

ExpEcon Methods: Robust SEs, Clustering, Fixed & Random Effects

ECON 8877

P.J. Healy

First version thanks to Han Wang

Updated 2026-01-28

Overview

- HUGE problem: the “repeated measures problem”
- Example: You run 10 sessions.
Each session: 12 subjects, 20 periods, random rematching
- This is **NOT** 2,800 independent observations!!!
 - Econometrician: “that’s panel data”
- But at what level are there problems?
 - Subject effects: i is riskier than j
 - Session effects: this group was more cooperative
- How to deal with these effects in a regression??
 - Clustering? Fixed effects? Random effects?
 - At the session level? Individual level? Both??

Let's start by reviewing the asymptotic results for OLS.

Model

For $i = 1, 2, \dots, n$,

$$y_i = x_i' \beta + \epsilon_i$$

where y_i and ϵ_i are scalar, and x_i and β are $k \times 1$ column vectors.

Matrix version:

$$y = X\beta + \epsilon$$

Assumptions

1. Linear model: $y = X\beta + \epsilon$
2. X has full rank (so $X'X$ is invertible)
3. Non-stochastic X
4. $\mathbb{E}[\epsilon|X] = 0$
5. $\mathbb{E}[\epsilon_i^2|X] = \sigma^2 \forall i$ and $\mathbb{E}[\epsilon_i \epsilon_j|X] = 0$ (**homoskedasticity**)
6. Even stronger: $\epsilon|X \sim N(0, \sigma^2 \mathbf{I})$ (Normality)

OLS Estimate:

- $\hat{\beta} = (X'X)^{-1}(X'y) = (X'X)^{-1}(X'(X\beta + \epsilon)) = \beta + (X'X)^{-1}X'\epsilon$
- $\mathbb{E}[\hat{\beta}] = \beta$ since $\mathbb{E}[\epsilon|X] = 0$
- $\mathbb{V}[\hat{\beta}] = \mathbb{E}[(X'X)^{-1}X'\epsilon][\epsilon'X(X'X)^{-1}] = (X'X)^{-1}X'\mathbb{E}[\epsilon\epsilon']X(X'X)^{-1}$

Homoskedasticity: $\text{Cov}(\epsilon_i, \epsilon_j) = 0$, and $\text{Var}(\epsilon_i|x_i) = \sigma^2$.

- Under homoskedasticity, the middle term $\mathbb{E}[\epsilon\epsilon'] = \sigma^2\mathbf{I}$. This simplifies our variance:

$$\mathbb{V}(\hat{\beta})_{\text{homoskedasticity}} = \sigma^2(X'X)^{-1}$$

- Unbiased estimator of σ^2 :

$$\hat{\sigma}^2 = \frac{e'e}{n - k}$$

where $e = y - X\hat{\beta}$ are the residuals

Heteroskedasticity: $\mathbb{E}[\epsilon\epsilon'] = \sigma^2\Omega$, where σ^2 is a scalar and

$$\Omega = \begin{bmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_n^2 \end{bmatrix}$$

So now

$$\begin{aligned} \mathbb{V}[\hat{\beta}] &= (X'X)^{-1}X'\mathbb{E}[\epsilon\epsilon']X(X'X)^{-1} \\ &= (X'X)^{-1}X'\sigma^2\Omega X(X'X)^{-1} \end{aligned}$$

- For large samples our $\hat{\sigma}^2$ will actually be unbiased for σ^2
- White (1980): a better (robust) estimator for finite samples is

$$\hat{\mathbb{V}}[\hat{\beta}]_{HW} = (X'X)^{-1} \left(\sum_{i=1}^n e_i^2 x_i x_i' \right) (X'X)^{-1}$$

Doing this in practice

- But there are many options (Long & Ervin, 2000)

$$\hat{\mathbb{V}}(\hat{\beta})_{HW} = (X'X)^{-1} \left(\sum_i e_i^2 x_i x_i' \right) (X'X)^{-1} \quad (\text{HCo})$$

$$\hat{\mathbb{V}}(\hat{\beta})_{robust} = (X'X)^{-1} \left(\sum_i \frac{n}{n-k} e_i^2 x_i x_i' \right) (X'X)^{-1} \quad (\text{HC1})$$

$$\hat{\mathbb{V}}(\hat{\beta})_{HC2} = (X'X)^{-1} \left(\sum_i (1 - h_{ii})^{-1} e_i^2 x_i x_i' \right) (X'X)^{-1} \quad (\text{HC2})$$

$$\hat{\mathbb{V}}(\hat{\beta})_{HC3} = (X'X)^{-1} \left(\sum_i (1 - h_{ii})^{-2} e_i^2 x_i x_i' \right) (X'X)^{-1} \quad (\text{HC3})$$

where h_{ii} is the diagonal element of the “hat matrix” $(X(X'X)^{-1}X')$.

Note: HC1 is just a d.o.f. adjustment to HCo ($n/(n-k)$)

$\hat{\mathbb{V}}$ determines our confidence intervals. Thus, our size & power

Doing this in practice

- Let's say $A \leq B$ if $B - A$ is PSD (B is more conservative towards 0)

$$HCo \leq HC1 \leq HC2 \leq HC3$$

Woodridge ran simulations to show $HC2 - HC1$ is PSD ($n=200$, $k=3$, 1,000,000 replications, always true).

- Simulation studies show that $HC2$ and $HC3$ lead to better—with small n , possibly much better—confidence intervals than $HC1$.
- The Stata default with `vce(robust)` uses $HC1$.
- The R default with `sandwich` uses $HC3$. For R, see `estimateR`, `clubSandwich` and Kolesar's github repo.

Related: Young (2019)

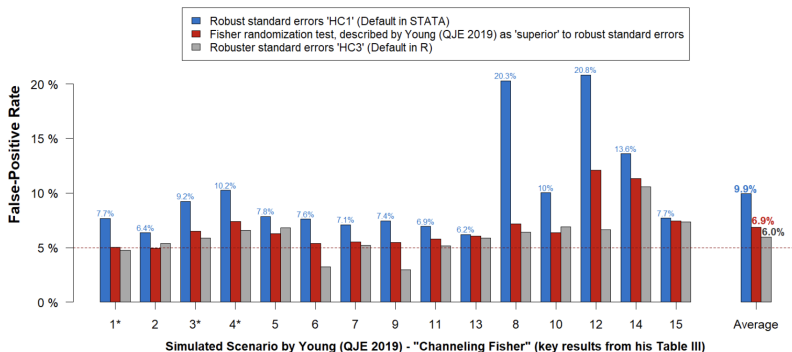
- A QJE paper (595 Google cites): use permutation tests instead
- look at 53 experimental papers from the journals of the AEA
- compare permutation tests to conventional tests
 - individual significance results: 13-22 percent fewer
 - joint significance results: 33-49 percent fewer

But what if we compare the different SE estimates?

Related: A post on Data Colada...

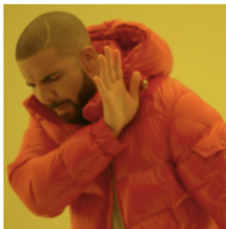
- The QJE study cited used HC1
- The datacolada post shows that using HC1 and HC3 can be very different when sample sizes are not large.
- But HC3 turns out to work quite well even with pretty small n.

False-Positive Rates with Randomization and Robuster Standard Errors are VERY Similar



Main Takeaway

For Stata users:



`reg y x, robust`



`reg y x, vce(hc3)`

Clustering and generalizing $\mathbb{E}[\epsilon\epsilon'|X]$

$$\Omega_{homoskedasticity} = \begin{bmatrix} \sigma^2 & & \\ & \ddots & \\ & & \sigma^2 \end{bmatrix}$$

$$\Omega_{heteroskedasticity} = \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_n^2 \end{bmatrix}$$

- We've ignored any correlation structure in Ω .
- In many cases, we don't have that. Instead, Ω has clusters.
 - units are people, and clusters are cities, states or countries
 - units are choices, and clusters are subjects, groups or sessions

Clustering and generalizing $\mathbb{E}[\epsilon\epsilon'|X]$

Let C_i denote unit i 's cluster assignment.

- A simple example:

$$\Omega_{ij} = \begin{cases} \sigma^2 & \text{if } i = j \\ \rho\sigma^2 & \text{if } C_i = C_j \text{ \& } i \neq j \\ 0 & \text{if } C_i \neq C_j \text{ \& } i \neq j \end{cases}$$

$$\Omega_{cluster} = \begin{bmatrix} \sigma^2 & \rho\sigma^2 & 0 & 0 \\ \rho\sigma^2 & \sigma^2 & 0 & 0 \\ & & \ddots & \\ 0 & 0 & \sigma^2 & \rho\sigma^2 \\ 0 & 0 & \rho\sigma^2 & \sigma^2 \end{bmatrix}$$

e.g., if we study individual choices, it might be ok to assume away the correlation between different subjects.

- A more unstructured example: $\Omega_{ij} = \sigma_{ij}$ if $C_i = C_j$.

Let the number of clusters be G , indexed by g

$$\hat{\mathbb{V}}(\hat{\beta})_{LZ} = (X'X)^{-1} \left(\sum_g X'_{g,n} e_{g,n} e'_{g,n} X_{g,n} \right) (X'X)^{-1}$$

(Liang & Zeger, 1986)

- This makes us think more generally, it's about getting the structure of Ω right. (So better to err on the conservative side)
- However, A recent QJE paper (Abadie, Athey, Imbens & Wooldridge, 2023) argues that this intuition is not correct.

Related: Abadie, Athey, Imbens & Wooldridge (2023)

A model of research design:

- There are m_k clusters (eg, states)
- Potential (unobserved) outcomes: $y_i(0)$ and $y_i(1)$ (control, trt)
- Goal: estimate $\tau := \frac{1}{n} \sum_{i=1}^n (y_i(1) - y_i(0))$
- True Avg Trt for cluster m : $\tau_m = \frac{1}{n_m} \sum_{i=1}^n \mathbf{1}_{\{m_i=m\}} (y_i(1) - y_i(0))$
- So $\tau = \sum_m \frac{n_m}{n} \tau_m$
- Sampling Process:
 1. Each cluster is sampled with probability $q \in (0, 1]$
 2. Each person is then sampled with probability $p \in (0, 1]$
- Treatment Assignment Process:
 1. Cluster- m assignment prob. $A_m \sim F$ with mean μ , var σ^2
 - Random assign.: $A_m = A_{m'} \forall m, m'$ (so $\sigma^2 = 0$)
 - Clustered assign: $A_m \in \{0, 1\}$ (so $\sigma^2 = \mu(1 - \mu)$)
 2. Each i in m treated iid with probability A_m
- Regress: $Y_i = \alpha + \beta W_i + \varepsilon_i$ (Y_i outcome, W_i treatment indicator)
- True standard error of $\hat{\beta}$ is v . How to estimate it?

Different cases:

1. Random sampling ($q = 1, p \leq 1$) and random assign ($\sigma^2 = 0$)
 - v = robust SE + correction factor
 - Correction vanishes if
 - 1.1 No heterogeneity in treatment effects across i , or
 - 1.2 Small sample ($p \approx 0$)
 - So, no need to cluster if you observe all clusters equally!
2. Clustered sampling ($q < 1, p \leq 1$) and rand assign ($\sigma^2 = 0$)
 - v = robust SE + correction factor + across-cluster variance
 - New term = 0 if avg trt effect is same across clusters
 - If so, then no need for cluster-robust SE
3. Clustered sampling plus clustered assignment ($\sigma^2 > 0$)
 - Two new terms added to v
 - Now cluster-robust SE is important

Abadie et al. provide an estimator \hat{v} , assumes you know the # of clusters.

A bootstrapped version is also available

Misconceptions:

- “The presence of within-cluster correlation implies the need for clustering.”
- “Being as conservative as necessary.”
 - Suppose we want to use the sample average to estimate the population mean. Suppose the population can be partitioned into clusters, e.g., in geographical units. If outcomes are positively correlated in clusters, the cluster variance will be larger than the robust variance.
But there's no need to cluster...
- “Researchers have only two choices: to cluster or not to cluster.”

Main Takeaways:

- “The decision on when and how to cluster standard errors depends on the nature of the sampling and the assignment processes only, not on the presence of within-cluster error components in the outcome variable.”
- The traditional advice of being as conservative as necessary is likely misguided.
- They suggest new ways to estimate variance: causal cluster variance (CCV) and two-stage cluster bootstrap (TSCB).
 - These are designed for applications with large number of observations and substantial variation in treatment assignment within clusters.
- Fixed effects do NOT remove need for clustering.

Doing this in practice

- There are ongoing debates on clustering...
- If we know the appropriate cluster level, we can implement this using the cluster command in Stata:

reg y x, cluster(g)

For experimentalists: Cluster by subject or by session?

- Frechette (2012): if session-level interactions are of interest (e.g., markets, or subjects make effectively one “choice” per session), then cluster by **session**
- Kim (2022): if subjects make repeated decisions in the same setting, cluster by **subject**. But make your sessions as similar as possible! Alternating order, etc.

Fixed Effects vs Random Effects

Suppose there are subject-specific effects:

$$y_{it} = x'_{it}\beta + u_i + e_{it}$$

where y_{it} , u_i and e_{it} are scalar, and x_{it} and β are $k \times 1$ column vectors.

Key Difference:

- Random effects: u_i is part of the error.
Need to assume no correlation between u_i and x_{it} .
So $\mathbb{E}[u_i | x_{i1}, \dots, x_{iT}] = \mathbb{E}[u_i]$
- Fixed effects: u_i is part of the intercept.
 u_i can be arbitrarily correlated with x_{it} .

Random effects

RE approach exploits the implied correlation structure of errors.

Let $v_{it} = u_i + e_{it}$. Stacking for T periods, we have $y_i = x_i\beta + v_i$. Define $\Omega = \mathbb{E}[v_i v_i' | x_i]$.

Assumptions

(RE 1) $\mathbb{E}[e_{it} | x_{i1}, \dots, x_{iT}] = 0$ and $\mathbb{E}[u_i | x_{i1}, \dots, x_{iT}] = 0$.

(RE 2) $\mathbb{E}[e_i e_i' | x_i, u_i] = \sigma_e^2 I_T$ and $\mathbb{E}[u_i^2 | x_i] = \sigma_u^2$

$$\Omega = \begin{bmatrix} \sigma_u^2 + \sigma_e^2 & \sigma_u^2 & \cdots & \sigma_u^2 \\ \sigma_u^2 & \sigma_u^2 + \sigma_e^2 & & \sigma_u^2 \\ \sigma_u^2 & & \ddots & \sigma_u^2 \\ \sigma_u^2 & \cdots & \sigma_u^2 & \sigma_u^2 + \sigma_e^2 \end{bmatrix}$$

$$\bullet \hat{\beta}_{RE} = (\sum_i x_i' \hat{\Omega}^{-1} x_i)^{-1} (\sum_i x_i' \hat{\Omega}^{-1} y_i).$$

Feasible GLS estimation of RE model

Step 1. Run a pooled OLS of y_{it} on x_{it} and get the residuals \hat{v}_{it} .

Step 2. Estimate $\sigma_v^2 = \sigma_u^2 + \sigma_e^2$ by $\hat{\sigma}_v^2 = \frac{1}{nT-k} \sum_i \sum_t \hat{v}_{it}^2$.

Step 3. Estimate σ_u^2 using cross terms only:

$$\hat{\sigma}_u^2 = \frac{1}{nT(T-1)/2-k} \sum_{i=1}^n \sum_{t=1}^{T-1} \sum_{s=t+1}^T \hat{v}_{it} \hat{v}_{is}.$$

Step 4. Form $\hat{\Omega}$ using $\hat{\sigma}_v^2$ and $\hat{\sigma}_u^2$.

Step 5. Estimate β by GLS: $\hat{\beta}_{RE} = (\sum_i x_i' \hat{\Omega}^{-1} x_i)^{-1} (\sum_i x_i' \hat{\Omega}^{-1} y_i)$

Assumptions

$$(\text{FE 1}) \mathbb{E}[e_{it} | x_{i1}, \dots, x_{iT}] = 0.$$

$$(\text{FE 2}) \mathbb{E}[e_i e_i' | x_i, u_i] = \sigma_e^2 I_T.$$

There are several derivations of the estimator.

- Add individual specific dummies: $y = X\beta + Du + e$. Then OLS estimation of β proceeds by the Frisch–Waugh–Lovell theorem. Define $y^* = y - D(D'D)^{-1}D'y$ and $X^* = X - D(D'D)^{-1}D'X$.

$$\hat{\beta}_{FE} = (X^{*'}X^*)^{-1}X^{*'}y^*$$

- De-mean/differencing: $\hat{\beta}_{FE} = \hat{\beta}_{within}$

Doing this in practice

- We can use the Hausman test to choose RE vs FE. (H_0 is in favor of “random effects”)
- In stata, RE or FE estimation:

xtset

xtreg y x, re

xtreg y x, fe

Note that the default panel structure in Stata has two dimensions (individual i and time t). There are packages for higher dimensions, e.g. in Changkuk's MPL paper, he has “individual”, “product” and “round”.

- Estimating FE using dummies is very flexible when we want to control different levels of fixed effects. But the # of regressors can be very large.

RE vs FE in experimental economics:

- Randomly assigning subjects to treatments $\Rightarrow u_i$ is uncorrelated with treatment (x_i)
- Thus, random effects are appropriate
 - Unless correlated with some other x_i column!
- Random effects estimators are more efficient than fixed effects

Merrett (2012): Cross-validation to test various methods. RE wins.

Final thoughts

- Many ways of estimating variance: analytical/ bootstrap
- With iid data, if we worry about heteroskedasticity, there are HCO, HC1 (HW), HC2, HC3... When sample size is small, we'd better use HC2 or HC3.
- With data that is not iid, clustering can adjust the variance. We need to motivate why and how to cluster.
- Random effects or fixed effects are on the model level. It's helpful, e.g., when we want to control some individual-specific effects.
- Individual-specific effects are treated as part of the error in RE models, while as part of the intercept in FE models.

References

- *Econometrics* textbook by Bruce Hansen
- “Yale Applied Empirical Methods PhD Courses” by Paul Goldsmith-Pinkham [link]
- “[99] Hying Fisher: The Most Cited 2019 QJE Paper Relied on an Outdated Stata Default to Conclude Regression p -values Are Inadequate” by Uri Simonsohn, on Data Colada [link]
- Jeffery Wooldridge’s comments on the Data Colada post [link]
- Alwyn Young, Channeling Fisher: Randomization Tests and the Statistical Insignificance of Seemingly Significant Experimental Results, *The Quarterly Journal of Economics*, Volume 134, Issue 2, May 2019, Pages 557–598, <https://doi.org/10.1093/qje/qjyo29>
- Alberto Abadie and others, When Should You Adjust Standard Errors for Clustering?, *The Quarterly Journal of Economics*, Volume 138, Issue 1, February 2023, Pages 1–35, <https://doi.org/10.1093/qje/qjaco38>