

Assignment 2: IV and Asymptotic (b13)

A. Colin Cameron U.C.-Davis

Data are at <http://cameron.econ.ucdavis.edu/bgpe2013>

1. Use data in file `mus06data.dta`

We will do analysis similar to that in the slides, but with the only regressors the endogenous variable `hi_empunion` and the exogenous variable `age`. The available instruments are `ssiratio` and `multlc`. Use robust standard errors throughout (option `vce(robust)`).

(a) According to the OLS estimates, what is the impact of being insured through a union on the level of medical expenditures?

(b) Estimate the same model by instrumental variables, using just the single instrument `ssiratio`. Is the coefficient of `hi_empunion` plausible? Explain.

(c) Compare the efficiency of IV in part (b) with that of OLS.

(d) Estimate the model by 2SLS using both available instruments.

(e) Estimate the model by optimal GMM using both available instruments. Is there much difference in the estimates and estimate efficiency compared to 2SLS?

2. Consider the same example as question 1 using both instruments.

(a) Perform a Hausman test of endogeneity (after 2SLS estimation) using command `estat endogenous`. What do you conclude?

(b) Perform an overidentifying restrictions test (after optimal GMM estimation) using command `estat overid`. What do you conclude?

(c) Harder as not covered: Obtain weak instruments diagnostics (after 2SLS estimation) using command `estat firststage`. What do you conclude?

(d) Harder as not covered: Use Stata addon command `condivreg` (you need to load this up)

```
condivreg ldrugexp (hi_empunion = multlc ssiratio) age, lm ar 2sls test(0)
```

This gives three confidence intervals from 2SLS for the coefficient of `hi_empunion` that are felt to be more robust to weak instruments. How do these compare to the usual confidence intervals based on robust standard errors?

3. Consider IV estimation in the linear regression model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$, where \mathbf{X} is $N \times K$.

Premultiply by the $N \times m$ matrix of instruments \mathbf{Z} that satisfy $E[\mathbf{u}|\mathbf{Z}] = \mathbf{0}$ and $E[\mathbf{u}\mathbf{u}'|\mathbf{Z}] = \boldsymbol{\Omega}$ we have the model

$$\mathbf{Z}'\mathbf{y} = \mathbf{Z}'\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}'\mathbf{u}.$$

We assume that $m > K$, so the model is over-identified.

(a) Show that the $m \times 1$ “error” $\mathbf{Z}'\mathbf{u}$ has mean $\mathbf{0}$ and variance $\mathbf{Z}'\boldsymbol{\Omega}\mathbf{Z}$ (conditional on \mathbf{Z}).

(b) Show that GLS estimator of $\boldsymbol{\beta}$ in this model yields

$$\hat{\boldsymbol{\beta}} = [\mathbf{X}'\mathbf{Z}(\mathbf{Z}'\boldsymbol{\Omega}\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{X}]^{-1}\mathbf{X}'\mathbf{Z}(\mathbf{Z}'\boldsymbol{\Omega}\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{y}.$$

(c) Show that this is the 2SLS estimator if $\boldsymbol{\Omega} = \sigma^2\mathbf{I}$.

(d) When errors are heteroskedastic show that $\mathbf{Z}'\boldsymbol{\Omega}\mathbf{Z} = \sum_i \sigma_i^2 \mathbf{z}_i \mathbf{z}_i'$. Hence give the formula for this estimator if errors are heteroskedastic, and relate it to optimal GMM in this case.

(e) Show that the answer in part (b) simplifies to the IV estimator in the just-identified case. Hint: When $m = K$ the matrix $\mathbf{Z}'\mathbf{X}$ is square and invertible. Also $(\mathbf{ABC})^{-1} = \mathbf{C}^{-1}\mathbf{B}^{-1}\mathbf{A}^{-1}$.

4. Consider IV estimation in the linear regression model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$, where \mathbf{X} is $N \times K$, with instruments satisfying $\mathbf{E}[\mathbf{Z}'\mathbf{u}] = \mathbf{0}$ where \mathbf{Z} is $N \times m$.

(a) If $m = K$ find the formula of the estimator that solves $\mathbf{Z}'\mathbf{u} = \mathbf{0}$.

(b) If $m > K$ find the formula of the estimator that solves $(\mathbf{Z}'\mathbf{u})'\mathbf{W}(\mathbf{Z}'\mathbf{u})$.

[Hint: Use the chain rule of differentiation $\partial \mathbf{u}'\mathbf{C}\mathbf{u} / \partial \boldsymbol{\beta} = (\partial \mathbf{u} / \partial \boldsymbol{\beta})' \times \partial \mathbf{u}'\mathbf{C}\mathbf{u} / \partial \mathbf{u}$ and the result that for quadratic forms $\partial \mathbf{u}'\mathbf{C}\mathbf{u} / \partial \mathbf{u} = 2\mathbf{C}\mathbf{u}$.]

(c) What is the best choice of \mathbf{W} ?

5. This is a harder question to show consistency and asymptotic normality.

Consider a model with a scalar regressor: $y_i = \beta x_i + u_i$.

The OLS estimator is $\hat{\beta} = (\sum_{i=1}^N x_i^2)^{-1} \sum_{i=1}^N x_i u_i$.

The assumed d.g.p. is simple random sampling: (x_i, u_i) are i.i.d. with x_i i.i.d. with $\mathbf{E}[x_i] = \mu_x$ and $\mathbf{E}[x_i^2] = \mathbf{E}[x^2]$, and u_i i.i.d. with $\mathbf{E}[u_i] = 0$ and $\mathbf{V}[u_i] = \sigma_u^2$.

We use the following results.

Khinchine's Theorem: Let $\{X_i\}$ be i.i.d. (independent and identically distributed). If and only if $\mathbf{E}[X_i] = \mu$ exists, then $(\bar{X}_N - \mu) \xrightarrow{p} 0$.

Lindeberg-Levy CLT: Let $\{X_i\}$ be i.i.d. with $\mathbf{E}[X_i] = \mu$ and $\mathbf{V}[X_i] = \sigma^2$. Then $Z_N = \sqrt{N}(\bar{X}_N - \mu) / \sigma \xrightarrow{d} \mathcal{N}[0, 1]$.

(a) As $x_i u_i$ are i.i.d. apply Khinchine's Theorem to show that $N^{-1} \sum_i x_i u_i \xrightarrow{p} 0$.

(b) As x_i^2 are i.i.d. apply Khinchine's Theorem to show that $N^{-1} \sum_i x_i^2 \xrightarrow{p} \mathbf{E}[x^2]$.

(c) Combine these results to show that OLS estimator $\hat{\beta} \xrightarrow{p} \beta$.

(d) As $x_i u_i$ are i.i.d. apply Lindeberg-Levy CLT to $N^{-1} \sum_{i=1}^N x_i u_i$ to show that

$$\left(\frac{1}{\sqrt{N}} \sum_{i=1}^N x_i u_i \right) / \left(\sqrt{\sigma^2 \mathbf{E}[x^2]} \right) \xrightarrow{d} \mathcal{N}[0, 1].$$

(e) Hence show that $\frac{1}{\sqrt{N}} \sum_{i=1}^N x_i u_i \xrightarrow{d} \mathcal{N}[0, \sigma^2 \mathbf{E}[x^2]]$, using $a_N \times b_N \xrightarrow{d} a \times b$ if $a_N \xrightarrow{d} a$, $b_N \xrightarrow{p} b$.

(f) Combine (d) and (e) to show $\sqrt{N}(\hat{\beta} - \beta) \xrightarrow{d} \mathcal{N}[0, \sigma^2 (\mathbf{E}[x^2])^{-1}]$.