

Instrumental Variables

Instrumental variables is predicated on the idea of a system of equations. In our running example, we'll use a standard earnings/education equation, where x indicates the number of years of education for an individual. We're interested in figuring out the impact of an additional year of education on log earnings.

The problem is the familiar one of selection bias, or endogeneity. Individuals who are more likely to earn more may seek out more years of education. This means that those individuals with more education also have more earnings, but that it could be the case that giving more people more education would not increase earnings.

To sort this out, we need to find something about the person's environment or traits that could increase their years of education but *only* impacts earnings through the mechanism of education.

If we had an experiment where some people were assigned to get more education, then we could use the experimental assignment variable to predict years of education, and there would be no reason that experimental assignment predicted earnings, except through the mechanism of

Two-Stage Least Squares

The primary estimator for instrumental variables is Two Stage Least Squares, often shatterer to 2SLS. It works this way. Let's say we have an endogenous regressor x (years of education), an instrument z and a set of controls c .

$$x_i = \gamma_0 + \gamma_1 z_i + \boldsymbol{\gamma} \mathbf{c}_i + \mu_i \quad (1)$$

This is called the first stage. We then take predictions of x_i from the first stage:

$$\hat{x}_i = \gamma_0 + \gamma_1 z_i + \boldsymbol{\gamma} \mathbf{c}_i \quad (2)$$

And plug those into the equation we're actually interested in estimating:

$$y_i = \beta_0 + \beta_1 \hat{x}_i + \boldsymbol{\beta} \mathbf{c}_i + \epsilon_i \quad (3)$$

There's an adjustment that needs to be made to the standard errors in the above regression, but with that, β_1 is our 2SLS estimator of the impact of x on y . Notice what the first stage does: it takes everything out of x that isn't a function of z or c .

This “version” of x will have no correlation with the error term in the second stage by construction, which is exactly what we want.

IV estimates are easy to run– the hard part is establish that the instruments are valid.

In a theoretical sense, what you need is an instrument that’s truly exogenous– something that people are not choosing (or aren’t choosing with the outcome in mind), but affects assignment to treatment, and ONLY affects the outcome through assignment to treatment. Most of the time you need to know the conditions of the individuals in your sample quite well to establish the validity of the instrument.

In an empirical sense, there are two equally important properties of an instrument. First, it must strongly predict assignment to treatment. Second, it must not impact the outcome except through the treatment variable. We can’t always demonstrate both of these assumptions are held.

Establishing Endogeneity

Even suspected endogeneity is sufficient to use IV, but you can try to establish endogeneity by seeing whether the residuals from the first stage predict the outcome in the basic form of the second stage. This is known as the Durbin-Wu-Hausman test. If it’s not significant (by a lot), then there’s less evidence of endogeneity between x and y . However, this test is predicated on the assumption that the instrument is valid.

Establishing the Strength of the Instrument(s)

If the instruments aren’t strongly related to the endogenous regressor, we have “weak” instruments. These will result in both high levels of variance in the second stage *AND* biased estimates in the second stage. That’s bad. What you need to show is that the instruments strongly predict the endogenous regressor. A basic tests is the F test for the first stage, testing the linear restriction when the instruments are left out of this equation. A more advanced test is the Stock & Yogo minimum eigenvalue test.

Establishing Overidentification

When you have just one instrument, the system of equations is said to be “exactly identified.” When you have more instruments than endogenous regressors, the system is “overidentified.” Overidentification can be helpful because it tests whether the instruments have an impact on the outcome through another mechanism beyond just the endogenous regressor. The basic Sargan test works by taking the residuals from the first stage, then using them as a predictor for the outcome. The R^2 from that regression times N is distributed χ^2 with degrees of freedom equal to the number of instruments minus 1.

Local Average Treatment Effects

Our understanding of IV has become more sophisticated over time. One of the key insights from recent literature is that IV estimates are Local Average Treatment Effects: they apply only to that part of the population that was induced into treatment by the instrument. The literature has identified the key assumptions that are crucial for this understanding:

1. Independence assumption: z is as good as randomly assigned, conditional on x
2. Exclusion: y is impacted by z only through x .
3. Monotonicity: x can only stay the same or increase as a function of z , z can never decrease x . (or vice-versa)

What's important to remember in practice is that the external validity of any IV estimate is limited to the group of people who could have been induced by the instrument into treatment.