

Topic 5: Two-Way Cluster-Robust (TWCR) Standard Errors

by Sai Zhang

Key points: The validity of Two-Way Cluster-Robust (TWCR) standard errors

Disclaimer: This note is compiled by Sai Zhang.

5.1 One-Way Clustering

First, consider the case of one-way clustering. The linear model with one-way clustering

$$y_{ig} = \mathbf{x}_{ig}\boldsymbol{\beta} + u_{ig}$$

where i denotes the i th of the N individuals in the sample, j denotes the g th of the G clusters, assume that

- $\mathbb{E}[u_{ig} | \mathbf{x}_{ig}] = 0$
- error independence across clusters: for $i \neq j$

$$\mathbb{E}[u_{ig}u_{jg'} | \mathbf{x}_{ig}, \mathbf{x}_{jg'}] = 0 \quad (5.1)$$

unless $g = g'$, that is, errors for individuals within the same cluster may be correlated.

Grouping observations by cluster, get

$$\mathbf{y}_g = \mathbf{X}_g\boldsymbol{\beta} + \mathbf{u}$$

where \mathbf{X}_g has dimension $N_g \times K$ and \mathbf{y}_g has dimension $N_g \times 1$, with N_g observations in cluster g . Stacking over cluster, get the matrix form of the model

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

with \mathbf{y}, \mathbf{u} being $N \times 1$ vectors, \mathbf{X} being an $N \times K$ matrix. OLS estimator gives

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} = \left(\sum_{g=1}^G \mathbf{X}_g' \mathbf{X}_g \right)^{-1} \sum_{g=1}^G \mathbf{X}_g' \mathbf{y}_g \quad (5.2)$$

then, by CLT, we have that $\sqrt{G}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{d} \mathcal{N}(0, \boldsymbol{\Sigma})$ where the variance matrix of the limit normal distribution $\boldsymbol{\Sigma}$ is

$$\left(\lim_{G \rightarrow \infty} \frac{1}{G} \sum_{g=1}^G \mathbb{E}[\mathbf{X}_g' \mathbf{X}_g] \right)^{-1} \left(\lim_{G \rightarrow \infty} \frac{1}{G} \sum_{g=1}^G \mathbb{E}[\mathbf{X}_g' \mathbf{u}_g' \mathbf{u}_g \mathbf{X}_g] \right) \times \left(\lim_{G \rightarrow \infty} \frac{1}{G} \sum_{g=1}^G \mathbb{E}[\mathbf{X}_g' \mathbf{X}_g] \right)^{-1} \quad (5.3)$$

If the primary source of clustering is due to group-level common shocks, a useful approximation is that for the j th regressor, the default OLS variance estimate based on $s^2(\mathbf{X}'\mathbf{X})^{-1}$ should be inflated by $\tau_j \approx 1 + \rho_{x_j}\rho_u(\bar{N}_g - 1)$, where

- s is the estimated standard deviation of the error

- ρ_{x_j} is a measure of within-cluster correlation of x_j
- ρ_u is the within-cluster error correlation
- \bar{N}_g is the average cluster size

It's easy to see the τ_j can be large even with small ρ_u (Kloek, 1981; Scott and Holt, 1982; Moulton, 1990). If assume the model for the cluster error variance matrices $\mathbf{\Omega}_g = \mathbb{V}[\mathbf{u}_g | \mathbf{X}_g] = \mathbb{E}[\mathbf{u}_g \mathbf{u}_g' | \mathbf{X}_g]$, and there is a consistent estimate $\hat{\mathbf{\Omega}}_g$ of $\mathbf{\Omega}_g$, we can estimate $\mathbb{E}[\mathbf{X}_g' \mathbf{u}_g \mathbf{u}_g' \mathbf{X}_g] = \mathbb{E}[\mathbf{X}_g' \mathbf{\Omega}_g \mathbf{X}_g]$ via GLS.

Cluster-robust variance matrix estimate consider

$$\hat{\mathbb{V}}[\hat{\beta}] = (\mathbf{X}'\mathbf{X})^{-1} \left(\sum_{g=1}^G \mathbf{X}_g' \hat{\mathbf{u}}_g \hat{\mathbf{u}}_g' \mathbf{X}_g \right) (\mathbf{X}'\mathbf{X})^{-1} \quad (5.4)$$

where $\hat{\mathbf{u}}_g = \mathbf{y}_g - \mathbf{X}_g \hat{\beta}$. This estimate is consistent if

$$G^{-1} \sum_{g=1}^G \mathbf{X}_g' \hat{\mathbf{u}}_g \hat{\mathbf{u}}_g' \mathbf{X}_g - G^{-1} \sum_{g=1}^G \mathbb{E}[\mathbf{X}_g' \mathbf{u}_g \mathbf{u}_g' \mathbf{X}_g] \xrightarrow{P} \mathbf{0}$$

as $G \rightarrow \infty$. An informal presentation of Eq.(5.4) is to rewrite the central matrix as

$$\hat{\mathbf{B}} = \sum_{g=1}^G \mathbf{X}_g' \hat{\mathbf{u}}_g \hat{\mathbf{u}}_g' \mathbf{X}_g = \mathbf{X}' \begin{bmatrix} \hat{\mathbf{u}}_1 \hat{\mathbf{u}}_1' & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \hat{\mathbf{u}}_2 \hat{\mathbf{u}}_2' & & \vdots \\ \vdots & & \ddots & \mathbf{0} \\ \mathbf{0} & \cdots & & \hat{\mathbf{u}}_G \hat{\mathbf{u}}_G' \end{bmatrix} \mathbf{X} = \mathbf{X}' (\hat{\mathbf{u}} \hat{\mathbf{u}}' \otimes \mathbf{S}^G) \mathbf{X} \quad (5.5)$$

where \otimes denotes element-wise multiplication. The (p, q) th element of this matrix is

$$\sum_{i=1}^N \sum_{j=1}^N x_{ia} x_{jb} \hat{u}_i \hat{u}_j \cdot \mathbf{1}(i, j \text{ in the same cluster})$$

with $\hat{u}_i = y_i - \mathbf{x}_i' \hat{\beta}$.

\mathbf{S}^G is an $N \times N$ indicator matrix with $\mathbf{S}_{ij}^G = 1$ only if the i th and j th observation belong to the same cluster: it zeros out a large amount of $\hat{\mathbf{u}} \hat{\mathbf{u}}'$ (asymptotically equivalently, $\mathbf{u} \mathbf{u}'$), specifically, only $\sum_{g=1}^G N_g^2$ out of $N^2 = \left(\sum_{g=1}^G N_g \right)^2$ terms are not zero (sub-matrices on the diagonal). Asymptotically

- for fixed N_g , $\frac{1}{N^2} \sum_{g=1}^G N_g^2 \xrightarrow{G \rightarrow \infty} 0$
- for balanced clusters $N_g = N/G$, $\frac{1}{N^2} \sum_{g=1}^G N_g^2 = \frac{1}{G} \xrightarrow{G \rightarrow \infty} 0$

A strand of literature popularizes this method:

- Liang and Zeger (1986): in a generalized estimatin equations setting
- Arellano (1987): fixed effects estimator in linear panel models
- Hansen (2007): asymptotic theory for panel data where $T \rightarrow \infty$ in addition to $N \rightarrow \infty$ (or $N_g \rightarrow \infty$ in addition to $G \rightarrow \infty$ in the notation above).

5.2 Two-Way Clustering

Now, consider the case of two-way clustering,

$$y_{i,gh} = \mathbf{x}'_{i,gh} \boldsymbol{\beta} + u$$

where each observation may belong to **two** dimension of groups: group $g \in \{1, \dots, G\}$ and $h \in \{1, \dots, H\}$, and for $i \neq j$

$$\mathbb{E} [u_{i,gh} u_{j,g'h'} \mid \mathbf{x}_{i,gh}, \mathbf{j}, \mathbf{g}'\mathbf{h}'] = 0 \quad (5.6)$$

unless $g = g'$ or $h = h'$, that is, errors for individuals within the same group (along either g or h) may be correlated.

Cluster-robust variance matrix estimate extending the one-way clustering case, keep elements of $\hat{\mathbf{u}}\hat{\mathbf{u}}'$ where the i th and j th observations share a cluster in **any** dimension, then similar to Eq.(5.5)

$$\hat{\mathbf{B}} = \mathbf{X}' \left(\hat{\mathbf{u}}\hat{\mathbf{u}}' \otimes \mathbf{S}^{GH} \right) \mathbf{X} \quad (5.7)$$

here \mathbf{S}^{GH} is an $N \times N$ indicator matrix with $S_{ij}^{GH} = 1$ only if the i th and j th observation share any cluster, the (p, q) th element of this matrix is

$$\sum_{i=1}^N \sum_{j=1}^N x_{ia} x_{jb} \hat{u}_i \hat{u}_j \cdot \mathbf{1}(i, j \text{ share any cluster})$$

$\hat{\mathbf{B}}$ can also be presented in one-way cluster-robust fashion:

$$\hat{\mathbf{B}} = \mathbf{X}' \left(\hat{\mathbf{u}}\hat{\mathbf{u}}' \otimes \mathbf{S}^{GH} \right) \mathbf{X} = \mathbf{X}' \left(\hat{\mathbf{u}}\hat{\mathbf{u}}' \otimes \mathbf{S}^G \right) \mathbf{X} + \mathbf{X}' \left(\hat{\mathbf{u}}\hat{\mathbf{u}}' \otimes \mathbf{S}^H \right) \mathbf{X} - \mathbf{X}' \left(\hat{\mathbf{u}}\hat{\mathbf{u}}' \otimes \mathbf{S}^{G \cap H} \right) \mathbf{X} \quad (5.8)$$

where $\mathbf{G}^{GH} = \mathbf{G}^G + \mathbf{G}^H - \mathbf{G}^{G \cap H}$, with

- \mathbf{G}^G : $G_{ij}^G = 1$ only if the i th and j th observation belong to the same cluster $g \in \{1, 2, \dots, G\}$
- \mathbf{G}^H : $G_{ij}^H = 1$ only if the i th and j th observation belong to the same cluster $h \in \{1, 2, \dots, H\}$
- $\mathbf{G}^{G \cap H}$: $G_{ij}^{G \cap H} = 1$ only if the i th and j th observation belong to **both** the same cluster $g \in \{1, 2, \dots, G\}$ and the same cluster $h \in \{1, 2, \dots, H\}$

then, similar to one-way clustering case,

$$\begin{aligned} \hat{\mathbf{V}}[\hat{\boldsymbol{\beta}}] &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \left(\hat{\mathbf{u}}\hat{\mathbf{u}}' \otimes \mathbf{S}^G \right) \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \\ &\quad + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \left(\hat{\mathbf{u}}\hat{\mathbf{u}}' \otimes \mathbf{S}^H \right) \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \\ &\quad - (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \left(\hat{\mathbf{u}}\hat{\mathbf{u}}' \otimes \mathbf{S}^{G \cap H} \right) \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \end{aligned} \quad (5.9)$$

that is,

$$\hat{\mathbf{V}}[\hat{\boldsymbol{\beta}}] = \hat{\mathbf{V}}^G[\hat{\boldsymbol{\beta}}] + \hat{\mathbf{V}}^H[\hat{\boldsymbol{\beta}}] - \hat{\mathbf{V}}^{G \cap H}[\hat{\boldsymbol{\beta}}] \quad (5.10)$$

each of Eq.(5.10) can be separately computed by OLS of \mathbf{y} on \mathbf{X} , with variance matrix estimates $\hat{\mathbf{V}}$ based on

- clustering on $g \in \{1, 2, \dots, G\}$
- clustering on $h \in \{1, 2, \dots, H\}$
- clustering on $(g, h) \in \{(1, 1), \dots, (G, H)\}$

Practical considerations It is required to know what *ways* will be potentially important for clustering, which can be tested via checking the dimension of correlations in the errors. There are several ways to test

- estimate sample covariances of $\mathbf{X}'\hat{\mathbf{u}}$ within dimensions, test the null that the **average** of such covariances is 0: rejecting this null is sufficient (not necessary) to reject the null of no clustering (White, 1980)
- for **small samples**, Eq. (5.4) is biased downwards. This is corrected (in Stata) by replacing $\hat{\mathbf{u}}_g$ with $\sqrt{c}\hat{\mathbf{u}}_g$, where $c = \frac{G}{G-1} \frac{N-1}{N-K} \simeq \frac{G}{G-1}$. For two-way clustering (Eq. 5.8), there are 2 ways of correction:
 - choose correction terms for each of the 3 components:

$$c_1 = \frac{G}{G-1} \frac{N-1}{N-K}, c_2 = \frac{H}{H-1} \frac{N-1}{N-K}, c_3 = \frac{I}{I-1} \frac{N-1}{N-K}$$

with I being the number of unique clusters determined by $G \cap H$

- choose a constant terms for all components:

$$c = \frac{J}{J-1} \frac{N-1}{N-K}$$

with $J = \min(G, H)$

- **Var-cov matrix not positive-semidefinite**: $\hat{\mathbf{V}}[\hat{\boldsymbol{\beta}}]$ might have negative elements on the diagonal (Eq. 5.10), informly, this is more likely to arise when clustering is done over the same groups as the fixed effects. One way to address this issue is using *eigendecomposition* technique:

$$\hat{\mathbf{V}}[\hat{\boldsymbol{\beta}}] = \mathbf{U}\mathbf{\Lambda}\mathbf{U}'$$

where

- \mathbf{U} containing the eigenvectors of $\hat{\mathbf{V}}$
- $\mathbf{\Lambda} = \text{diag}[\lambda_1, \dots, \lambda_d]$ contains the eigenvalues of $\hat{\mathbf{V}}$

then create $\mathbf{\Lambda}^+ = \text{diag}[\lambda_1^+, \dots, \lambda_d^+]$ with $\lambda_j^+ = \max(0, \lambda_j)$ and use $\hat{\mathbf{V}}^+[\hat{\boldsymbol{\beta}}] = \mathbf{U}\mathbf{\Lambda}^+\mathbf{U}'$ as the estimate

5.3 Multiway Clustering

Cameron et al. (2011) extended the framework¹ to allow clustering in D dimensions, then we can do the following reframing

- G_d : the number of clusters in dimension $d \in \{1, 2, \dots, D\}$
- D -vector $\boldsymbol{\delta}_i = \delta(i)$, with function $\delta : \{1, 2, \dots, N\} \rightarrow \times_{d=1}^D \{1, 2, \dots, G_d\}$ lists the cluster membership in each dimension of each observation

then we have

$$\mathbf{1}[i, j \text{ shares a cluster}] = 1 \Leftrightarrow \delta_{id} = \delta_{jd}$$

for some $d \in \{1, 2, \dots, D\}$, where δ_{id} denotes the d th element of $\boldsymbol{\delta}_i$. Also

- D -vector \mathbf{r} : define the set

$$R \equiv \{\mathbf{r} : r_d \in \{0, 1\}, d = 1, 2, \dots, D, \mathbf{r} \neq \mathbf{0}\}$$

elements of the set R can be used to index all cases where 2 observations share a cluster in at least one dimension. Define the function

$$\mathbf{I}_r(i, j) \equiv \mathbf{1}[r_d \delta_{id} = r_d \delta_{jd}, \forall d]$$

¹Also proposed by Thompson (2011).

which indicates whether observations i and j have identical cluster membership for **all** dimensions d s.t. $r_d = 1$. Then we have a *aggregate* identifier

$$\mathbf{I}(i, j) = 1 \Leftrightarrow \mathbf{I}_r(i, j) = 1 \text{ for some } \mathbf{r} \in R$$

i.e., 2 observations share **at least** one dimension.

The define the $2^D - 1$ matrices

$$\tilde{\mathbf{B}}_r \equiv \sum_{i=1}^N \sum_{j=1}^N \mathbf{x}_i \mathbf{x}_j' \hat{u}_i \hat{u}_j \mathbf{I}_r(i, j) \quad (5.11)$$

with $\mathbf{r} \in R$.

Var-cov matrix estimator consider, similarly, an estimator

$$\hat{\mathbb{V}}[\hat{\beta}] = (\mathbf{X}'\mathbf{X})^{-1} \tilde{\mathbf{B}} (\mathbf{X}'\mathbf{X})^{-1} \equiv (\mathbf{X}'\mathbf{X})^{-1} \left(\sum_{\|\mathbf{r}\|=k, \mathbf{r} \in R} (-1)^{k+1} \tilde{\mathbf{B}}_r \right) (\mathbf{X}'\mathbf{X})^{-1} \quad (5.12)$$

where cases of clustering on an odd number of dimensions are added, those of clustering on an even number of dimensions are subtracted. Consider the case of $D = 3$,

$$(\tilde{\mathbf{B}}_{(1,0,0)} + \tilde{\mathbf{B}}_{(0,1,0)} + \tilde{\mathbf{B}}_{(0,0,1)}) - (\tilde{\mathbf{B}}_{(1,1,0)} + \tilde{\mathbf{B}}_{(1,0,1)} + \tilde{\mathbf{B}}_{(0,1,1)}) + \tilde{\mathbf{B}}_{(1,1,1)}$$

$\tilde{\mathbf{B}}$ is identical to $\hat{\mathbf{B}}$ defined analogically as in Eq.(5.8), since

- no observation pair with $\mathbf{I}(i, j) = 0$: this is immediate, since $\mathbf{I}(i, j) = 0 \Leftrightarrow \mathbf{I}_r(i, j) = 0, \forall \mathbf{r}$
- the covariance term corresponding to each observation pair with $\mathbf{I}(i, j) = 1$ is included **exactly once** in $\tilde{\mathbf{B}}$: by inclusion-exclusion principle for set cardinality

$$\mathbf{I}(i, j) \Rightarrow \sum_{\|\mathbf{r}\|=k, \mathbf{r} \in R} (-1)^{k+1} \mathbf{I}_r(i, j) = 1$$

Curse of dimensionality this could arise in a setting with **many dimensions** of clustering, and in which one or more dimensions have **few** clusters². **Cameron et al. (2011)** suggested an ad-hoc rule of thumb for approximating sufficient numbers of clusters.

5.3.1 Non-linear Estimators

m-Estimators Consider an m -estimator that solves

$$\sum_{i=1}^N \mathbf{h}_i(\hat{\theta}) = \mathbf{0}$$

under standard assumptions, $\hat{\theta}$ is asymptotically normal with estimated variance matrix

$$\hat{\mathbb{V}}[\hat{\theta}] = \hat{\mathbf{A}}^{-1} \hat{\mathbf{B}} \hat{\mathbf{A}}'^{-1} \quad (5.13)$$

where $\hat{\mathbf{A}} = \sum_i \frac{\partial \mathbf{h}_i}{\partial \theta'} \Big|_{\hat{\theta}}$ and $\hat{\mathbf{B}}$ is an estimate of $\mathbb{V}[\sum_i \mathbf{h}_i]$.

²The square design (each dimension has the same number of clusters) with orthogonal dimensions has the **least** independence of observations.

- **one-way clustering** $\hat{\mathbf{B}} = \sum_{g=1}^G \hat{\mathbf{h}}_g \hat{\mathbf{h}}_g'$ where $\hat{\mathbf{h}}_g = \sum_{i=1}^{N_g} \hat{\mathbf{h}}_{ig}$, clustering may not lead to parameter inconsistency, depending on whether $\mathbb{E}[\mathbf{h}_i(\boldsymbol{\theta})] = \mathbf{0}$ with clustering
 - **population-averaged approach**: assume $\mathbb{E}[y_{ig} | \mathbf{x}_{ig}] = \Phi(\mathbf{x}_{ig}'\boldsymbol{\beta})$
 - **random effects approach**: let $y_{ig} = 1$ if $y_{ig}^* > 0$ where $y_{ig}^* = \mathbf{x}_{ig}'\boldsymbol{\beta} + \epsilon_g + \epsilon_{ig}$, where
 - * idiosyncratic error $\epsilon_{ig} \sim \mathcal{N}(0, 1)$
 - * cluster-specific error $\epsilon_g \sim \mathcal{N}(0, \sigma_g^2)$
 then we have the alternative moment condition

$$\mathbb{E}[y_{ig} | \mathbf{x}_{ig}] = \Phi\left(\frac{\mathbf{x}_{ig}'\boldsymbol{\beta}}{\sqrt{1 + \sigma_g^2}}\right)$$

- **multiway clustering** replacing $\hat{\mathbf{u}}_i \mathbf{x}_i$ in Eq.(5.11) with $\hat{\mathbf{h}}_i$, then we have

$$\hat{\mathbb{V}}[\hat{\boldsymbol{\theta}}] = \hat{\mathbf{A}}^{-1} \hat{\mathbf{B}} \hat{\mathbf{A}}^{-1}$$

where

$$\hat{\mathbf{A}} = \sum_i \frac{\partial \mathbf{h}_i}{\partial \boldsymbol{\theta}'} \bigg|_{\hat{\boldsymbol{\theta}}} \quad \hat{\mathbf{B}} = \sum_{\|\mathbf{r}\|=k, \mathbf{r} \in R} (-1)^{k+1} \tilde{\mathbf{B}}_{\mathbf{r}} \quad \tilde{\mathbf{B}}_{\mathbf{r}} \equiv \sum_{i=1}^N \sum_{j=1}^N \hat{\mathbf{h}}_i \hat{\mathbf{h}}_j' \mathbb{I}_{\mathbf{r}}(i, j)$$

with $\mathbf{r} \in R^3$.

GMM estimation Consider an example of over-identified models: linear two stage least squares with more instruments than endogenous regressors, we have

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}) = \arg \min_{\boldsymbol{\theta}} \left(\sum_{i=1}^N \mathbf{h}_i(\boldsymbol{\theta}) \right)' \mathbf{W} \left(\sum_{i=1}^N \mathbf{h}_i(\boldsymbol{\theta}) \right)$$

where \mathbf{W} is a symmetric positive definite weighting matrix. Under standard regularity conditions, $\hat{\boldsymbol{\theta}}$ is asymptotically normal, with estimated variance matrix

$$\hat{\mathbb{V}}[\hat{\boldsymbol{\theta}}] = (\hat{\mathbf{A}}' \mathbf{W} \hat{\mathbf{A}})^{-1} \hat{\mathbf{A}}' \mathbf{W} \hat{\mathbf{B}} \mathbf{W} \hat{\mathbf{A}} (\hat{\mathbf{A}}' \mathbf{W} \hat{\mathbf{A}})^{-1}$$

again, $\hat{\mathbf{A}} = \sum_i \frac{\partial \mathbf{h}_i}{\partial \boldsymbol{\theta}'} \bigg|_{\hat{\boldsymbol{\theta}}}$, and $\hat{\mathbf{B}}$ is an estimate of $\mathbb{V}[\sum_i \mathbf{h}_i]$.

5.4 Menzel (2021): Asymptotic Gaussianity

One key of TWCR inference is the asymptotic Gaussianity, [Menzel \(2021\)](#) pointed out the potential non-Gaussianity of the limit distribution. Still, consider a random array (Y_{it}) indexed by two dimensions by $i = 1, \dots, N$ and $t = 1, \dots, T$. Clusters are sampled independently at random from an infinite population, but otherwise **unrestricted** in dependence within each row $\mathbf{Y}_i := (Y_{i1} \dots, Y_{iT})$ and within each column $\mathbf{Y}_{\cdot t} := (Y_{1t}, \dots, Y_{Nt})$.

³This multiway clustering can be implemented using several one-way clustered bootstraps. Each of the one-way cluster robust matrices is estimated by a pairs cluster bootstrap that resamples with replacement from the appropriate cluster dimension. They are then combined as if they had been estimated analytically ([Cameron et al., 2011](#)).

5.4.1 Distribution of Sample Average

First, consider

$$\bar{Y}_{NT} := \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T Y_{it}$$

and approximate the asymptotic distribution regardless of whether, or what type of, cluster-dependence is present.

3 scenarios of the array (Y_{it})

- **no cluster-dependence**: (Y_{it}) are mutually independent, CLT at a rate of $(NT)^{-1/2}$ applies (under regularity conditions)
- **correlation within clusters**: the convergence rate of (Y_{it}) is determined by the number of relevant clusters
- **non-separable models of heterogeneity (dependence with clusters, even uncorrelated)**⁴: The asymptotic behavior is non-standard

Consider 2 examples:

- **Additive factor model**

$$Y_{it} = \mu + \alpha_i + \gamma_t + \epsilon_{it}$$

where μ is a constant, and $\alpha_i, \gamma_t, \epsilon_{it}$ are zero-mean i.i.d. random variables for $i = 1, \dots, N$ and $t = 1, \dots, T$ with bounded second moments, and $N = T$. Based on a standard central limit theory, we have

- in the non-degenerate case with $\text{Var}(\alpha_i) > 0$ or $\gamma_t > 0$, the sample distribution

$$\sqrt{N} \left(\bar{Y}_{NT} - \mathbb{E}[Y_{it}] \right) \xrightarrow{d} \mathcal{N}(0, \text{Var}(\alpha_i) + \text{Var}(\gamma_t))$$

- in the degenerate case of no clustering with $\text{Var}(\alpha_i) = \text{Var}(\gamma_t) = 0$, the sample distribution

$$\sqrt{NT} \left(\bar{Y}_{NT} - \mathbb{E}[Y_{it}] \right) \xrightarrow{d} \mathcal{N}(0, \text{Var}(\epsilon_{it}))$$

if marginal distributions of $\alpha_i, \gamma_t, \epsilon_{it}$ are known, we can simulate from the joint distribution of (Y_{it}) by sampling the individual components at random, a bootstrap procedure would be consistent. If **unknown**, consider estimators

$$\begin{aligned} \hat{\alpha}_i &:= \frac{1}{T} \sum_{t=1}^T (Y_{it} - \bar{Y}_{NT}) = \alpha_i + \frac{1}{T} \sum_{t=1}^T (\epsilon_{it} - \bar{\epsilon}_{NT}) \\ \hat{\gamma}_t &:= \frac{1}{N} \sum_{i=1}^N (Y_{it} - \bar{Y}_{NT}) = \gamma_t + \frac{1}{N} \sum_{i=1}^N (\epsilon_{it} - \bar{\epsilon}_{NT}) \\ \hat{\epsilon}_{it} &:= Y_{it} - \bar{Y}_{NT} - \hat{\alpha}_i - \hat{\gamma}_t \end{aligned}$$

then use these empirical distributions for estimation and form a bootstrap sample

$$Y_{it}^* := \bar{Y}_{NT} + \alpha_i^* + \gamma_t^* + \epsilon_{it}^*$$

⁴This is specific to clustering in 2 or more dimensions.

by drawing from these estimators and obtain $\bar{Y}_{NT}^* := \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T Y_{it}^*$, and verify the conditional variances of the bootstrap distribution given the sample:

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N \left(\hat{\alpha}_i - \frac{1}{N} \sum_{j=1}^N \hat{\alpha}_j \right)^2 - \left[\text{Var}(\alpha_i) + \frac{\text{Var}(\epsilon_{it})}{T} \right] &\xrightarrow{p} 0 \\ \frac{1}{T} \sum_{t=1}^T \left(\hat{\gamma}_t - \frac{1}{T} \sum_{s=1}^T \hat{\gamma}_s \right)^2 - \left[\text{Var}(\gamma_t) + \frac{\text{Var}(\epsilon_{it})}{N} \right] &\xrightarrow{p} 0 \end{aligned}$$

then the bootstrap distribution is

– in the non-degenerate case,

$$\sqrt{N} \left(\bar{Y}_{NT}^* - \bar{Y}_{NT} \right) \xrightarrow{d} \mathcal{N} \left(0, \text{Var}(\alpha_i) + \text{Var}(\gamma_t) \right)$$

the estimation error $\hat{\alpha}_i$ does **NOT** affect the asymptotic variance.

– in the degenerate case,

$$\sqrt{NT} \left(\bar{Y}_{NT}^* - \bar{Y}_{NT} \right) \xrightarrow{d} \mathcal{N} \left(0, 3\text{Var}(\epsilon_{it}) \right)$$

asymptotically overestimates the variance of the sampling distribution, leading to inconsistency of this naive bootstrapping procedure.

- **Non-Gaussian limit distribution**

$$Y_{it} = \alpha_i \gamma_t + \epsilon_{it}$$

where $\alpha_i, \gamma_t, \epsilon_{it}$ are independently distributed with $\mathbb{E}[\epsilon_{it}] = 0$, $\text{Var}(\alpha_i) = \sigma_\alpha^2$, $\text{Var}(\gamma_t) = \sigma_\gamma^2$, $\text{Var}(\epsilon_{it}) = \sigma_\epsilon^2$.

If $\mathbb{E}[\alpha_i] = \mathbb{E}[\gamma_t] = 0$, then CLT and Continuous Mapping Theorem (CMT) imply

$$\begin{aligned} \sqrt{NT} \cdot \bar{Y}_{NT} &= \frac{1}{\sqrt{NT}} \sum_{i=1}^N \sum_{t=1}^T (\alpha_i \gamma_t + \epsilon_{it}) \\ &= \left(\frac{1}{\sqrt{N}} \sum_{i=1}^N \alpha_i \right) \left(\frac{1}{\sqrt{T}} \sum_{t=1}^T \gamma_t \right) + \frac{1}{\sqrt{NT}} \sum_{i=1}^N \sum_{t=1}^T \epsilon_{it} \\ &\xrightarrow{d} \sigma_\alpha \sigma_\gamma Z_1 Z_2 + \sigma_\epsilon Z_3 \end{aligned}$$

then even without correlation within clusters, non-separable heterogeneity can still generate dependence in 2nd or higher moments in the limiting distribution⁵.

5.4.2 Menzel (2021)'s Bootstrap procedure

5.4.2.1 Notation

For the array $(Y_{it})_{i,t}$, denote

- \mathbb{P} : joint distribution of $(Y_{it})_{i,t}$

⁵2 major issues arise:

- The limiting distribution needs **not** be Gaussian: plug-in asymptotic inference based on the normal distribution is invalid
- It only comes from two-or-more-dimension cluster dependence, not single-dimension cluster dependence.

- \mathbb{P}_{NT} : drifting DGP indexed by N, T
- \mathbb{P}_{NT}^* : bootstrap distribution for (Y_{it}^*) given the realizations $(Y_{it} : i = 1, \dots, N; t = 1, \dots, T)$
- respective distributions $\mathbb{E}, \mathbb{E}_{NT}, \mathbb{E}_{NT}^*$

5.4.2.2 Inference: Sample Mean

First, consider the assumption of *separate exchangeability*

Assumption 5.4.1: Separate Exchangeability

- A **separately exchangeable** array is an infinite array $(Y_{it})_{i,t}$ such that for any integers \tilde{N}, \tilde{T} and permutations $\pi_1 : \{1, \dots, \tilde{N}\} \rightarrow \{1, \dots, \tilde{N}\}$ and $\pi_2 : \{1, \dots, \tilde{T}\} \rightarrow \{1, \dots, \tilde{T}\}$, we have

$$(Y_{\pi_1(i), \pi_2(t)})_{i,t} \stackrel{d}{=} (Y_{it})_{i,t}$$

such an array is called **dissociated** if for any $N_0, T_0 \geq 1$, $(Y_{it})_{i=1, t=1}^{i=N_0, t=T_0}$ is independent of $(Y_{it})_{i>N_0, t>T_0}$.

- For dyadic data, consider the alternative assumption **jointly exchangeable** arrays $(Y_{ij})_{i,j}$ satisfying

$$(Y_{\pi(i), \pi(j)})_{i,j} \stackrel{d}{=} (Y_{ij})_{i,j}$$

for any permutation π on $\{1, \dots, \tilde{N}\}$, in addition, $(Y_{ij})_{i,j=1}^{N_0}$ is independent of $(Y_{ij})_{i,j>N_0}$

This assumption can be interpreted as rows (and columns) corresponding to units that are drawn independently from a common population, where we then observe the joint outcome for every row-column pair, consider the requirements in the following applications

- **DiD/matched data**: the units corresponding to either dimension of the sample to represent independent draws from a common, infinite population
- **non-exhaustively matched data**: only observe joint outcomes for a possibly self-selected subset of unit pairs, sample selection should be (jointly or separately) exchangeable
- **U-/V-statistics**: the kernel $Y_{i_1, \dots, i_D} := h(X_{i_1}, \dots, X_{i_D})$ evaluated at i.i.d. observations X_1, \dots, X_N forms a dissociated, jointly exchangeable array
- **Network**: unlabeled⁶ data implies finite exchangeability, the sampled graph has joint (*infinite*) exchangeability if it is a subgraph of an infinite graph

Directly from Assumption 5.4.1, any dissociated separately exchangeable array can be represented as

$$Y_{it} = f(\alpha_i, \gamma_t, \epsilon_{it})$$

for some function $f(\cdot)$ where $\alpha_1, \dots, \alpha_N, \gamma_1, \dots, \gamma_T, \epsilon_{11}, \dots, \epsilon_{NT}$ are mutually independent, uniformly distributed random variables.

Projection now, decompose the array $(Y_{it})_{i,t}$ as

$$Y_{it} = b + a_i + g_t + w_{it}$$

$$\mathbb{E}[w_{it} \mid a_i, g_t] = 0$$

⁶Unlabeled: model identifiers do not carry any significance for the statistical model.

where a_i and g_t are mean-zero and mutually independent, s.t. the joint distribution of Y_{it} can then be expanded as

$$\begin{aligned} Y_{it} &= \mathbb{E}[Y_{it}] + (\mathbb{E}[Y_{it} | \alpha_i] - \mathbb{E}[Y_{it}]) + (\mathbb{E}[Y_{it} | \gamma_t] - \mathbb{E}[Y_{it}]) \\ &\quad + (\mathbb{E}[Y_{it} | \alpha_i, \gamma_t] - \mathbb{E}[Y_{it} | \alpha_i] - \mathbb{E}[Y_{it} | \gamma_t] + \mathbb{E}[Y_{it}]) + (Y_{it} - \mathbb{E}[Y_{it} | \alpha_i, \gamma_t]) \\ &=: b + a_i + g_t + v_{it} + e_{it} \end{aligned}$$

with

- $e_{it} = Y_{it} - \mathbb{E}[Y_{it} | \alpha_i, \gamma_t]$
- $a_i = \mathbb{E}[Y_{it} | \alpha_i] - \mathbb{E}[Y_{it}]$, $g_t = \mathbb{E}[Y_{it} | \gamma_t] - \mathbb{E}[Y_{it}]$
- $v_{it} = \mathbb{E}[Y_{it} | \alpha_i, \gamma_t] - \mathbb{E}[Y_{it} | \alpha_i] - \mathbb{E}[Y_{it} | \gamma_t] + \mathbb{E}[Y_{it}]$
- $b = \mathbb{E}[Y_{it}]$

here,

- temporal and cross-sectional units were drawn independently: a_1, \dots, a_N and g_1, \dots, g_T are independent of each other.
- by construction, $\mathbb{E}[e_{it} | a_i, g_t, v_{it}] = 0$, $\mathbb{E}[v_{it} | a_i] = \mathbb{E}[v_{it} | g_t] = 0$
- e_{it} , (a_i, g_t) and v_{it} are **uncorrelated**

then, rewrite the sample mean as

$$\begin{aligned} \hat{Y}_{NT} &= b + \bar{a}_N + \bar{g}_T + \bar{v}_{NT} + \bar{e}_{NT} \\ &=: b + \frac{1}{N} \sum_{i=1}^N a_i + \frac{1}{T} \sum_{t=1}^T g_t + \frac{1}{NT} \sum_{t=1}^T \sum_{i=1}^N v_{it} + \frac{1}{NT} \sum_{t=1}^T \sum_{i=1}^N e_{it} \end{aligned}$$

and the unconditional variances of the projections with

$$\sigma_a^2 := \text{Var}(a_i) \quad \sigma_g^2 := \text{Var}(g_t) \quad \sigma_v^2 := \text{Var}(v_{it}) \quad \sigma_e^2 := \text{Var}(e_{it})$$

let $w_{it} := v_{it} + e_{it}$, and denote its variance by $\sigma_w^2 = \text{Var}(w_{it})$. Then, assume integrability

Assumption 5.4.2: Integrability

Let $Y_{it} = f(\alpha_i, \gamma_t, \epsilon_{it})$, where $\alpha_i, \gamma_t, \epsilon_{it}$ are random arrays with elements i.i.d. drawn from $[0, 1]$ uniform distribution, assume

- $a_i/\sigma_a, g_t/\sigma_g, v_{it}/\sigma_v, e_{it}/\sigma_e$ are well-defined and have bounded moments up to the order $4 + \delta$ for some $\delta > 0$, whenever the respective variances $\sigma_a^2, \sigma_g^2, \sigma_v^2, \sigma_e^2$ are non-zero.
- $\sigma_a^2 + \sigma_g^2 > 0$, or $\sigma_v^2 + \sigma_e^2 > 0$

Low-rank approximation Consider the row/column projection

$$\bar{v}_{NT} \equiv \frac{1}{NT} \sum_{t=1}^T \sum_{i=1}^N (\mathbb{E}[Y_{it} | \alpha_i, \gamma_t] - \mathbb{E}[Y_{it} | \alpha_i] - \mathbb{E}[Y_{it} | \gamma_t] + \mathbb{E}[Y_{it}]) =: \frac{1}{NT} \sum_{t=1}^T \sum_{i=1}^N v(\alpha_i, \gamma_t)$$

as a generalized U-statistic with a kernel $v(\alpha, \gamma)$ evaluated at the samples $\alpha_1, \dots, \alpha_N$ and $\gamma_1, \dots, \gamma_T$. There are 2 major issues w.r.t. characterizing the distribution of \bar{Y}_{NT}

- the presence of the projection error e_{it}
- the factors α_i, γ_t are not observable

Define,

$$v(\alpha, \gamma) := \mathbb{E}[Y_{it} \mid \alpha_i = \alpha, \gamma_t = \gamma] - \mathbb{E}[Y_{it} \mid \alpha_i = \alpha] - \mathbb{E}[Y_{it} \mid \gamma_t = \gamma] + \mathbb{E}[Y_{it}]$$

under Assumption 5.4.2, we have compact integral operators

$$S(u)(g) = \int v(a, g)u(a)F_\alpha(da) \quad S^*(u)(a) = \int v(a, g)u(g)F_\gamma(dg)$$

where F_α, F_γ are the marginal distributions corresponding to the joint $F_{\alpha\gamma}$ of α_i, γ_t . Then the low-rank approximation is

$$v(\alpha, \gamma) = \sum_{k=1}^{\infty} c_k \phi_k(\alpha) \psi_k(\gamma) \quad (5.14)$$

under the $L_2(F_{\alpha\gamma})$ norm on the space of smooth functions of $(\alpha, \gamma) \in [0, 1]^2$. Here

- $(c_k)_{k \geq 1}$: a sequence of singular values, $\lim |c_k| \rightarrow 0$
- $(\phi_k(\cdot))_{k \geq 1}$ and $(\psi_k(\cdot))_{k \geq 1}$: orthonormal bases for $L_2([0, 1], F_\alpha)$ and $L_2([0, 1], F_\gamma)$:
 - By construction:

$$\mathbb{E}[v(a, \gamma_t)] = \mathbb{E}[v(\alpha_i, g)] = 0, \forall a, g \in [0, 1] \Rightarrow \mathbb{E}[\phi_k(\alpha_i)] = \mathbb{E}[\psi_k(\gamma_t)] = 0, \forall k = 1, 2, \dots$$

- the basis functions are orthonormal and α_i and γ_t are independent, then $\forall K < \infty$

$$\text{Cov}[(\phi_1(\alpha_i), \psi_1(\gamma_t), \dots, \phi_K(\alpha_i), \psi_K(\gamma_t))]$$

is the $2K$ -dimensional identity matrix

- $(\phi_1(\alpha_i), \dots, \phi_K(\alpha_i))$ can be correlated with a_i : $\sigma_{ak} := \text{Cov}(a_i, \phi_k(\alpha_i))$
- $(\psi_1(\gamma_t), \dots, \psi_K(\gamma_t))$ can be correlated with g_t : $\sigma_{gk} := \text{Cov}(g_t, \psi_k(\gamma_t))$

with this representation of Eq.(5.14), we have⁷

$$\frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T v(\alpha_i, \gamma_t) = \sum_{k=1}^{\infty} c_k \left(\frac{1}{N} \sum_{i=1}^N \phi_k(\alpha_i) \right) \left(\frac{1}{T} \sum_{t=1}^T \psi_k(\gamma_t) \right)$$

and the second-order projection term can also be represented as a function of **countably many** sample averages of **i.i.d. mean-zero** random variables.

Assumption 5.4.3: Eigenfucntions and coefficients in the spectral representation (5.14)

The function $v(\alpha, \gamma) := \mathbb{E}[Y_{it} \mid \alpha_i = \alpha, \gamma_t = \gamma] - \mathbb{E}[Y_{it} \mid \alpha_i = \alpha] - \mathbb{E}[Y_{it} \mid \gamma_t = \gamma] + \mathbb{E}[Y_{it}]$ admits a spectral representation

$$v(\alpha, \gamma) = \sum_{k=1}^{\infty} c_k \phi_k(\alpha) \psi_k(\gamma)$$

under the $L_2(F_{\alpha\gamma})$ norm. And

- the singular values are uniformly bounded by a square summable null sequence \bar{c}_k : $c_k \leq \bar{c}_k, \forall k = 1, 2, \dots$, where $\sum_{k=1}^{\infty} \bar{c}_k^2 < \infty$

⁷The limiting distribution of this term is not Gaussian, but can be represented as a linear combination of independent chi-squared random variables. This type of distributions is known as Wiener/Gaussian chaos.

- $\forall k = 1, 2, \dots$, the first 3 moments of the eigenfunctions $\phi_k(\alpha_i)$ and $\psi_k(\gamma_t)$ are bounded by a constant $B > 0$

To summarize the two assumptions

- Assumption 5.4.1 guarantees the pointwise consistency of the bootstrap
- Assumption 5.4.3 gives the uniform consistency of the bootstrap: it imposes common bounds on moments and singular values and restricts the set of joint distribution F to a **uniformity** class⁸.

5.4.2.3 Bootstrap procedure

For the sample mean $\bar{Y}_{NT} - \mathbb{E}[Y_{it}]$, the limiting distribution depends on the scale parameters:

- If observations are independent across rows and columns: $\sqrt{NT} \left(\bar{Y}_{NT} - \mathbb{E}[Y_{it}] \right) \xrightarrow{d} \mathcal{N}(0, \sigma_e^2)$
- If $N = T$, within-cluster covariances are bounded from 0 in **at least one dimension**: $\sqrt{N} \left(\bar{Y}_{NT} - \mathbb{E}[Y_{it}] \right) \xrightarrow{d} \mathcal{N}(0, \sigma_a^2 + \sigma_g^2)$

The bootstrap procedure should then be adaptive for both degenerate and non-degenerate cases. For the expansion

$$\begin{aligned} Y_{it} &= \mathbb{E}[Y_{it}] + (\mathbb{E}[Y_{it} | \alpha_i] - \mathbb{E}[Y_{it}]) + (\mathbb{E}[Y_{it} | \gamma_t] - \mathbb{E}[Y_{it}]) \\ &\quad + (\mathbb{E}[Y_{it} | \alpha_i, \gamma_t] - \mathbb{E}[Y_{it} | \alpha_i] - \mathbb{E}[Y_{it} | \gamma_t] + \mathbb{E}[Y_{it}]) + (Y_{it} - \mathbb{E}[Y_{it} | \alpha_i, \gamma_t]) \\ &=: b + a_i + g_t + v_{it} + e_{it} \end{aligned} \quad (5.15)$$

the sample analogs are:

$$\hat{a}_i := \frac{1}{T} \sum_{t=1}^T Y_{it} - \bar{Y}_{NT} \quad \hat{g}_t := \frac{1}{N} \sum_{i=1}^N Y_{it} - \bar{Y}_{NT} \quad \hat{w}_{it} := Y_{it} - \hat{a}_i - \hat{g}_t - \bar{Y}_{NT}$$

Evaluating bootstrap performance it is crucial at what rates these estimators are consistent depending on the extent of clustering in the true DGP. The variance of the projection terms are:

$$\text{Var}(\hat{a}_i) = \sigma_a^2 + \frac{\sigma_w^2}{T} \quad \text{Var}(\hat{g}_t) = \sigma_g^2 + \frac{\sigma_w^2}{N}$$

s.t. the **convolution error** depending on σ_w^2 dominates in the degenerate case. Therefore, to correct for the contribution of the row/column averages of w_{it} , consider the scalar for the distribution of \hat{a}_i, \hat{g}_t by

$$\lambda_a = \frac{T\sigma_a^2}{T\sigma_a^2 + \sigma_w^2} \quad \lambda_g = \frac{N\sigma_g^2}{N\sigma_g^2 + \sigma_w^2}$$

⁸Here, the sequence $c := (\tilde{c})_{k \geq 0}$ controls the magnitude of the error from a finite-dimensional approximation to $v(\alpha, \gamma)$.

Component variance estimator let

$$\begin{aligned}\hat{s}_a^2 &:= \frac{1}{N-1} \sum_{i=1}^N \left(\hat{a}_i - \bar{Y}_{NT} \right)^2 \\ \hat{s}_g^2 &:= \frac{1}{T-1} \sum_{t=1}^T \left(\hat{g}_t - \bar{Y}_{NT} \right)^2 \\ \hat{s}_w^2 &:= \frac{1}{NT - N - T} \sum_{i=1}^N \sum_{t=1}^T \left(Y_{it} - \hat{a}_i - \hat{g}_t - \bar{Y}_{NT} \right)^2\end{aligned}$$

then form the estimators as

$$\hat{\sigma}_a^2 = \max \left\{ 0, \hat{s}_a^2 - \frac{1}{T} \hat{s}_w^2 \right\} \quad \hat{\sigma}_g^2 = \max \left\{ 0, \hat{s}_g^2 - \frac{1}{N} \hat{s}_w^2 \right\} \quad \hat{\sigma}_w^2 := \hat{s}_w^2 \quad (5.16)$$

A Theoretical

Chiang and Sasaki (2023)

References

- Manuel Arellano. Computing robust standard errors for within-groups estimators. *Oxford bulletin of Economics and Statistics*, 49(4):431–434, 1987.
- A Colin Cameron, Jonah B Gelbach, and Douglas L Miller. Robust inference with multiway clustering. *Journal of Business & Economic Statistics*, 29(2):238–249, 2011.
- Harold D Chiang and Yuya Sasaki. On using the two-way cluster-robust standard errors. *arXiv preprint arXiv:2301.13775*, 2023.
- Christian B Hansen. Asymptotic properties of a robust variance matrix estimator for panel data when t is large. *Journal of Econometrics*, 141(2):597–620, 2007.
- Teunis Kloek. Ols estimation in a model where a microvariable is explained by aggregates and contemporaneous disturbances are equicorrelated. *Econometrica: Journal of the Econometric Society*, pages 205–207, 1981.
- Kung-Yee Liang and Scott L Zeger. Longitudinal data analysis using generalized linear models. *Biometrika*, 73(1):13–22, 1986.
- Konrad Menzel. Bootstrap with cluster-dependence in two or more dimensions. *Econometrica*, 89(5):2143–2188, 2021.
- Brent R Moulton. An illustration of a pitfall in estimating the effects of aggregate variables on micro units. *The review of Economics and Statistics*, pages 334–338, 1990.
- Andrew J Scott and D Holt. The effect of two-stage sampling on ordinary least squares methods. *Journal of the American statistical Association*, 77(380):848–854, 1982.
- Samuel B Thompson. Simple formulas for standard errors that cluster by both firm and time. *Journal of financial Economics*, 99(1):1–10, 2011.
- Halbert White. A heteroskedasticity-consistent covariance matrix estimator and a direct test for heteroskedasticity. *Econometrica: journal of the Econometric Society*, pages 817–838, 1980.