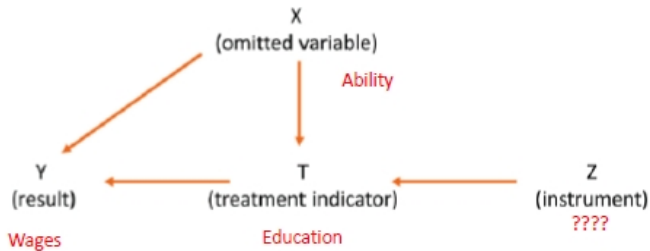
Intuition: total variation of X can be divided in two parts:

- one part correlated with ε ;
- other uncorrelated with ε .

IV strategy is to use just the second part of X variation.







$Z (N \times I)$ is the instruments matrix;

$X (N \times k)$ is the explanatory variables matrix;

If $k = I$, just-identified case;

If $k < I$, over-identified case;

If $k > I$ under-identified case - cannot be estimated.

All exogenous variables in X should be included in Z .

All other variables in X are endogenous.

Variables in Z not included in X are called exogenous variables.

Conditions for instrument validity

Must be correlated with X - testable:

$$COV(Z_i, X_i) \neq 0,$$

Must be uncorrelated with 'error' – untestable if $k = l$ – have to argue case for this assumption:

$$COV(Z_i, \varepsilon_i) = 0.$$

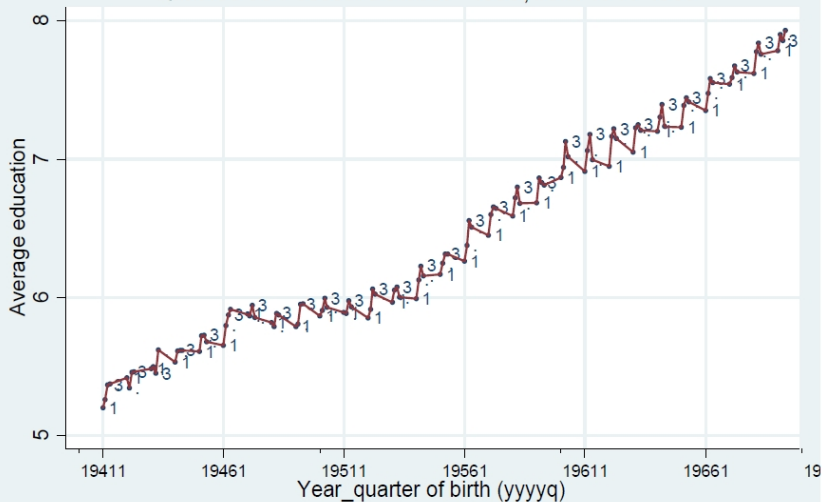
These conditions guaranteed with instrument for experimental data.
But more problematic for data from quasi-experiments.

Example: Causal effect of education on wages
Instrument: quarter of birth
(QP dataset, 1999, born 1940-1970)

Quarter of birth	Educ Mean	Freq.
1	6.555	383,593
2	6.669	360,218
3	6.749	352,951
4	6.684	367,036
Total	6.662	1,463,798

Average Education by quarter of birth

Quadros de Pessoal Dataset - 1999, born 1941-1969



2SLS (Two stage least squares method)

First Stage

Regress X on Z :

$$X = Z\Pi + v, \quad (2)$$

therefore

$$\hat{\Pi} = (Z'Z)^{-1}Z'X, \quad (3)$$

predict \hat{X} :

$$\hat{X} = Z\hat{\Pi} = Z(Z'Z)^{-1}Z'X, \quad (4)$$

2SLS (Two stage least squares method)

Second Stage

Regress Y on \hat{X} :

$$Y = \hat{X}\beta_{IV} + u, \quad (5)$$

$$\begin{aligned} \hat{\beta}_{IV} &= (\hat{X}'\hat{X})^{-1}\hat{X}'Y = \\ &= ((Z(Z'Z)^{-1}Z'X)'(Z(Z'Z)^{-1}Z'X))^{-1}(Z(Z'Z)^{-1}Z'X)'Y = \\ &= (X'Z(Z'Z)^{-1}Z'Z(Z'Z)^{-1}Z'X)^{-1}(X'Z(Z'Z)^{-1}Z')Y = \\ &= (X'Z(Z'Z)^{-1}Z'X)^{-1}X'Z(Z'Z)^{-1}Z'Y \end{aligned} \quad (6)$$

2SLS (Two stage least squares method)

If $l = k$, then $X'Z$ is an invertible square matrix:

$$\hat{\beta}_{IV} = (Z'X)^{-1}(Z'Z)(X'Z)^{-1}X'Z(Z'Z)^{-1}Z'Y = \quad (7)$$

$$= (Z'X)^{-1}Z'Y \quad (8)$$

Asymptotic variance of IV estimator

$$Y = \hat{X}\beta_{IV} + u, \quad (9)$$

and

$$\hat{X} = Z(Z'Z)^{-1}Z'X, \quad (10)$$

then:

$$\begin{aligned} \text{Var}(\hat{\beta}_{IV}) &= \hat{\sigma}^2(\hat{X}'\hat{X})^{-1} = \hat{\sigma}^2((Z(Z'Z)^{-1}Z'X)'(Z(Z'Z)^{-1}Z'X))^{-1} = \\ &= \hat{\sigma}^2(X'Z(Z'Z)^{-1}Z'Z(Z'Z)^{-1}Z'X)^{-1} = \hat{\sigma}^2(X'Z(Z'Z)^{-1}Z'X)^{-1}. \end{aligned}$$

If $l = k$:

$$\text{Var}(\hat{\beta}_{IV}) = \hat{\sigma}^2(X'Z)^{-1}Z'Z(Z'X)^{-1}.$$

Asymptotic variance of IV estimator

The variance of $\hat{\beta}_{IV}$ will be large if the correlation between X with Z is small. This is the "weak instruments" problem. Note: if $X = Z$, then $Var(\hat{\beta}_{IV}) = Var(\hat{\beta}_{OLS})$. Therefore $Var(\hat{\beta}_{IV}) \geq Var(\hat{\beta}_{OLS})$.

Note that $\hat{\sigma}^2$ is estimated using:

$$\hat{\sigma}^2 = (Y - X\hat{\beta}_{IV})'((Y - X\hat{\beta}_{IV}))/ (N - K),$$

not:

$$\hat{\sigma}^2 = (Y - \hat{X}\hat{\beta}_{IV})'((Y - \hat{X}\hat{\beta}_{IV}))/ (N - K).$$

The Problem of Weak Instruments

- Say that instruments are 'weak' if correlation between X and Z low (after inclusion of other exogenous variables).
- Rule of thumb - if F-statistic on instruments in first-stage less than 10 ($t < 3.33$) then may be problem.
- A whole range of problems tend to arise if instruments are weak.

The Problem of Weak Instruments

- Asymptotic problems:

The Problem of Weak Instruments

- Asymptotic problems:
 - High asymptotic variance.

The Problem of Weak Instruments

- Asymptotic problems:
 - High asymptotic variance.
 - Small departures from instrument exogeneity lead to big inconsistencies.

The Problem of Weak Instruments

- Asymptotic problems:
 - High asymptotic variance.
 - Small departures from instrument exogeneity lead to big inconsistencies.
- Finite-Sample Problems:

The Problem of Weak Instruments

- Asymptotic problems:
 - High asymptotic variance.
 - Small departures from instrument exogeneity lead to big inconsistencies.
- Finite-Sample Problems:
 - Small-sample distribution may be very different from asymptotic one.

The Problem of Weak Instruments

- Asymptotic problems:
 - High asymptotic variance.
 - Small departures from instrument exogeneity lead to big inconsistencies.
- Finite-Sample Problems:
 - Small-sample distribution may be very different from asymptotic one.
 - May be large bias.

The Problem of Weak Instruments

- Asymptotic problems:
 - High asymptotic variance.
 - Small departures from instrument exogeneity lead to big inconsistencies.
- Finite-Sample Problems:
 - Small-sample distribution may be very different from asymptotic one.
 - May be large bias.
 - Computed variance may be wrong.

The Problem of Weak Instruments

- Asymptotic problems:
 - High asymptotic variance.
 - Small departures from instrument exogeneity lead to big inconsistencies.
- Finite-Sample Problems:
 - Small-sample distribution may be very different from asymptotic one.
 - May be large bias.
 - Computed variance may be wrong.
 - Distribution may be very different from normal.

Testing overidentification restrictions

If the number of excluded instruments is equal to the number of endogenous variables it is not possible to test if these instruments are correlated with the error term.

In the overidentified case, we can test if some instruments are correlated with the error, but we don't know which.

To construct this test you:

- estimate the model by IV and get the residuals;
- run a regression of these residuals on all exogenous variables and get $n \times R^2$.
- under the null that all instruments are not correlated with the disturbance term, the LM statistic of $n \times R^2$ follows a $\chi^2(l - k)$ distribution where $l - k$ is the number of restrictions of overidentification (the number of instruments in excess).