

Instrumental Variables

When any of the regressors is correlated with the disturbances, OLS estimates are biased and inconsistent.

Most common solution is **instrumental variables**, or **IV**, estimation.

Suppose there is just one column of \mathbf{X} , say \mathbf{x} , that is correlated with \mathbf{u} .

Let \mathbf{W} denote a matrix of instruments with the same dimension as \mathbf{X} , where \mathbf{x} has been replaced by an instrument \mathbf{w} that is assumed to be uncorrelated with the disturbances but correlated with \mathbf{x} .

The remaining columns of \mathbf{W} and \mathbf{X} are identical.

The **simple IV estimator** is

$$\hat{\beta}_{IV} = (\mathbf{W}^\top \mathbf{X})^{-1} \mathbf{W}^\top \mathbf{y}. \quad (1)$$

This is consistent whenever $\text{plim } N^{-1} \mathbf{W}^\top \mathbf{u} = \mathbf{0}$.

Errors in Variables

Correlations between regressors and disturbances arise for two main reasons. The first is that regressors may be measured with error.

Consider the simple linear regression model

$$y_i^\circ = \beta_1 + \beta_2 x_i^\circ + u_i^\circ, \quad u_i^\circ \sim \text{IID}(0, \sigma^2), \quad (2)$$

where the variables x_i° and y_i° are not actually observed. These are sometimes called **latent variables**. Instead, we observe

$$x_i = x_i^\circ + v_{1i}, \quad \text{and} \quad y_i = y_i^\circ + v_{2i}. \quad (3)$$

Here v_{1i} and v_{2i} are measurement errors which are assumed (not always realistically) to be IID with variances ω_1^2 and ω_2^2 , respectively, and to be independent of x_i° , y_i° , and u_i° .

From (3), we see that $x_i^\circ = x_i - v_{1i}$ and $y_i^\circ = y_i - v_{2i}$.

In terms of observables, (2) is replaced by

$$y_i = \beta_1 + \beta_2(x_i - v_{1i}) + u_i^\circ + v_{2i} \quad (4)$$

$$= \beta_1 + \beta_2 x_i + u_i^\circ + v_{2i} - \beta_2 v_{1i} \quad (5)$$

$$= \beta_1 + \beta_2 x_i + u_i, \quad (6)$$

where $u_i = u_i^\circ + v_{2i} - \beta_2 v_{1i}$. Thus

$$\text{Var}(u_i) = \sigma^2 + \omega_2^2 + \beta_2^2 \omega_1^2 > \sigma^2. \quad (7)$$

Because $x_i = x_i^\circ + v_{1i}$, and u_i depends on v_{1i} , u_i must be correlated with x_i whenever $\beta_2 \neq 0$ and $\omega_1^2 > 0$.

Since $E(u_i | x_i) = E(u_i | v_{1i}) = -\beta_2 v_{1i}$,

$$\text{Cov}(x_i, u_i) = E(x_i u_i) = E(x_i E(u_i | x_i)) \quad (8)$$

$$= -E((x_i^\circ + v_{1i})\beta_2 v_{1i}) = -\beta_2 \omega_1^2. \quad (9)$$

The correlation between x_i and u_i has sign opposite to that of β_2 .

Because of the negative correlation between x_i and u_i , the OLS estimate $\hat{\beta}_2$ is biased towards zero.

$\hat{\beta}_1$ is also biased, as are any other coefficients.

- Friedman's **permanent income hypothesis** is essentially a measurement error story. Consumption depends on permanent income, but we only observe current income.
- Regressing consumption on current income leads to an estimated **marginal propensity to consume** that is much too small.

If we knew how large the measurement errors were, we could remove the bias in $\hat{\beta}_2$.

Of course, this will require the assumption that they are uncorrelated with the x_i° , which we made in (3).

We can also obtain consistent estimates if we can find an instrument w_i that is correlated with x_i° but uncorrelated with v_{1i} .

Simultaneous Equations

Suppose that q_i is quantity and p_i is price, both of which might be in logarithms. A linear (or loglinear) model of demand and supply is

$$q_i = \gamma_d p_i + \mathbf{X}_i^d \boldsymbol{\beta}_d + u_i^d \quad (10)$$

$$q_i = \gamma_s p_i + \mathbf{X}_i^s \boldsymbol{\beta}_s + u_i^s, \quad (11)$$

where (10) is the demand function and (11) is the supply function. These are the two **structural equations** of the system.

Here \mathbf{X}_i^d and \mathbf{X}_i^s are row vectors of observations on exogenous or predetermined variables.

Equations (10) and (11) are two linear equations that jointly determine p_i and q_i . They constitute a **linear simultaneous equations model**.

As written, quantity depends on price in both equations. But either or both equations could be rewritten so that price depends on quantity.

We can rewrite the two equations in matrix notation as

$$\begin{bmatrix} 1 & -\gamma_d \\ 1 & -\gamma_s \end{bmatrix} \begin{bmatrix} q_i \\ p_i \end{bmatrix} = \begin{bmatrix} \mathbf{X}_i^d \boldsymbol{\beta}_d \\ \mathbf{X}_i^s \boldsymbol{\beta}_s \end{bmatrix} + \begin{bmatrix} u_i^d \\ u_i^s \end{bmatrix}. \quad (12)$$

Provided $\gamma_d \neq \gamma_s$, a solution exists. It is

$$\begin{bmatrix} q_i \\ p_i \end{bmatrix} = \begin{bmatrix} 1 & -\gamma_d \\ 1 & -\gamma_s \end{bmatrix}^{-1} \left(\begin{bmatrix} \mathbf{X}_i^d \boldsymbol{\beta}_d \\ \mathbf{X}_i^s \boldsymbol{\beta}_s \end{bmatrix} + \begin{bmatrix} u_i^d \\ u_i^s \end{bmatrix} \right). \quad (13)$$

This is the **restricted reduced form** of the system.

Here p_i and q_i depend on both u_i^d and u_i^s , and on every exogenous or predetermined variable in \mathbf{X}_i^d and/or \mathbf{X}_i^s .

Therefore p_i , which appears on the r.h.s. of both (10) and (11), must be correlated with the disturbances in both of those equations.

This is also true of q_i . Rewriting one or both equations will not eliminate the problem.

Instrumental Variables (IV) Estimators

We will focus on the linear regression model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} = \boldsymbol{\beta}\mathbf{y}_2 + \mathbf{Z}\boldsymbol{\gamma} + \mathbf{u}, \quad \mathbb{E}(\mathbf{u}\mathbf{u}^\top) = \sigma^2\mathbf{I}, \quad (14)$$

where \mathbf{X} is $N \times k$, and \mathbf{Z} is $N \times (k - 1)$. Here every column of \mathbf{Z} is exogenous or predetermined, but \mathbf{y}_2 is not predetermined.

Let \mathbf{W} be an $N \times k$ matrix consisting of \mathbf{Z} plus one other column, say w_2 , that is predetermined and therefore a **valid instrument**.

Since $\mathbb{E}(\mathbf{W}^\top\mathbf{u}) = \mathbf{0}$, we can employ the moment conditions

$$\mathbf{W}^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{0}. \quad (15)$$

These lead to the **simple IV estimator**

$$\hat{\boldsymbol{\beta}}_{\text{IV}} = (\mathbf{W}^\top\mathbf{X})^{-1}\mathbf{W}^\top\mathbf{y}. \quad (16)$$

The simple IV estimator is consistent whenever

$$\text{plim}_{N \rightarrow \infty} \frac{1}{N} \mathbf{W}^\top \mathbf{X} = \mathbf{S}_{\mathbf{W}^\top \mathbf{X}} \quad \text{and} \quad \text{plim}_{N \rightarrow \infty} \frac{1}{N} \mathbf{W}^\top \mathbf{u} = \mathbf{0}, \quad (17)$$

where $\mathbf{S}_{\mathbf{W}^\top \mathbf{X}}$ is deterministic and nonsingular. Just replace \mathbf{y} in (16) by $\mathbf{X}\boldsymbol{\beta}_0 + \mathbf{u}$, divide both factors by N , and take plims.

The asymptotic covariance matrix of $\hat{\boldsymbol{\beta}}_{\text{IV}}$ is

$$\text{Var}(N^{1/2}(\hat{\boldsymbol{\beta}}_{\text{IV}} - \boldsymbol{\beta}_0)) = \sigma^2 (\mathbf{S}_{\mathbf{W}^\top \mathbf{X}})^{-1} \mathbf{S}_{\mathbf{W}^\top \mathbf{W}} (\mathbf{S}_{\mathbf{W}^\top \mathbf{X}})^{-1}, \quad (18)$$

where $\mathbf{S}_{\mathbf{W}^\top \mathbf{W}} = \text{plim}_{N \rightarrow \infty} N^{-1} \mathbf{W}^\top \mathbf{W}$. In practice, we use

$$\widehat{\text{Var}}(\hat{\boldsymbol{\beta}}_{\text{IV}}) = \hat{\sigma}^2 (\mathbf{X}^\top \mathbf{P}_\mathbf{W} \mathbf{X})^{-1}, \quad (19)$$

where

$$\hat{\sigma}^2 = \frac{1}{N} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{\text{IV}})^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{\text{IV}}). \quad (20)$$

Of course, if we assumed that the u_i are heteroskedastic rather than i.i.d., (19) would be replaced by

$$\widehat{\text{Var}}_h(\hat{\beta}_{\text{IV}}) = (\mathbf{W}^\top \mathbf{X})^{-1} \left(\sum_{i=1}^N \hat{u}_i^2 \mathbf{W}_i^\top \mathbf{W}_i \right) (\mathbf{W}^\top \mathbf{X})^{-1}. \quad (21)$$

We could also allow for clustering. Since the i.i.d. assumption leads to some interesting results on efficiency, we will maintain it for now.

The simple IV estimator $\hat{\beta}_{\text{IV}}$ can also be written as $(\mathbf{X}^\top \mathbf{P}_W \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{P}_W \mathbf{y}$.

Because $\mathbf{X}^\top \mathbf{W}$ is square, $(\mathbf{X}^\top \mathbf{P}_W \mathbf{X})^{-1} = (\mathbf{W}^\top \mathbf{X})^{-1} \mathbf{W}^\top \mathbf{W} (\mathbf{X}^\top \mathbf{W})^{-1}$. Thus

$$\hat{\beta}_{\text{IV}} = (\mathbf{X}^\top \mathbf{P}_W \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{P}_W \mathbf{y} \quad (22)$$

$$= (\mathbf{W}^\top \mathbf{X})^{-1} \mathbf{W}^\top \mathbf{W} (\mathbf{X}^\top \mathbf{W})^{-1} \mathbf{X}^\top \mathbf{W} (\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{y} \quad (23)$$

$$= (\mathbf{W}^\top \mathbf{X})^{-1} \mathbf{W}^\top \mathbf{y} \quad (24)$$

The first expression is also the **generalized IV estimator**, or **GIVE**.

In many cases, we have $l > k$ instruments. The GIVE finds a linear combination of them that is optimal.

Suppose that

$$\mathbf{X} = \bar{\mathbf{X}} + \mathbf{V}, \text{ where } \bar{X}_i = E(X_i | \Omega_i). \quad (25)$$

Then \mathbf{W} should, ideally, be proportional to $\bar{\mathbf{X}}$. In this case

$$\text{plim}_{N \rightarrow \infty} \frac{1}{N} \mathbf{X}^\top \mathbf{P}_W \mathbf{X} = \text{plim}_{N \rightarrow \infty} \frac{1}{N} \bar{\mathbf{X}}^\top \mathbf{P}_W \bar{\mathbf{X}} = \text{plim}_{N \rightarrow \infty} \frac{1}{N} \bar{\mathbf{X}}^\top \bar{\mathbf{X}}. \quad (26)$$

Of course, if $\text{Var}(\mathbf{V})$ is large, $\bar{\mathbf{X}}^\top \bar{\mathbf{X}}$ may be much smaller than $\mathbf{X}^\top \mathbf{X}$, and the IV estimator may be much less efficient than the OLS estimator.

The optimal instruments are given by $\bar{\mathbf{X}}$, which we do not observe.

But if we redefine \mathbf{W} as an $N \times l$ matrix, we can use $\mathbf{P}_W \mathbf{X}$ to estimate $\bar{\mathbf{X}}$:

$$\hat{\beta}_{\text{IV}} = (\mathbf{X}^\top \mathbf{P}_W \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{P}_W \mathbf{y}. \quad [\text{GIVE again}] \quad (27)$$

Standard asymptotic results for GIVE are similar to those for simple IV. $S_{W^T X}$ is now $l \times k$ instead of $k \times k$, and $S_{W^T W}$ is now $l \times l$. The former should have rank k , the latter rank l .

It is still essential that

$$\text{plim}_{N \rightarrow \infty} \frac{1}{N} W^T u = \mathbf{0}. \quad (28)$$

Including even one invalid instrument in W will destroy consistency.

The most efficient possible GIV estimator uses \bar{X} as the matrix of instruments. It is less efficient than OLS because

$$\text{plim}_{N \rightarrow \infty} \frac{1}{N} X^T X = \text{plim}_{N \rightarrow \infty} \frac{1}{N} \bar{X}^T \bar{X} + \text{plim}_{N \rightarrow \infty} \frac{1}{N} V^T V. \quad (29)$$

This follows from (25). Terms involving $V^T \bar{X}$ vanish.

Adding additional (valid) instruments to W brings $X^T P_W X$ closer to $\bar{X}^T \bar{X}$ and thus increases asymptotic efficiency. But $X^T P_W X$ may still be much “smaller” than $X^T X$, so that $\hat{\beta}_{IV}$ is much less efficient than $\hat{\beta}_{OLS}$.

The generalized IV estimator can be obtained by minimizing the **IV criterion function**

$$Q(\boldsymbol{\beta}, \mathbf{y}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \mathbf{P}_W (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}). \quad (30)$$

The resulting moment conditions are

$$\mathbf{X}^\top \mathbf{P}_W (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{0}. \quad (31)$$

Intuitively, we only minimize the part of the residuals that can be explained by W .

In the i.i.d. case, $\hat{\boldsymbol{\beta}}_{IV}$ will be asymptotically normal if it is root- N consistent and

$$N^{-1/2} \mathbf{X}^\top \mathbf{P}_W \mathbf{u} \stackrel{a}{\sim} N \left(\mathbf{0}, \sigma_0^2 \text{plim}_{N \rightarrow \infty} \frac{1}{N} \mathbf{X}^\top \mathbf{P}_W \mathbf{X} \right). \quad (32)$$

With heteroskedasticity and/or clustering, the covariance matrix of $N^{-1/2} \mathbf{X}^\top \mathbf{P}_W \mathbf{u}$ will be more complicated.

Two-Stage Least Squares

The (generalized) IV estimator is also commonly known as the **two-stage least-squares**, or **2SLS**, estimator.

It can be calculated using OLS regressions in two stages.

In the first stage, each column x_i , $i = 1, \dots, k$, of X is regressed on W , if necessary. If x_i is a valid instrument, it is already one of the columns of W , so it serves as its own instrument without a regression.

The second-stage regression is

$$y = P_W X \beta + \text{residuals.} \quad (33)$$

Because P_W is idempotent, the OLS estimate of β from this second-stage regression is

$$\hat{\beta}_{2sls} = (X^\top P_W X)^{-1} X^\top P_W y, \quad (34)$$

which is identical to $\hat{\beta}_{IV}$.

Estimating the Variance of the Disturbances

If 2SLS is used, it is tricky to estimate the standard error of the regression and the covariance matrix of the parameter estimates.

The OLS estimate of σ^2 from (33) is

$$s^2 = \frac{\|\mathbf{y} - \mathbf{P}_W \mathbf{X} \hat{\boldsymbol{\beta}}_{IV}\|^2}{N - k}. \quad (35)$$

In contrast, the estimate that was used in the estimated IV covariance matrix is

$$\hat{\sigma}^2 = \frac{\|\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}_{IV}\|^2}{N}. \quad (36)$$

These two estimates of σ^2 are not asymptotically equivalent, and s^2 is not consistent.

Residuals from (33) do not tend to \mathbf{u} as $N \rightarrow \infty$.

This happens because the regressors $P_W X$ are not the true explanatory variables X .

We divided by N in (36) because the correct IV residuals, $y - X\hat{\beta}_{IV}$, are not necessarily too small.

Cluster-Robust Inference

When the disturbances display intra-cluster correlation, the scores almost always will too.

We need to use the CRVE

$$(X^\top P_W X)^{-1} \left(\sum_{g=1}^G [X^\top P_W]_g \hat{u}_g \hat{u}_g^\top [P_W X]_g \right) (X^\top P_W X)^{-1}, \quad (37)$$

where $[P_W X]_g$ is the $N_g \times k$ submatrix of $P_W X$ for cluster g , and \hat{u}_g is the corresponding $N_g \times 1$ subvector of \hat{u} .

Heteroskedasticity-robust variance estimator is a special case of (37) with $G = N$.

Tests of Over-Identifying Restrictions

A model with k regressors and l instruments implicitly incorporates $l - k$ **over-identifying restrictions**. We cannot solve l equations uniquely for k unknowns, so we only solve k moment conditions.

We can always write

$$\mathcal{S}(W) = \mathcal{S}(P_W X, W^*), \quad (38)$$

where W^* has $l - k$ columns. Recall that

$$P_W X = W(W^\top W)^{-1} W^\top X \quad (39)$$

is just a linear combination of the columns of W . Luckily, we don't have to decide which columns of W belong to W^* .

Although we cannot test the moment conditions that $E(X^\top P_W u) = \mathbf{0}$, we can and should test the moment conditions that $W^{*\top} u = \mathbf{0}$.

In the i.i.d. case, we can test restrictions on β by evaluating the IV criterion function (30) at the restricted and unrestricted estimates to get $Q(\tilde{\beta})$ and $Q(\hat{\beta})$. Then

$$\frac{Q(\tilde{\beta}) - Q(\hat{\beta})}{\sigma_{\text{IV}}^2} \stackrel{a}{\sim} \chi^2(r). \quad (40)$$

Now consider the IV regression models

$$y = X\beta + u, \quad u \sim \text{IID}(\mathbf{0}, \sigma^2 \mathbf{I}), \quad E(W^\top u) = \mathbf{0}, \quad \text{and} \quad (41)$$

$$y = X\beta + W^* \gamma + u, \quad u \sim \text{IID}(\mathbf{0}, \sigma^2 \mathbf{I}), \quad E(W^\top u) = \mathbf{0}. \quad (42)$$

Call the estimates from (41) $\hat{\beta}_{\text{IV}}$ and the ones from (42) $\hat{\beta}_{\text{U}}$. For (41),

$$Q(\hat{\beta}_{\text{IV}}) = \hat{u}_{\text{IV}}^\top P_W \hat{u}_{\text{IV}}. \quad (43)$$

For (42), by contrast, $Q(\hat{\beta}_{\text{U}}) = \hat{u}_{\text{U}}^\top P_W \hat{u}_{\text{U}} = 0$. This holds because $\mathcal{S}(W) = \mathcal{S}(P_W X, W^*)$, so that \hat{u}_{U} is orthogonal to W .

Thus we don't need to estimate the model (42) at all. We only used it to justify a particular test statistic for over-identifying restrictions.

In this case, the test statistic (40) reduces to

$$\frac{Q(\hat{\beta}_{IV})}{\hat{\sigma}_{IV}^2} = \frac{\hat{u}_{IV}^\top \mathbf{P}_W \hat{u}_{IV}}{\hat{\sigma}_{IV}^2} \stackrel{a}{=} \frac{\mathbf{u}^\top (\mathbf{P}_W - \mathbf{P}_{P_W X}) \mathbf{u}}{\sigma^2}, \quad (44)$$

because $\hat{u}_{IV} = \mathbf{u} - \mathbf{X}(\mathbf{X}^\top \mathbf{P}_W \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{P}_W \mathbf{u}$. The numerator is a quadratic form in a matrix that projects onto a space of dimension $l - k$.

The denominator ensures that, asymptotically, the elements of \mathbf{u}/σ have variance 1. Thus (44) is asymptotically $\chi^2(l - k)$.

The **Sargan statistic** to test for over-identifying restrictions is the middle expression in (44).

It is simply the ESS from a regression of \hat{u}_{IV} on \mathbf{W} , divided by $\hat{\sigma}_{IV}^2$.

This works because $\hat{u}_{IV}^\top \mathbf{P}_W \hat{u}_{IV} = \mathbf{u}^\top (\mathbf{P}_W - \mathbf{P}_{P_W X}) \mathbf{u}$; see (44).

The model (42) represents two conceptually different alternatives. The “true” parameter vector γ could be nonzero in two situations.

- 1 The model (14) is correctly specified, but some of the instruments are asymptotically correlated with the disturbances and are therefore not valid instruments.
- 2 The model (14) is misspecified, and some of the instruments (or, possibly, other variables that are correlated with them) have incorrectly been omitted from the regression function.

In either case, the over-identification test statistic should lead us to reject the null hypothesis whenever the sample size is large enough.

The Sargan statistic can be generalized to allow for heteroskedasticity and/or clustering. For GMM, it becomes the **Hansen-Sargan statistic**, which is the minimized value of the GMM criterion function.

These tests may have poor finite-sample properties, so it can be good to bootstrap them; see Davidson and MacKinnon (2015).

Durbin-Wu-Hausman Tests

When do we actually need to use instrumental variables?

Some variables may be measured with error, but how large are the errors? Will they cause enough inconsistency to worry about?

Is a certain explanatory variable actually endogenous? If so, how much inconsistency will this cause?

It may be perfectly reasonable to employ OLS estimation if the disturbances are not very correlated with the regressors.

We can test whether there is correlation by using a test due to Durbin (1954), Wu (1974), and Hausman (1978).

The null and alternative hypotheses for the DWH test are

$$H_0: \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}, \quad \mathbf{u} \sim \text{IID}(\mathbf{0}, \sigma^2 \mathbf{I}), \quad \text{E}(\mathbf{X}^\top \mathbf{u}) = \mathbf{0}, \quad \text{and} \quad (45)$$

$$H_1: \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}, \quad \mathbf{u} \sim \text{IID}(\mathbf{0}, \sigma^2 \mathbf{I}), \quad \text{E}(\mathbf{W}^\top \mathbf{u}) = \mathbf{0}. \quad (46)$$

Under H_1 , $\hat{\beta}_{IV}$ is consistent, but $\hat{\beta}_{OLS}$ is not. Under H_0 , both are consistent, but $\hat{\beta}_{OLS}$ is more efficient.

Evidently, $\text{plim}(\hat{\beta}_{IV} - \hat{\beta}_{OLS})$ is zero under the null and nonzero under the alternative.

The **vector of contrasts** is

$$\hat{\beta}_{IV} - \hat{\beta}_{OLS} = (\mathbf{X}^\top \mathbf{P}_W \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{P}_W \mathbf{y} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad (47)$$

$$= (\mathbf{X}^\top \mathbf{P}_W \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{P}_W \mathbf{y} - \mathbf{X}^\top \mathbf{P}_W \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}) \quad (48)$$

$$= (\mathbf{X}^\top \mathbf{P}_W \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{P}_W (\mathbf{I} - \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \mathbf{y} \quad (49)$$

$$= (\mathbf{X}^\top \mathbf{P}_W \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{P}_W \mathbf{M}_X \mathbf{y}. \quad (50)$$

Here $(\mathbf{X}^\top \mathbf{P}_W \mathbf{X})^{-1}$ is a positive definite matrix.

Testing whether $\hat{\beta}_{IV} - \hat{\beta}_{OLS} = \mathbf{0}$ is equivalent to testing whether $\text{E}(\mathbf{X}^\top \mathbf{P}_W \mathbf{M}_X \mathbf{y}) = \mathbf{0}$. Since $\mathbf{M}_X \mathbf{y} = \mathbf{M}_X \mathbf{u}$, it does if $\mathbf{P}_W \mathbf{X}$ is uncorrelated with $\mathbf{M}_X \mathbf{u}$.

Let us partition X as $[Z \ Y]$, where the k_1 columns of Z belong to W , and the $k_2 = k - k_1$ columns of Y are treated as potentially endogenous.

It is always the case that $Z^\top \hat{u} = Z^\top P_W M_X \hat{u} = \mathbf{0}$; here $\hat{u} = \hat{u}_{OLS}$.

But it will almost never be the case $Y^\top P_W M_X \hat{u} = \mathbf{0}$. We can test whether this vector differs significantly from $\mathbf{0}$.

The easiest way to test whether $E(Y^\top P_W M_X y) = \mathbf{0}$ is to use an F test for the k_2 restrictions $\delta = \mathbf{0}$ in the OLS regression

$$y = X\beta + P_W Y\delta + u. \quad (51)$$

The OLS estimates of δ are, by the FWL Theorem, the same as those from the FWL regression of $M_X y$ on $M_X P_W Y$, that is,

$$\hat{\delta} = (Y^\top P_W M_X P_W Y)^{-1} Y^\top P_W M_X y. \quad (52)$$

Testing whether $\delta = \mathbf{0}$ is equivalent to testing whether $E(Y^\top P_W M_X y) = \mathbf{0}$. We can use an F test with k_2 and $N - k - k_2$ degrees of freedom.

It is not hard to see that an F test of $\zeta = \mathbf{0}$ in the regression

$$y = X\beta + M_W Y \zeta + u \quad (53)$$

is numerically identical to the F test for $\delta = \mathbf{0}$ in (51).

This follows from the fact that $\mathcal{S}(Y, P_W Y)$ and $\mathcal{S}(Y, M_W Y)$ are the same subspace, so that (51) and (53) fit the same.

If the DWH test rejects, there are two possibilities.

- ① At least one column of Y is endogenous.
- ② One or more of columns of W have incorrectly been omitted from Z ; some instruments should have been treated as regressors.

If the DWH test does not reject, then we may feel justified in using OLS. But classic pre-test issues arise when we use this test (Guggenberger, 2010a, 2010b).

Tests Based on Vectors of Contrasts

Hausman (1978) proposed a test that can be used whenever there are two root- N consistent estimators.

One of them, say $\hat{\theta}_I$, is like $\hat{\beta}_{IV}$. It is inefficient but consistent under relatively weak conditions.

The other, say $\hat{\theta}_E$, is like $\hat{\beta}_{OLS}$. It may be inconsistent, but it is more efficient when it is consistent.

In a broad range of cases, we can write

$$N^{1/2}(\hat{\theta}_I - \theta_0) \stackrel{a}{=} N^{1/2}(\hat{\theta}_E - \theta_0) + v, \quad (54)$$

where v is a random k -vector that is uncorrelated with $N^{1/2}(\hat{\theta}_E - \theta_0)$.

Recall the proof of the Gauss-Markov Theorem.

The vector v is asymptotically equal to $N^{1/2}$ times the vector of contrasts, which is just $\hat{\theta}_I - \hat{\theta}_E$.

Because v is not correlated with $N^{1/2}(\hat{\theta}_E - \theta_0)$,

$$\text{Var}(v) \stackrel{a}{=} \text{Var}(N^{1/2}(\hat{\theta}_I - \theta_0)) - \text{Var}(N^{1/2}(\hat{\theta}_E - \theta_0)). \quad (55)$$

This is Hausman's key result. Therefore, a suitable test statistic is

$$(\hat{\theta}_I - \hat{\theta}_E)^\top (\widehat{\text{Var}}(\hat{\theta}_I) - \widehat{\text{Var}}(\hat{\theta}_E))^{-1} (\hat{\theta}_I - \hat{\theta}_E). \quad (56)$$

However, Hausman forgot that the covariance matrix that is inverted in the middle of (56) may not have full rank.

It does not have full rank in the DWH case if we define θ as β . We need to define it as the coefficient vector for Y , which has k_2 elements.

Care needs to be taken when using Hausman tests based on (56). Ideally, the rank of $\text{Var}(v)$ should be known. It is the number of degrees of freedom for the test based on (56).

The numerical rank of $\widehat{\text{Var}}(\hat{\theta}_I) - \widehat{\text{Var}}(\hat{\theta}_E)$ may be misleading.

GIV Estimation Using Control Functions

The estimate of β from the DWH test regression (53) is simply $\hat{\beta}_{IV}$.

Here the regressor(s) $M_W Y$ are called **control functions**.

There is no advantage to computing $\hat{\beta}_{IV}$ in this way for linear regression models, but control function estimators can be useful for nonlinear models, such as logit and probit.

Care must be taken to obtain valid standard errors. The usual standard errors from (51) are not valid, because the control functions are **generated regressors**.

The easiest approach is often to use the bootstrap. The bootstrap samples may be obtained by resampling the rows of $[y, Y, X, W]$.

But bootstrap standard errors will be very misleading if $\hat{\beta}$ is exactly identified. They may perform poorly if there are not several over-identifying restrictions. Bootstrap IQR approach should work better.

We can demonstrate numerically that the DWH regression yields $\hat{\beta}_{IV}$. It is easiest to prove it if we assume that there is just one endogenous right-hand-side variable and one instrument:

$$\mathbf{y}_1 = \beta \mathbf{y}_2 + \mathbf{u} \quad (57)$$

$$\mathbf{y}_2 = \gamma \mathbf{w} + \mathbf{v}. \quad (58)$$

In this case, it is obvious that

$$\hat{\beta}_{IV} = (\mathbf{y}_2^\top \mathbf{P}_w \mathbf{y}_2)^{-1} \mathbf{y}_2^\top \mathbf{P}_w \mathbf{y}_1. \quad (59)$$

The DWH regression is

$$\mathbf{y}_1 = \beta \mathbf{y}_2 + \eta \mathbf{M}_w \mathbf{y}_2 + \mathbf{u}. \quad (60)$$

The OLS estimate of β , which is the control function estimate, is then easily seen to be

$$\hat{\beta}_{\text{CF}} = (\mathbf{y}_2^\top \mathbf{M}_{M_w \mathbf{y}_2} \mathbf{y}_2)^{-1} \mathbf{y}_2^\top \mathbf{M}_{M_w \mathbf{y}_2} \mathbf{y}_1. \quad (61)$$

Now observe that

$$\mathbf{M}_{M_w \mathbf{y}_2} = \mathbf{P}_{P_w \mathbf{y}_2}, \quad (62)$$

because $\mathbf{M}_{M_w \mathbf{y}_2}$ and $\mathbf{P}_{P_w \mathbf{y}_2}$ are orthogonal projection matrices. Projecting off the former is equivalent to projecting onto the latter. Thus

$$\hat{\beta}_{\text{CF}} = (\mathbf{y}_2^\top \mathbf{P}_{P_w \mathbf{y}_2} \mathbf{y}_2)^{-1} \mathbf{y}_2^\top \mathbf{P}_{P_w \mathbf{y}_2} \mathbf{y}_1. \quad (63)$$

Now observe that

$$\mathbf{P}_{P_w \mathbf{y}_2} = \mathbf{P}_{P_w \mathbf{y}_2} (\mathbf{y}_2^\top \mathbf{P}_{P_w \mathbf{y}_2})^{-1} \mathbf{y}_2^\top \mathbf{P}_w. \quad (64)$$

It follows that

$$\mathbf{y}_2^\top \mathbf{P}_{P_w \mathbf{y}_2} \mathbf{y}_2 = \mathbf{y}_2^\top \mathbf{P}_w \mathbf{y}_2 \quad \text{and} \quad \mathbf{y}_2^\top \mathbf{P}_{P_w \mathbf{y}_2} \mathbf{y}_1 = \mathbf{y}_2^\top \mathbf{P}_w \mathbf{y}_1. \quad (65)$$

When we substitute (65) into (63), we find that

$$\hat{\beta}_{\text{CF}} = (\mathbf{y}_2^\top \mathbf{P}_w \mathbf{y}_2)^{-1} \mathbf{y}_2^\top \mathbf{P}_w \mathbf{y}_1 = \hat{\beta}_{\text{IV}} \quad (66)$$

Because there is only one instrument, this expression can be simplified further. Observe that

$$(\mathbf{y}_2^\top \mathbf{P}_w \mathbf{y}_2)^{-1} = (\mathbf{y}_2^\top \mathbf{w})^{-1} \mathbf{w}^\top \mathbf{w} (\mathbf{w}^\top \mathbf{y}_2)^{-1}, \text{ and} \quad (67)$$

$$\mathbf{y}_2^\top \mathbf{P}_w \mathbf{y}_1 = \mathbf{y}_2^\top \mathbf{w} (\mathbf{w}^\top \mathbf{w})^{-1} \mathbf{w}^\top \mathbf{y}_1. \quad (68)$$

Thus we see that

$$\hat{\beta}_{\text{CF}} = \hat{\beta}_{\text{IV}} = \frac{\mathbf{y}_1^\top \mathbf{w}}{\mathbf{y}_2^\top \mathbf{w}} = (\mathbf{w}^\top \mathbf{y}_2)^{-1} \mathbf{w}^\top \mathbf{y}_1. \quad (69)$$

This is, of course, just the simple IV estimator for (57).

This result also holds when the structural equation includes a matrix of exogenous regressors, say \mathbf{Z} . We just replace \mathbf{w} by $\mathbf{M}_Z \mathbf{w}$ in (69).

The simple IV estimator in (69) has an interesting interpretation. If we regress y_1 on Z and w , we get the coefficient estimate

$$\frac{w^\top M_Z y_1}{w^\top M_Z w} \quad (70)$$

which is just the OLS estimate for the reduced-form equation for y_1 . If we regress y_2 on Z and w , we get the coefficient estimate

$$\frac{w^\top M_Z y_2}{w^\top M_Z w} \quad (71)$$

which is just the OLS estimate for the reduced-form equation for y_2 . Since these two OLS estimators have the same denominator, their ratio is just the IV estimator $w^\top M_Z y_1 / w^\top M_Z y_2$. So the IV estimate will be smaller than the OLS estimate whenever the coefficient (71) is larger than the coefficient (70).