EC 425/525, Lab 8

Edward Rubin 24 May 2019

Prologue

Schedule

Last time

Helpful tips and tricks in ${\tt R}$

Today

Inference (in R)

Motivation

So far, we've focused on carefully **obtaining causal estimates** of the effect of some treatment D_i on our outcome Y_i .

Our discussion of research designs and their requirements/assumptions has centered on avoiding selection and securing unbiased and/or consistent estimates for τ .

In other words, we've concentrated on **point estimates**.

What about **inference**?

Shminference

Q Why care about inference?

A I'll give you two reasons.

- 1. We often want to **test theories/hypotheses**. Point estimates (i.e., $\hat{\beta}$) can't do this alone. Inference finishes the job.
- 1. Other times, we want to **measure the effect** of a treatment. Inference helps us think about the **precision** of our estimates.

Note: Similar reasoning can apply to bounding forecasting/predictions.

If you want answers, then you need to do inference correctly.

What's so complicated?

Angrist and Pischke told us that "correcting" our standard errors for heteroskedasticity may increase the standard errors up to 25%.

What else are we worried about?

What we're worried about

- Transformations of estimators, i.e., $\mathrm{Var}\Big[f\left(\hat{eta}\right)\Big]
 eq f\Big(\mathrm{Var}\Big[\hat{eta}\Big]\Big)$
- Dependence/correlation in our disturbance, i.e., $\mathrm{Cov}(\varepsilon_i,\,\varepsilon_j) \neq 0$
 - \circ Autocorrelation $arepsilon_t =
 ho arepsilon_{t-1} + arepsilon_t$
 - \circ Correlated shocks within groups $arepsilon_i = arepsilon_{g(i)} + arepsilon_i$
- Finite-sample properties vs. asymptotic properties
- Power and minimal detectable effects
- Multiple-hypothesis testing and p-hacking

In other words: We've got a lot to worry/think about.

Setup

Many studies—observational and experimental—have a treatment that is assigned to all/most individuals within a group.

- Classrooms/schools
- Households
- Villages/counties/states

Furthermore, we might imagine individuals within the same group may have correlated disturbances. For i and j in group g

$$\mathrm{Cov}ig(arepsilon_i,\,arepsilon_jig)=Eig[arepsilon_iarepsilon_jig]=
ho_arepsilon\sigma_arepsilon^2$$

where ρ_{ε} gives the within-group correlation of disturbances—what MHE calls the intraclass correlation coefficient.

Setup

In other words, we have a regression

$$y_i = eta_0 + eta_1 x_{g(i)} + arepsilon_i$$

where individual i is in group g, and $X_{q(i)}$ only varies across groups.

For within-group correlation, we can use an additive random-effects model

$$arepsilon_i =
u_{g(i)} + \eta_i$$

meaning group members all receive a common shock $\nu_{g(i)}$, and individuals receive independent shocks η_i .

Note We assume η_i is independent of η_j $(i \neq j)$ and ν_g $(\forall g)$.

Additive random effects

Based upon this model we've set up

$$arepsilon_i =
u_{g(i)} + \eta_i$$

the covariance between individuals i and j in group g is

$$egin{split} \operatorname{Cov}ig(arepsilon_i,\,arepsilon_jig) &= Eig[arepsilon_iarepsilon_jig] = Eig[ig(
u_g + \eta_iig)\,ig(
u_g + \eta_jig)ig] = Eig[
u_g^2ig] = \sigma_
u^2 \ &=
ho_arepsilon\sigma_arepsilon^2 \ &=
ho_arepsilon\left(\sigma_
u^2 + \sigma_\eta^2
ight) \end{split}$$

Thus, we can write the intraclass correlation coefficient as

$$ho_arepsilon = rac{\sigma_
u^2}{\sigma_arepsilon^2} = rac{\sigma_
u^2}{\sigma_
u^2 + \sigma_\eta^2}$$

What is ρ_{ε} ?

Let's review what we know.

$$arepsilon_i =
u_{g(i)} + \eta_i \qquad ext{and} \qquad
ho_arepsilon = rac{\sigma_
u^2}{\sigma_arepsilon^2} = rac{\sigma_
u^2}{\sigma_
u^2 + \sigma_\eta^2}$$

One way to think about ρ_{ε} is as the share of the variance of the disturbance ε_i accounted for by the shared disurbance $\nu_{g(i)}$.

As $u_{g(i)}$ accounts for more and more of the variation in ε_i , $\rho_{\varepsilon} \to 1$.

So...

Q Why do we care about ρ_{ε} ?

A It tells us by how wrong our standard errors can be if we treat all observations as independent.

Let ${
m Var}_o (\hat{eta}_1)$ denote the conventional variance formula for OLS estimator. †

Let $\operatorname{Var}(\hat{\beta}_1)$ denote the actual variance of $\hat{\beta}_1$.

[†] which treats all disturbances as independent (and identically distributed).

So....

With (1) nonstochasic regressors fixed by group and (2) groups of size n

$$rac{ ext{Var}ig(\hat{eta}_1ig)}{ ext{Var}_oig(\hat{eta}_1ig)} = 1 + (n-1)
ho_arepsilon \qquad \Longrightarrow \qquad rac{ ext{S.E.}ig(\hat{eta}_1ig)}{ ext{S.E.}_oig(\hat{eta}_1ig)} = \sqrt{1 + (n-1)
ho_arepsilon}$$

The term $\sqrt{1+(n-1)\rho_{\varepsilon}}$ is called the **Moulton factor**[†].

The **Moulton factor** tells us by what factor standard errors will be wrong if we ignore within-group correlation (conditional on assumptions **1** and **2**).

- Q What happens if $\rho = 1$? What if you duplicated your dataset?
- \mathbf{Q} What happens as n increases?

[†] After Moulton (1986).

The Moulton factor

The Moulton factor

$$rac{ ext{S.E.} \left(\hat{eta}_1
ight)}{ ext{S.E.}_o \left(\hat{eta}_1
ight)} = \sqrt{1 + (n-1)
ho_arepsilon}$$

shows even when ρ_{ε} is small, we can have vary large standard error issues.

Ex An experiment on 400 schools, each with 1,000 students.

If
$$ho_arepsilon=0.01$$
, the Moulton factor is $\sqrt{1+(1,000-1) imes0.01}pprox3.32$.

Test statistics

Recall
$$t_{ ext{stat}} = rac{\hat{eta}_1}{ ext{S.E.} \left(\hat{eta}_1
ight)}.$$

$$\therefore \frac{t_o}{t} = \frac{\hat{\beta}_1/\operatorname{S.E.}_o\left(\hat{\beta}_1\right)}{\hat{\beta}_1/\operatorname{S.E.}\left(\hat{\beta}_1\right)} = \frac{\operatorname{S.E.}\left(\hat{\beta}_1\right)}{\operatorname{S.E.}_o\left(\hat{\beta}_1\right)} = \text{the Moulton factor.}$$

Ex Thus, in our example of 400 schools with 1,000 students, ignoring within-school correlation of $\rho_{\varepsilon}=$ 0.01 would lead us test statistics that are more than 3 times as large as they should be.

This is why economics seminars have standard-error police.

Relaxing assumptions

If we allow regressors to vary by individual and groups to differ in size (n_g) ,

$$rac{ ext{Var}\Big(\hat{eta}_1\Big)}{ ext{Var}_o\Big(\hat{eta}_1\Big)} = 1 + \left[rac{ ext{Var}ig(n_gig)}{\overline{n}} + \overline{n} - 1
ight]
ho_x
ho_arepsilon$$

where ho_x denotes the intraclass (within-group) correlation of x_i . †

Important The Moulton factor for this general model depends upon the amount of within-group correlation in x_i and ε_i .

The special case is also important, as treatment is often fixed at some level.

[†] See MHF for mathematical definitions and the derivation.

The answer

Q So what do we do now?

A We've got options (as usual)

- 1. Parametrically model the random effects
- 2. Cluster-robust standard error (estimator)
- 3. Aggregate up to the group (or a similar method)
- 4. Block (group-based) bootstrap
- 5. GLS/MLE modeling y_i and ε_i

Most common: Cluster-robust standard errors

Runner up: Block bootstrap

Second runner up: Group-level analysis

Cluster-robust standard errors

Liang and Zeger (1986) extend White's heteroskedasticity-robust covariance matrix to allow for both clustering and heteroskedasticity.[†]

$$\hat{\Omega}_{ ext{cl}} = \left(ext{X}' ext{X}
ight)^{-1} \left(\sum_g ext{X}_g' \hat{\Psi}_g ext{X}_g
ight) \left(ext{X}' ext{X}
ight)^{-1}$$

$$\hat{\Psi}_g = ae_ge_g' = a egin{bmatrix} e_{1g}^2 & e_{1g}e_{2g} & \cdots & e_{1g}e_{n_gg} \ e_{1g}e_{2g} & e_{2g}^2 & e_{2g} \cdots & e_{2g}e_{n_gg} \ dots & dots & \ddots & dots \ e_{1g}e_{n_gg} & e_{2g}e_{n_gg} & \cdots & e_{n_gg}^2 \end{bmatrix}$$

such that e_g are the OLS residuals for group g, e_{ig} is the residual for individual i in group g, and a is a degrees-of-freedom adjustment.

[†] When people say clustering, they typically mean correlated disturbances within a group.

Cluster-robust standard errors

So now you know what lm_robust(), iv_robust(), etc. are doing when you specify a variable for clustering (e.g., clusters = var).

Time for a simulation.

Cluster simulation

Cluster simulation

The DGP

Let's opt for a simple-ish example.[†]

$$egin{aligned} y_{ig} &= 1 + 2x_g + arepsilon_{ig} \ arepsilon_{ig} &=
u_g + \eta_i \end{aligned}$$

where the $\eta_i \perp \eta_j$, $\eta_i \perp \nu_g$, and $\nu_g \perp \nu_h$.

Let's assume $\eta_i \sim N(0,1)$ and $u_g \sim N(0,1)$. And $x_g \sim N(0,1)$.

Plus n=100 with 10 groups.

[†] So we have more room for problem sets.

First we need to write the data generating process for one iteration.

```
# The DGP
sim dgp \leftarrow function(n = 100, n grps = 10, \sigma v = 1, \sigma \eta = 1) 
  # Create the right number of observations
  sample df \leftarrow expand.grid(i = 1:n, g = 1:n grps) %>% as tibble()
  # Create a unique ID (from 1 to number of observations)
  sample df \%% mutate(id = 1:(n * n grps))
  # Sample v at the group level
  # NOTE: Ungroup after grouping
  sample df %<>% group by(g) %>%
    mutate(v = rnorm(1, sd = \sigma v)) \% > \% ungroup()
  # Sample n at the individual level
  sample df \%% mutate(\eta = rnorm(n * n grps, sd = \sigma\eta))
  # Sample x g from N(0,1)
  sample_df %<>% group_by(g) %>% mutate(x = rnorm(1)) %>% ungroup()
  # Calculate v
  sample_df \%% mutate(y = 1 + 2 * x + v + \eta)
  # Return
  return(sample df)
```

Now we analyze.

```
# Analyze 'data'
sim analyze ← function(data) {
  # Conventional SEs
  se_ols ← lm_robust(y ~ x, data = data, se_type = "classical") %>%
   tidy() %>% extract2(2,3)
 # Cluster-robust SEs
  se cl \leftarrow lm robust(y \sim x, data = data, clusters = g) %>%
   tidy() %>% extract2(2,3)
  # Return a data frame of results
  data.frame(
    se = c(se ols, se cl),
   type = c("conventional", "clustered")
```

Now put them together with another function.

```
# Join sim_dgp and sim_analyze
sim_iter ← function(n = 100, n_grps = 10, σv = 1, ση = 1) {
    # Run the analysis in sim_analyze on the output of sim_dgp
    sim_dgp(n = 100, n_grps = 10, σv = 1, ση = 1) %>% sim_analyze()
}
```

And we run the simulation.

```
# Load and set up 'furrr'
p_load(furrr)
plan(multiprocess, workers = 8)
# Set a seed
set.seed(1234)
# Run the simulation 1,000 times
sim df \leftarrow future map dfr(
  # Repeat sample size 100 for 1000 times
 rep(100, 1000),
 # Our function
 sim_iter,
  # Let furrr know we want to set a seed
  .options = future_options(seed = T)
```

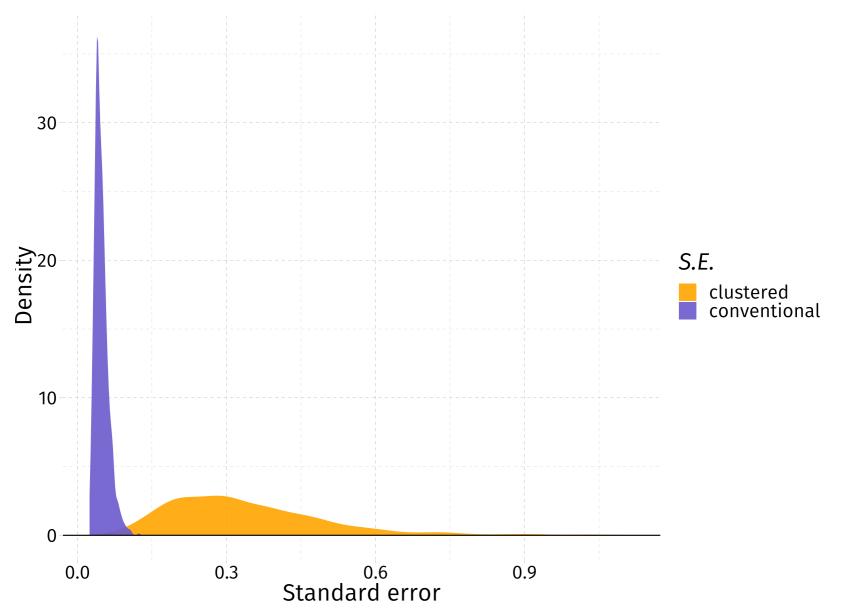


Table of contents

Inference

- 1. Motivation
- 2. Clustering
- 3. Moulton factors
 - Example
 - Test statistics
- 4. Answers
- 5. Cluster-robust S.E.s
- 6. Simulation