

# Jackknife Inference with Two-Way Clustering

James G. MacKinnon (Queen's University and ACE)  
Morten Ørregaard Nielsen (Aarhus University)  
Matthew D. Webb (Carleton University)

Workshop to Honour Lynda Khalaf, October 17, 2025



# Introduction

# Introduction

For models with cross-section or panel data, the disturbances may be clustered in two dimensions.

# Introduction

For models with cross-section or panel data, the disturbances may be clustered in two dimensions.

Unfortunately, the finite-sample properties of two-way cluster-robust tests and confidence intervals are often poor.

# Introduction

For models with cross-section or panel data, the disturbances may be clustered in two dimensions.

Unfortunately, the finite-sample properties of two-way cluster-robust tests and confidence intervals are often poor.

There can also be undefined standard errors when a cluster-robust variance estimator (CRVE) is not positive definite.

# Introduction

For models with cross-section or panel data, the disturbances may be clustered in two dimensions.

Unfortunately, the finite-sample properties of two-way cluster-robust tests and confidence intervals are often poor.

There can also be undefined standard errors when a cluster-robust variance estimator (CRVE) is not positive definite.

We discuss ways to improve inference with two-way clustering.

# Introduction

For models with cross-section or panel data, the disturbances may be clustered in two dimensions.

Unfortunately, the finite-sample properties of two-way cluster-robust tests and confidence intervals are often poor.

There can also be undefined standard errors when a cluster-robust variance estimator (CRVE) is not positive definite.

We discuss ways to improve inference with two-way clustering.

- Two existing methods for dealing with undefined standard errors.

# Introduction

For models with cross-section or panel data, the disturbances may be clustered in two dimensions.

Unfortunately, the finite-sample properties of two-way cluster-robust tests and confidence intervals are often poor.

There can also be undefined standard errors when a cluster-robust variance estimator (CRVE) is not positive definite.

We discuss ways to improve inference with two-way clustering.

- Two existing methods for dealing with undefined standard errors.
- A new method that completely avoids undefined standard errors.

# Introduction

For models with cross-section or panel data, the disturbances may be clustered in two dimensions.

Unfortunately, the finite-sample properties of two-way cluster-robust tests and confidence intervals are often poor.

There can also be undefined standard errors when a cluster-robust variance estimator (CRVE) is not positive definite.

We discuss ways to improve inference with two-way clustering.

- Two existing methods for dealing with undefined standard errors.
- A new method that completely avoids undefined standard errors.
- Several new two-way CRVEs based on the cluster jackknife.

# Introduction

For models with cross-section or panel data, the disturbances may be clustered in two dimensions.

Unfortunately, the finite-sample properties of two-way cluster-robust tests and confidence intervals are often poor.

There can also be undefined standard errors when a cluster-robust variance estimator (CRVE) is not positive definite.

We discuss ways to improve inference with two-way clustering.

- Two existing methods for dealing with undefined standard errors.
- A new method that completely avoids undefined standard errors.
- Several new two-way CRVEs based on the cluster jackknife.

Simulations for models with two-way fixed effects suggest that a cluster-jackknife CRVE based on the new method often yields surprisingly accurate inferences.

# Literature

# Literature

- One-way CRVEs: Liang and Zeger (1986); Arellano (1987); Bester, Conley, and Hansen (2011); Hansen and Lee (2019); Djogbenou, MacKinnon, and Nielsen (2019).

# Literature

- One-way CRVEs: Liang and Zeger (1986); Arellano (1987); Bester, Conley, and Hansen (2011); Hansen and Lee (2019); Djogbenou, MacKinnon, and Nielsen (2019).
- Two-way CRVEs: Cameron, Gelbach, and Miller (2011); Miglioretti and Heagerty (2006); Thompson (2011).

# Literature

- One-way CRVEs: Liang and Zeger (1986); Arellano (1987); Bester, Conley, and Hansen (2011); Hansen and Lee (2019); Djogbenou, MacKinnon, and Nielsen (2019).
- Two-way CRVEs: Cameron, Gelbach, and Miller (2011); Miglioretti and Heagerty (2006); Thompson (2011).
- Theory for two-way CRVEs: Davezies, D'Haultfoeuille, and Guyonvarch (2021, 2025); Menzel (2021); Chiang, Kato, and Sasaki (2023), Yap (2025).

# Literature

- One-way CRVEs: Liang and Zeger (1986); Arellano (1987); Bester, Conley, and Hansen (2011); Hansen and Lee (2019); Djogbenou, MacKinnon, and Nielsen (2019).
- Two-way CRVEs: Cameron, Gelbach, and Miller (2011); Miglioretti and Heagerty (2006); Thompson (2011).
- Theory for two-way CRVEs: Davezies, D'Haultfoeuille, and Guyonvarch (2021, 2025); Menzel (2021); Chiang, Kato, and Sasaki (2023), Yap (2025).
- Wild bootstrap methods: MacKinnon, Neilsen, and Webb (2021); Hounyo and Lin (2024). `boottest`. [Not entirely satisfactory]

# Literature

- One-way CRVEs: Liang and Zeger (1986); Arellano (1987); Bester, Conley, and Hansen (2011); Hansen and Lee (2019); Djogbenou, MacKinnon, and Nielsen (2019).
- Two-way CRVEs: Cameron, Gelbach, and Miller (2011); Miglioretti and Heagerty (2006); Thompson (2011).
- Theory for two-way CRVEs: Davezies, D'Haultfoeuille, and Guyonvarch (2021, 2025); Menzel (2021); Chiang, Kato, and Sasaki (2023), Yap (2025).
- Wild bootstrap methods: MacKinnon, Neilsen, and Webb (2021); Hounyo and Lin (2024). `boottest`. [Not entirely satisfactory]
- Jackknife variance estimation: Tukey (1958); Efron (1981); Efron and Stein (1981); MacKinnon and White (1985); Bell and McCaffrey (2002); MacKinnon, Nielsen, and Webb (JAE 2023, SJ 2023); Hansen (2025 WP & JAE).

# Linear Regression with Two-Way Clustering

# Linear Regression with Two-Way Clustering

With two-way clustering,  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$  can be written as

$$\mathbf{y}_{gh} = \mathbf{X}_{gh}\boldsymbol{\beta} + \mathbf{u}_{gh}, \quad g = 1, \dots, G, \quad h = 1, \dots, H. \quad (2)$$

# Linear Regression with Two-Way Clustering

With two-way clustering,  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$  can be written as

$$\mathbf{y}_{gh} = \mathbf{X}_{gh}\boldsymbol{\beta} + \mathbf{u}_{gh}, \quad g = 1, \dots, G, \quad h = 1, \dots, H. \quad (2)$$

The variance matrix of the score vector  $\mathbf{s} = \mathbf{X}^\top \mathbf{u}$  is

$$\boldsymbol{\Sigma} = \sum_{g,g'=1}^G \sum_{h,h'=1}^H \text{E}(\mathbf{s}_{gh}\mathbf{s}_{g'h'}^\top) = \sum_{g=1}^G \boldsymbol{\Sigma}_g + \sum_{h=1}^H \boldsymbol{\Sigma}_h - \sum_{g=1}^G \sum_{h=1}^H \boldsymbol{\Sigma}_{gh}. \quad (3)$$

# Linear Regression with Two-Way Clustering

With two-way clustering,  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$  can be written as

$$\mathbf{y}_{gh} = \mathbf{X}_{gh}\boldsymbol{\beta} + \mathbf{u}_{gh}, \quad g = 1, \dots, G, \quad h = 1, \dots, H. \quad (2)$$

The variance matrix of the score vector  $\mathbf{s} = \mathbf{X}^\top \mathbf{u}$  is

$$\boldsymbol{\Sigma} = \sum_{g,g'=1}^G \sum_{h,h'=1}^H \mathbb{E}(\mathbf{s}_{gh}\mathbf{s}_{g'h'}^\top) = \sum_{g=1}^G \boldsymbol{\Sigma}_g + \sum_{h=1}^H \boldsymbol{\Sigma}_h - \sum_{g=1}^G \sum_{h=1}^H \boldsymbol{\Sigma}_{gh}. \quad (3)$$

where

$$\boldsymbol{\Sigma}_g = \mathbb{E}(\mathbf{s}_g\mathbf{s}_g^\top), \quad \boldsymbol{\Sigma}_h = \mathbb{E}(\mathbf{s}_h\mathbf{s}_h^\top), \quad \text{and} \quad \boldsymbol{\Sigma}_{gh} = \mathbb{E}(\mathbf{s}_{gh}\mathbf{s}_{gh}^\top). \quad (4)$$

# Linear Regression with Two-Way Clustering

With two-way clustering,  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$  can be written as

$$\mathbf{y}_{gh} = \mathbf{X}_{gh}\boldsymbol{\beta} + \mathbf{u}_{gh}, \quad g = 1, \dots, G, \quad h = 1, \dots, H. \quad (2)$$

The variance matrix of the score vector  $\mathbf{s} = \mathbf{X}^\top \mathbf{u}$  is

$$\boldsymbol{\Sigma} = \sum_{g,g'=1}^G \sum_{h,h'=1}^H \text{E}(\mathbf{s}_{gh}\mathbf{s}_{g'h'}^\top) = \sum_{g=1}^G \boldsymbol{\Sigma}_g + \sum_{h=1}^H \boldsymbol{\Sigma}_h - \sum_{g=1}^G \sum_{h=1}^H \boldsymbol{\Sigma}_{gh}. \quad (3)$$

where

$$\boldsymbol{\Sigma}_g = \text{E}(\mathbf{s}_g\mathbf{s}_g^\top), \quad \boldsymbol{\Sigma}_h = \text{E}(\mathbf{s}_h\mathbf{s}_h^\top), \quad \text{and} \quad \boldsymbol{\Sigma}_{gh} = \text{E}(\mathbf{s}_{gh}\mathbf{s}_{gh}^\top). \quad (4)$$

The variance matrix of  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$  is

$$\mathbf{V}_{\hat{\boldsymbol{\beta}}} = (\mathbf{X}^\top \mathbf{X})^{-1} \boldsymbol{\Sigma} (\mathbf{X}^\top \mathbf{X})^{-1} = \mathbf{V}_G + \mathbf{V}_H - \mathbf{V}_I. \quad (5)$$

The first component of  $V_\beta$  is

$$V_G = (\mathbf{X}^\top \mathbf{X})^{-1} \left( \sum_{g=1}^G \Sigma_g \right) (\mathbf{X}^\top \mathbf{X})^{-1}. \quad (6)$$

The other two,  $V_H$  and  $V_I$ , are defined similarly.

The first component of  $V_\beta$  is

$$V_G = (\mathbf{X}^\top \mathbf{X})^{-1} \left( \sum_{g=1}^G \Sigma_g \right) (\mathbf{X}^\top \mathbf{X})^{-1}. \quad (6)$$

The other two,  $V_H$  and  $V_L$ , are defined similarly.

The empirical analog of (5) is the three-term two-way CRVE

$$\hat{V}_1^{(3)} = \hat{V}_G + \hat{V}_H - \hat{V}_L, \quad (7)$$

The first component of  $V_\beta$  is

$$V_G = (\mathbf{X}^\top \mathbf{X})^{-1} \left( \sum_{g=1}^G \Sigma_g \right) (\mathbf{X}^\top \mathbf{X})^{-1}. \quad (6)$$

The other two,  $V_H$  and  $V_I$ , are defined similarly.

The empirical analog of (5) is the three-term two-way CRVE

$$\hat{V}_1^{(3)} = \hat{V}_G + \hat{V}_H - \hat{V}_I, \quad (7)$$

where

$$\hat{V}_I = \frac{I(N-1)}{(I-1)(N-k)} (\mathbf{X}^\top \mathbf{X})^{-1} \left( \sum_{g=1}^G \sum_{h=1}^H \hat{\mathbf{s}}_{gh} \hat{\mathbf{s}}_{gh}^\top \right) (\mathbf{X}^\top \mathbf{X})^{-1}, \quad (8)$$

and likewise for  $V_G$  and  $V_H$ . Here  $I \leq GH$  is the number of *intersections*.

The first component of  $V_\beta$  is

$$V_G = (\mathbf{X}^\top \mathbf{X})^{-1} \left( \sum_{g=1}^G \Sigma_g \right) (\mathbf{X}^\top \mathbf{X})^{-1}. \quad (6)$$

The other two,  $V_H$  and  $V_I$ , are defined similarly.

The empirical analog of (5) is the three-term two-way CRVE

$$\hat{V}_1^{(3)} = \hat{V}_G + \hat{V}_H - \hat{V}_I, \quad (7)$$

where

$$\hat{V}_I = \frac{I(N-1)}{(I-1)(N-k)} (\mathbf{X}^\top \mathbf{X})^{-1} \left( \sum_{g=1}^G \sum_{h=1}^H \hat{\mathbf{s}}_{gh} \hat{\mathbf{s}}_{gh}^\top \right) (\mathbf{X}^\top \mathbf{X})^{-1}, \quad (8)$$

and likewise for  $V_G$  and  $V_H$ . Here  $I \leq GH$  is the number of *intersections*.

But  $\hat{V}_1^{(3)}$  may not be positive definite!

We could instead use the two-term estimator

$$\hat{V}_1^{(2)} = \hat{V}_G + \hat{V}_H. \quad (9)$$

We could instead use the two-term estimator

$$\hat{V}_1^{(2)} = \hat{V}_G + \hat{V}_H. \quad (9)$$

It is easier to compute than  $\hat{V}_1^{(3)}$  and must be positive semi-definite.

We could instead use the two-term estimator

$$\hat{V}_1^{(2)} = \hat{V}_G + \hat{V}_H. \quad (9)$$

It is easier to compute than  $\hat{V}_1^{(3)}$  and must be positive semi-definite.

But  $\hat{V}_1^{(2)}$  is asymptotically invalid when the scores are independent or only correlated at the intersection level, so that  $V_\beta = V_I$ .

We could instead use the two-term estimator

$$\hat{V}_1^{(2)} = \hat{V}_G + \hat{V}_H. \quad (9)$$

It is easier to compute than  $\hat{V}_1^{(3)}$  and must be positive semi-definite.

But  $\hat{V}_1^{(2)}$  is asymptotically invalid when the scores are independent or only correlated at the intersection level, so that  $V_\beta = V_I$ .

In that case,  $\hat{V}_G \approx \hat{V}_H \approx \hat{V}_I$ . Therefore,

$$\hat{V}_1^{(2)} = \hat{V}_G + \hat{V}_H \approx 2\hat{V}_I, \quad (10)$$

We could instead use the two-term estimator

$$\hat{V}_1^{(2)} = \hat{V}_G + \hat{V}_H. \quad (9)$$

It is easier to compute than  $\hat{V}_1^{(3)}$  and must be positive semi-definite.

But  $\hat{V}_1^{(2)}$  is asymptotically invalid when the scores are independent or only correlated at the intersection level, so that  $V_\beta = V_I$ .

In that case,  $\hat{V}_G \approx \hat{V}_H \approx \hat{V}_I$ . Therefore,

$$\hat{V}_1^{(2)} = \hat{V}_G + \hat{V}_H \approx 2\hat{V}_I, \quad (10)$$

whereas

$$\hat{V}_1^{(3)} = \hat{V}_G + \hat{V}_H - \hat{V}_I \approx \hat{V}_I. \quad (11)$$

We could instead use the two-term estimator

$$\hat{V}_1^{(2)} = \hat{V}_G + \hat{V}_H. \quad (9)$$

It is easier to compute than  $\hat{V}_1^{(3)}$  and must be positive semi-definite.

But  $\hat{V}_1^{(2)}$  is asymptotically invalid when the scores are independent or only correlated at the intersection level, so that  $V_\beta = V_I$ .

In that case,  $\hat{V}_G \approx \hat{V}_H \approx \hat{V}_I$ . Therefore,

$$\hat{V}_1^{(2)} = \hat{V}_G + \hat{V}_H \approx 2\hat{V}_I, \quad (10)$$

whereas

$$\hat{V}_1^{(3)} = \hat{V}_G + \hat{V}_H - \hat{V}_I \approx \hat{V}_I. \quad (11)$$

Thus, in this case,  $\hat{V}_1^{(2)}$  is approximately twice as large as  $\hat{V}_1^{(3)}$ , and twice as large as it should be.

We could instead use the two-term estimator

$$\hat{V}_1^{(2)} = \hat{V}_G + \hat{V}_H. \quad (9)$$

It is easier to compute than  $\hat{V}_1^{(3)}$  and must be positive semi-definite.

But  $\hat{V}_1^{(2)}$  is asymptotically invalid when the scores are independent or only correlated at the intersection level, so that  $V_\beta = V_I$ .

In that case,  $\hat{V}_G \approx \hat{V}_H \approx \hat{V}_I$ . Therefore,

$$\hat{V}_1^{(2)} = \hat{V}_G + \hat{V}_H \approx 2\hat{V}_I, \quad (10)$$

whereas

$$\hat{V}_1^{(3)} = \hat{V}_G + \hat{V}_H - \hat{V}_I \approx \hat{V}_I. \quad (11)$$

Thus, in this case,  $\hat{V}_1^{(2)}$  is approximately twice as large as  $\hat{V}_1^{(3)}$ , and twice as large as it should be.

Unfortunately,  $\hat{V}_1^{(3)}$  is not necessarily positive semi-definite, and its diagonal elements may be negative.

Cameron, Gelbach, and Miller (2011) suggested a work-around to avoid negative diagonals, which Stata 18 now implements.

Cameron, Gelbach, and Miller (2011) suggested a work-around to avoid negative diagonals, which Stata 18 now implements.

Compute the eigenvalues of  $\hat{V}_1^{(3)}$ , say  $\lambda_1, \dots, \lambda_k$ .

Cameron, Gelbach, and Miller (2011) suggested a work-around to avoid negative diagonals, which Stata 18 now implements.

Compute the eigenvalues of  $\hat{\mathbf{V}}_1^{(3)}$ , say  $\lambda_1, \dots, \lambda_k$ .

When any of them is not positive, replace  $\hat{\mathbf{V}}_1^{(3)}$  by

$$\hat{\mathbf{V}}_1^{(3+)} = \mathbf{U}\mathbf{\Lambda}^+\mathbf{U}^\top,$$

where  $\mathbf{U}$  is the  $k \times k$  matrix of eigenvectors, and  $\mathbf{\Lambda}^+$  is a diagonal matrix with typical diagonal element  $\lambda_j^+ = \max\{\lambda_j, 0\}$ .

Cameron, Gelbach, and Miller (2011) suggested a work-around to avoid negative diagonals, which Stata 18 now implements.

Compute the eigenvalues of  $\hat{\mathbf{V}}_1^{(3)}$ , say  $\lambda_1, \dots, \lambda_k$ .

When any of them is not positive, replace  $\hat{\mathbf{V}}_1^{(3)}$  by

$$\hat{\mathbf{V}}_1^{(3+)} = \mathbf{U}\mathbf{\Lambda}^+\mathbf{U}^\top,$$

where  $\mathbf{U}$  is the  $k \times k$  matrix of eigenvectors, and  $\mathbf{\Lambda}^+$  is a diagonal matrix with typical diagonal element  $\lambda_j^+ = \max\{\lambda_j, 0\}$ .

Instead of 0, we use  $\eta = 10^{-12}$ , so that  $\hat{\mathbf{V}}_1^{(3+)}$  is positive definite.

Cameron, Gelbach, and Miller (2011) suggested a work-around to avoid negative diagonals, which Stata 18 now implements.

Compute the eigenvalues of  $\hat{\mathbf{V}}_1^{(3)}$ , say  $\lambda_1, \dots, \lambda_k$ .

When any of them is not positive, replace  $\hat{\mathbf{V}}_1^{(3)}$  by

$$\hat{\mathbf{V}}_1^{(3+)} = \mathbf{U}\mathbf{\Lambda}^+\mathbf{U}^\top,$$

where  $\mathbf{U}$  is the  $k \times k$  matrix of eigenvectors, and  $\mathbf{\Lambda}^+$  is a diagonal matrix with typical diagonal element  $\lambda_j^+ = \max\{\lambda_j, 0\}$ .

Instead of 0, we use  $\eta = 10^{-12}$ , so that  $\hat{\mathbf{V}}_1^{(3+)}$  is positive definite.

- Wald and  $t$ -statistics based on  $\hat{\mathbf{V}}_1^{(3+)}$  may be extremely large.

Cameron, Gelbach, and Miller (2011) suggested a work-around to avoid negative diagonals, which Stata 18 now implements.

Compute the eigenvalues of  $\hat{V}_1^{(3)}$ , say  $\lambda_1, \dots, \lambda_k$ .

When any of them is not positive, replace  $\hat{V}_1^{(3)}$  by

$$\hat{V}_1^{(3+)} = \mathbf{U}\mathbf{\Lambda}^+\mathbf{U}^\top,$$

where  $\mathbf{U}$  is the  $k \times k$  matrix of eigenvectors, and  $\mathbf{\Lambda}^+$  is a diagonal matrix with typical diagonal element  $\lambda_j^+ = \max\{\lambda_j, 0\}$ .

Instead of 0, we use  $\eta = 10^{-12}$ , so that  $\hat{V}_1^{(3+)}$  is positive definite.

- Wald and  $t$ -statistics based on  $\hat{V}_1^{(3+)}$  may be extremely large.
- Replacing  $\hat{V}_1^{(3)}$  by  $\hat{V}_1^{(3+)}$  can change **all** the standard errors.

Cameron, Gelbach, and Miller (2011) suggested a work-around to avoid negative diagonals, which Stata 18 now implements.

Compute the eigenvalues of  $\hat{V}_1^{(3)}$ , say  $\lambda_1, \dots, \lambda_k$ .

When any of them is not positive, replace  $\hat{V}_1^{(3)}$  by

$$\hat{V}_1^{(3+)} = \mathbf{U}\mathbf{\Lambda}^+\mathbf{U}^\top,$$

where  $\mathbf{U}$  is the  $k \times k$  matrix of eigenvectors, and  $\mathbf{\Lambda}^+$  is a diagonal matrix with typical diagonal element  $\lambda_j^+ = \max\{\lambda_j, 0\}$ .

Instead of 0, we use  $\eta = 10^{-12}$ , so that  $\hat{V}_1^{(3+)}$  is positive definite.

- Wald and  $t$ -statistics based on  $\hat{V}_1^{(3+)}$  may be extremely large.
- Replacing  $\hat{V}_1^{(3)}$  by  $\hat{V}_1^{(3+)}$  can change **all** the standard errors.
- $\text{se}(\hat{\beta}_j)$  is not invariant to nonsingular transformations of the remaining columns of the matrix  $\mathbf{X}$ .

Cameron, Gelbach, and Miller (2011) suggested a work-around to avoid negative diagonals, which Stata 18 now implements.

Compute the eigenvalues of  $\hat{V}_1^{(3)}$ , say  $\lambda_1, \dots, \lambda_k$ .

When any of them is not positive, replace  $\hat{V}_1^{(3)}$  by

$$\hat{V}_1^{(3+)} = \mathbf{U}\mathbf{\Lambda}^+\mathbf{U}^\top,$$

where  $\mathbf{U}$  is the  $k \times k$  matrix of eigenvectors, and  $\mathbf{\Lambda}^+$  is a diagonal matrix with typical diagonal element  $\lambda_j^+ = \max\{\lambda_j, 0\}$ .

Instead of 0, we use  $\eta = 10^{-12}$ , so that  $\hat{V}_1^{(3+)}$  is positive definite.

- Wald and  $t$ -statistics based on  $\hat{V}_1^{(3+)}$  may be extremely large.
- Replacing  $\hat{V}_1^{(3)}$  by  $\hat{V}_1^{(3+)}$  can change **all** the standard errors.
- $\text{se}(\hat{\beta}_j)$  is not invariant to nonsingular transformations of the remaining columns of the matrix  $\mathbf{X}$ .
- Precisely how fixed effects or other dummy variables are specified may affect  $\text{se}(\hat{\beta}_j)$ .

# A Better Way to Avoid Non-Positive Test Statistics

# A Better Way to Avoid Non-Positive Test Statistics

Compute three test statistics. Use the smallest one that is positive.  
Also suggested in [Davezies, D'Haultfoeuille, and Guyonvarch \(2025\)](#).

# A Better Way to Avoid Non-Positive Test Statistics

Compute three test statistics. Use the smallest one that is positive. Also suggested in [Davezies, D'Haultfoeuille, and Guyonvarch \(2025\)](#).

For the hypothesis that  $\mathbf{R}\boldsymbol{\beta} = \mathbf{r}$ , the three Wald statistics are

$$\begin{aligned} W_3 &= (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})^\top (\mathbf{R}\hat{\mathbf{V}}_1^{(3)}\mathbf{R}^\top)^{-1}(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r}), \\ W_G &= (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})^\top (\mathbf{R}\hat{\mathbf{V}}_G\mathbf{R}^\top)^{-1}(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r}), \text{ and} \\ W_H &= (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})^\top (\mathbf{R}\hat{\mathbf{V}}_H\mathbf{R}^\top)^{-1}(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r}). \end{aligned} \tag{12}$$

## A Better Way to Avoid Non-Positive Test Statistics

Compute three test statistics. Use the smallest one that is positive. Also suggested in [Davezies, D'Haultfoeuille, and Guyonvarch \(2025\)](#).

For the hypothesis that  $\mathbf{R}\boldsymbol{\beta} = \mathbf{r}$ , the three Wald statistics are

$$\begin{aligned} W_3 &= (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})^\top (\mathbf{R}\hat{\mathbf{V}}_1^{(3)}\mathbf{R}^\top)^{-1}(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r}), \\ W_G &= (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})^\top (\mathbf{R}\hat{\mathbf{V}}_G\mathbf{R}^\top)^{-1}(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r}), \text{ and} \\ W_H &= (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})^\top (\mathbf{R}\hat{\mathbf{V}}_H\mathbf{R}^\top)^{-1}(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r}). \end{aligned} \quad (12)$$

Our **max-se** procedure uses the statistic

$$W_{\min} = \min \{ \max\{W_3, 0\}, W_G, W_H \}, \quad (13)$$

where  $\max\{W_3, 0\} = 0$  if  $W_3$  is either negative or undefined.

## A Better Way to Avoid Non-Positive Test Statistics

Compute three test statistics. Use the smallest one that is positive. Also suggested in [Davezies, D'Haultfoeuille, and Guyonvarch \(2025\)](#).

For the hypothesis that  $R\beta = r$ , the three Wald statistics are

$$\begin{aligned} W_3 &= (R\hat{\beta} - r)^\top (R\hat{V}_1^{(3)}R^\top)^{-1}(R\hat{\beta} - r), \\ W_G &= (R\hat{\beta} - r)^\top (R\hat{V}_GR^\top)^{-1}(R\hat{\beta} - r), \text{ and} \\ W_H &= (R\hat{\beta} - r)^\top (R\hat{V}_HR^\top)^{-1}(R\hat{\beta} - r). \end{aligned} \quad (12)$$

Our **max-se** procedure uses the statistic

$$W_{\min} = \min \{ \max\{W_3, 0\}, W_G, W_H \}, \quad (13)$$

where  $\max\{W_3, 0\} = 0$  if  $W_3$  is either negative or undefined.

We denote the variance and standard error estimators based on  $\hat{V}_1^{(2)}$ ,  $\hat{V}_1^{(3)}$ , and  $\hat{V}_1^{(3+)}$  as  $CV_1^{(2)}$ ,  $CV_1^{(3)}$ , and  $CV_1^{(3+)}$ , respectively, and the one that is implicit in (13) as the  $CV_1^{(\max)}$  estimator.

# Two-Way Cluster Jackknife CRVEs

## Two-Way Cluster Jackknife CRVEs

Let  $J \in \{G, H, I\}$ , and let  $j$  denote the corresponding lower-case letter. The OLS estimates of  $\beta$  when each cluster in the  $J$  dimension is omitted in turn are

$$\hat{\beta}^{(j)} = (\mathbf{X}^\top \mathbf{X} - \mathbf{X}_j^\top \mathbf{X}_j)^{-1} (\mathbf{X}^\top \mathbf{y} - \mathbf{X}_j^\top \mathbf{y}_j), \quad j = 1, \dots, J. \quad (14)$$

## Two-Way Cluster Jackknife CRVEs

Let  $J \in \{G, H, I\}$ , and let  $j$  denote the corresponding lower-case letter. The OLS estimates of  $\beta$  when each cluster in the  $J$  dimension is omitted in turn are

$$\hat{\beta}^{(j)} = (\mathbf{X}^\top \mathbf{X} - \mathbf{X}_j^\top \mathbf{X}_j)^{-1} (\mathbf{X}^\top \mathbf{y} - \mathbf{X}_j^\top \mathbf{y}_j), \quad j = 1, \dots, J. \quad (14)$$

Then the component cluster jackknife variance matrix estimators are

$$\hat{\mathbf{V}}_J^{\text{JK}} = \frac{J-1}{J} \sum_{j=1}^J (\hat{\beta}^{(j)} - \hat{\beta})(\hat{\beta}^{(j)} - \hat{\beta})^\top \quad \text{for } \{j, J\} = \{g, G\}, \{h, H\}, \{i, I\}. \quad (15)$$

## Two-Way Cluster Jackknife CRVEs

Let  $J \in \{G, H, I\}$ , and let  $j$  denote the corresponding lower-case letter. The OLS estimates of  $\beta$  when each cluster in the  $J$  dimension is omitted in turn are

$$\hat{\beta}^{(j)} = (\mathbf{X}^\top \mathbf{X} - \mathbf{X}_j^\top \mathbf{X}_j)^{-1} (\mathbf{X}^\top \mathbf{y} - \mathbf{X}_j^\top \mathbf{y}_j), \quad j = 1, \dots, J. \quad (14)$$

Then the component cluster jackknife variance matrix estimators are

$$\hat{\mathbf{V}}_J^{\text{JK}} = \frac{J-1}{J} \sum_{j=1}^J (\hat{\beta}^{(j)} - \hat{\beta})(\hat{\beta}^{(j)} - \hat{\beta})^\top \quad \text{for } \{j, J\} = \{g, G\}, \{h, H\}, \{i, I\}. \quad (15)$$

Thus the three-term jackknife CRVE is

$$\hat{\mathbf{V}}_3^{(3)} = \hat{\mathbf{V}}_G^{\text{JK}} + \hat{\mathbf{V}}_H^{\text{JK}} - \hat{\mathbf{V}}_I^{\text{JK}}, \quad (16)$$

which is analogous to (7). Notation is based on  $\text{HC}_3$ .

# Computation

# Computation

First, calculate the cluster-level matrices and vectors

$$\mathbf{X}_j^\top \mathbf{X}_j \text{ and } \mathbf{X}_j^\top \mathbf{y}_j, \quad j = 1, \dots, J, \quad \text{for } \{j, J\} = \{g, G\}, \{h, H\}, \{i, I\}. \quad (17)$$

The ones for the intersections can be computed in a single pass over the  $N$  observations. The others are just sums of some of them.

# Computation

First, calculate the cluster-level matrices and vectors

$$\mathbf{X}_j^\top \mathbf{X}_j \text{ and } \mathbf{X}_j^\top \mathbf{y}_j, \quad j = 1, \dots, J, \text{ for } \{j, J\} = \{g, G\}, \{h, H\}, \{i, I\}. \quad (17)$$

The ones for the intersections can be computed in a single pass over the  $N$  observations. The others are just sums of some of them.

With two-way fixed effects in the  $G$  and  $H$  dimensions,

$$\mathbf{y} = \mathbf{Z}\boldsymbol{\beta}_p + \mathbf{D}^G\boldsymbol{\gamma} + \mathbf{D}^H\boldsymbol{\delta} + \mathbf{u}. \quad (18)$$

Now  $\mathbf{X} = [\mathbf{Z} \ \mathbf{D}^G \ \mathbf{D}^H]$ , and  $k = p + G + H - 1$ , and the matrices  $\mathbf{X}^\top \mathbf{X} - \mathbf{X}_g^\top \mathbf{X}_g$  and  $\mathbf{X}^\top \mathbf{X} - \mathbf{X}_h^\top \mathbf{X}_h$  cannot be inverted.

# Computation

First, calculate the cluster-level matrices and vectors

$$\mathbf{X}_j^\top \mathbf{X}_j \text{ and } \mathbf{X}_j^\top \mathbf{y}_j, \quad j = 1, \dots, J, \text{ for } \{j, J\} = \{g, G\}, \{h, H\}, \{i, I\}. \quad (17)$$

The ones for the intersections can be computed in a single pass over the  $N$  observations. The others are just sums of some of them.

With two-way fixed effects in the  $G$  and  $H$  dimensions,

$$\mathbf{y} = \mathbf{Z}\boldsymbol{\beta}_p + \mathbf{D}^G\boldsymbol{\gamma} + \mathbf{D}^H\boldsymbol{\delta} + \mathbf{u}. \quad (18)$$

Now  $\mathbf{X} = [\mathbf{Z} \ \mathbf{D}^G \ \mathbf{D}^H]$ , and  $k = p + G + H - 1$ , and the matrices  $\mathbf{X}^\top \mathbf{X} - \mathbf{X}_g^\top \mathbf{X}_g$  and  $\mathbf{X}^\top \mathbf{X} - \mathbf{X}_h^\top \mathbf{X}_h$  cannot be inverted.

Simplest approach is to replace the inverse in (14) by a generalized inverse. Then  $\hat{\mathbf{V}}_{JK}^{(3)}$  in (16) can only be calculated as a  $p \times p$  matrix.

# Computation

First, calculate the cluster-level matrices and vectors

$$\mathbf{X}_j^\top \mathbf{X}_j \text{ and } \mathbf{X}_j^\top \mathbf{y}_j, \quad j = 1, \dots, J, \text{ for } \{j, J\} = \{g, G\}, \{h, H\}, \{i, I\}. \quad (17)$$

The ones for the intersections can be computed in a single pass over the  $N$  observations. The others are just sums of some of them.

With two-way fixed effects in the  $G$  and  $H$  dimensions,

$$\mathbf{y} = \mathbf{Z}\boldsymbol{\beta}_p + \mathbf{D}^G\boldsymbol{\gamma} + \mathbf{D}^H\boldsymbol{\delta} + \mathbf{u}. \quad (18)$$

Now  $\mathbf{X} = [\mathbf{Z} \ \mathbf{D}^G \ \mathbf{D}^H]$ , and  $k = p + G + H - 1$ , and the matrices  $\mathbf{X}^\top \mathbf{X} - \mathbf{X}_g^\top \mathbf{X}_g$  and  $\mathbf{X}^\top \mathbf{X} - \mathbf{X}_h^\top \mathbf{X}_h$  cannot be inverted.

Simplest approach is to replace the inverse in (14) by a generalized inverse. Then  $\hat{\mathbf{V}}_{JK}^{(3)}$  in (16) can only be calculated as a  $p \times p$  matrix.

Computing  $\text{CV}_3^{(3)}$  and friends for (18) can be costly when  $G$  and  $H$  are not fairly small.

# Critical Values

# Critical Values

Student's  $t$  distribution with  $\min\{G, H\} - 1$  degrees of freedom is normally used for  $t$ -statistics based on  $CV_1$ .

# Critical Values

Student's  $t$  distribution with  $\min\{G, H\} - 1$  degrees of freedom is normally used for  $t$ -statistics based on  $CV_1$ .

We do this for  $t$ -statistics based on  $CV_3$  as well.

# Critical Values

Student's  $t$  distribution with  $\min\{G, H\} - 1$  degrees of freedom is normally used for  $t$ -statistics based on  $CV_1$ .

We do this for  $t$ -statistics based on  $CV_3$  as well.

There are at least two possible alternatives:

# Critical Values

Student's  $t$  distribution with  $\min\{G, H\} - 1$  degrees of freedom is normally used for  $t$ -statistics based on  $CV_1$ .

We do this for  $t$ -statistics based on  $CV_3$  as well.

There are at least two possible alternatives:

- 1 Employ a  $t$ -distribution with an estimated degrees-of-freedom parameter, and maybe an estimated scale parameter, as in Bell and McCaffrey (2002), Imbens and Kolesár (2016), and Hansen (2023).

# Critical Values

Student's  $t$  distribution with  $\min\{G, H\} - 1$  degrees of freedom is normally used for  $t$ -statistics based on  $CV_1$ .

We do this for  $t$ -statistics based on  $CV_3$  as well.

There are at least two possible alternatives:

- 1 Employ a  $t$ -distribution with an estimated degrees-of-freedom parameter, and maybe an estimated scale parameter, as in Bell and McCaffrey (2002), Imbens and Kolesár (2016), and Hansen (2023).
- 2 Use the wild cluster bootstrap (MacKinnon, Nielsen, and Webb 2021) or the pigeonhole bootstrap (Owen, 2007).

# Critical Values

Student's  $t$  distribution with  $\min\{G, H\} - 1$  degrees of freedom is normally used for  $t$ -statistics based on  $CV_1$ .

We do this for  $t$ -statistics based on  $CV_3$  as well.

There are at least two possible alternatives:

- 1 Employ a  $t$ -distribution with an estimated degrees-of-freedom parameter, and maybe an estimated scale parameter, as in Bell and McCaffrey (2002), Imbens and Kolesár (2016), and Hansen (2023).
- 2 Use the wild cluster bootstrap (MacKinnon, Nielsen, and Webb 2021) or the pigeonhole bootstrap (Owen, 2007).
  - For two-way clustering, no version of the wild cluster bootstrap can replicate the intra-cluster covariances in the residuals.

# Critical Values

Student's  $t$  distribution with  $\min\{G, H\} - 1$  degrees of freedom is normally used for  $t$ -statistics based on  $CV_1$ .

We do this for  $t$ -statistics based on  $CV_3$  as well.

There are at least two possible alternatives:

- 1 Employ a  $t$ -distribution with an estimated degrees-of-freedom parameter, and maybe an estimated scale parameter, as in Bell and McCaffrey (2002), Imbens and Kolesár (2016), and Hansen (2023).
- 2 Use the wild cluster bootstrap (MacKinnon, Nielsen, and Webb 2021) or the pigeonhole bootstrap (Owen, 2007).
  - For two-way clustering, no version of the wild cluster bootstrap can replicate the intra-cluster covariances in the residuals.
  - The pigeonhole bootstrap is an ingenious generalization of the ordinary pairs (resampling) bootstrap for one-way clustering.

# Critical Values

Student's  $t$  distribution with  $\min\{G, H\} - 1$  degrees of freedom is normally used for  $t$ -statistics based on  $CV_1$ .

We do this for  $t$ -statistics based on  $CV_3$  as well.

There are at least two possible alternatives:

- 1 Employ a  $t$ -distribution with an estimated degrees-of-freedom parameter, and maybe an estimated scale parameter, as in Bell and McCaffrey (2002), Imbens and Kolesár (2016), and Hansen (2023).
- 2 Use the wild cluster bootstrap (MacKinnon, Nielsen, and Webb 2021) or the pigeonhole bootstrap (Owen, 2007).
  - For two-way clustering, no version of the wild cluster bootstrap can replicate the intra-cluster covariances in the residuals.
  - The pigeonhole bootstrap is an ingenious generalization of the ordinary pairs (resampling) bootstrap for one-way clustering.
  - But pairs bootstrap typically performs worse than WCR bootstrap (MacKinnon and Webb, TPM 2017; MacKinnon, 2023).

# Consistency of the Cluster-Jackknife CRVE

# Consistency of the Cluster-Jackknife CRVE

## Theorem 1.

Let  $\hat{V}_3$  denote any of the three jackknife CRVEs —  $CV_3^{(2)}$ ,  $CV_3^{(3)}$ , and  $CV_3^{(\max)}$  — and let  $\text{Var}(\hat{\beta})$  be given in (5).

# Consistency of the Cluster-Jackknife CRVE

## Theorem 1.

Let  $\hat{V}_3$  denote any of the three jackknife CRVEs —  $CV_3^{(2)}$ ,  $CV_3^{(3)}$ , and  $CV_3^{(\max)}$  — and let  $\text{Var}(\hat{\beta})$  be given in (5).

Then, under suitable assumptions,  $(\text{Var}(\hat{\beta}))^{-1} \hat{V}_3 \xrightarrow{P} 1$ .

# Consistency of the Cluster-Jackknife CRVE

## Theorem 1.

Let  $\hat{V}_3$  denote any of the three jackknife CRVEs —  $CV_3^{(2)}$ ,  $CV_3^{(3)}$ , and  $CV_3^{(\max)}$  — and let  $\text{Var}(\hat{\beta})$  be given in (5).

Then, under suitable assumptions,  $(\text{Var}(\hat{\beta}))^{-1} \hat{V}_3 \xrightarrow{P} 1$ .

It follows that

$$(\text{Var}(\hat{\beta}))^{-1/2} (\hat{\beta} - \beta_0) \xrightarrow{d} \mathbf{N}(\mathbf{0}, \mathbf{I}_k). \quad (19)$$

The assumptions and proof follow [Yap \(2025\)](#).

# Consistency of the Cluster-Jackknife CRVE

## Theorem 1.

Let  $\hat{V}_3$  denote any of the three jackknife CRVEs —  $CV_3^{(2)}$ ,  $CV_3^{(3)}$ , and  $CV_3^{(\max)}$  — and let  $\text{Var}(\hat{\beta})$  be given in (5).

Then, under suitable assumptions,  $(\text{Var}(\hat{\beta}))^{-1} \hat{V}_3 \xrightarrow{P} 1$ .

It follows that

$$(\text{Var}(\hat{\beta}))^{-1/2}(\hat{\beta} - \beta_0) \xrightarrow{d} \mathbf{N}(\mathbf{0}, \mathbf{I}_k). \quad (19)$$

The assumptions and proof follow [Yap \(2025\)](#).

He demonstrates the consistency of  $\hat{V}_1$  without assuming that the disturbances are generated as separately exchangeable arrays, as [Davezies, D'Haultfoeuille, and Guyonvarch \(2018,2025\)](#) do.

# Consistency of the Cluster-Jackknife CRVE

## Theorem 1.

Let  $\hat{V}_3$  denote any of the three jackknife CRVEs —  $CV_3^{(2)}$ ,  $CV_3^{(3)}$ , and  $CV_3^{(\max)}$  — and let  $\text{Var}(\hat{\beta})$  be given in (5).

Then, under suitable assumptions,  $(\text{Var}(\hat{\beta}))^{-1} \hat{V}_3 \xrightarrow{P} 1$ .

It follows that

$$(\text{Var}(\hat{\beta}))^{-1/2}(\hat{\beta} - \beta_0) \xrightarrow{d} \mathbf{N}(\mathbf{0}, \mathbf{I}_k). \quad (19)$$

The assumptions and proof follow [Yap \(2025\)](#).

He demonstrates the consistency of  $\hat{V}_1$  without assuming that the disturbances are generated as separately exchangeable arrays, as [Davezies, D'Haultfoeuille, and Guyonvarch \(2018,2025\)](#) do.

Yap's assumptions are weaker, and his method of proof is simpler.

# Simulation Experiments

# Simulation Experiments

The disturbances are generated so that cluster fixed effects do not eliminate intra-cluster correlation. We use factor models of the form

$$\begin{aligned} z_{ghi} &= \sigma_g \tilde{\zeta}_g^1 + \sigma_h \tilde{\zeta}_h^1 + \sigma_\epsilon \zeta_{ghi} & \text{if } i \text{ is odd,} \\ z_{ghi} &= \sigma_g \tilde{\zeta}_g^2 + \sigma_h \tilde{\zeta}_h^2 + \sigma_\epsilon \zeta_{ghi} & \text{if } i \text{ is even.} \end{aligned} \tag{20}$$

# Simulation Experiments

The disturbances are generated so that cluster fixed effects do not eliminate intra-cluster correlation. We use factor models of the form

$$\begin{aligned} z_{ghi} &= \sigma_g \zeta_g^1 + \sigma_h \zeta_h^1 + \sigma_\epsilon \zeta_{ghi} & \text{if } i \text{ is odd,} \\ z_{ghi} &= \sigma_g \zeta_g^2 + \sigma_h \zeta_h^2 + \sigma_\epsilon \zeta_{ghi} & \text{if } i \text{ is even.} \end{aligned} \tag{20}$$

- $\zeta_g^1$  and  $\zeta_g^2$  are  $N(0, 1)$  random effects which apply to the odd-numbered and even-numbered observations within the  $g^{\text{th}}$  cluster.

# Simulation Experiments

The disturbances are generated so that cluster fixed effects do not eliminate intra-cluster correlation. We use factor models of the form

$$\begin{aligned} z_{ghi} &= \sigma_g \zeta_g^1 + \sigma_h \zeta_h^1 + \sigma_\epsilon \zeta_{ghi} & \text{if } i \text{ is odd,} \\ z_{ghi} &= \sigma_g \zeta_g^2 + \sigma_h \zeta_h^2 + \sigma_\epsilon \zeta_{ghi} & \text{if } i \text{ is even.} \end{aligned} \quad (20)$$

- $\zeta_g^1$  and  $\zeta_g^2$  are  $N(0, 1)$  random effects which apply to the odd-numbered and even-numbered observations within the  $g^{\text{th}}$  cluster.
- $\zeta_h^1$  and  $\zeta_h^2$  are  $N(0, 1)$  random effects which apply to the odd-numbered and even-numbered observations within the  $h^{\text{th}}$  cluster. The  $\zeta_{ghi}$  are independent standard normals.

# Simulation Experiments

The disturbances are generated so that cluster fixed effects do not eliminate intra-cluster correlation. We use factor models of the form

$$\begin{aligned} z_{ghi} &= \sigma_g \zeta_g^1 + \sigma_h \zeta_h^1 + \sigma_\epsilon \zeta_{ghi} & \text{if } i \text{ is odd,} \\ z_{ghi} &= \sigma_g \zeta_g^2 + \sigma_h \zeta_h^2 + \sigma_\epsilon \zeta_{ghi} & \text{if } i \text{ is even.} \end{aligned} \quad (20)$$

- $\zeta_g^1$  and  $\zeta_g^2$  are  $N(0, 1)$  random effects which apply to the odd-numbered and even-numbered observations within the  $g^{\text{th}}$  cluster.
- $\zeta_h^1$  and  $\zeta_h^2$  are  $N(0, 1)$  random effects which apply to the odd-numbered and even-numbered observations within the  $h^{\text{th}}$  cluster. The  $\zeta_{ghi}$  are independent standard normals.
- $\sigma_g$  and  $\sigma_h$  are specified as functions of correlations  $\rho_g$  and  $\rho_h$ , with  $\sigma_j = (\rho_j / (1 - \rho_j))^{1/2}$  for  $j = g, h$ .

# Simulation Experiments

The disturbances are generated so that cluster fixed effects do not eliminate intra-cluster correlation. We use factor models of the form

$$\begin{aligned} z_{ghi} &= \sigma_g \zeta_g^1 + \sigma_h \zeta_h^1 + \sigma_\epsilon \zeta_{ghi} & \text{if } i \text{ is odd,} \\ z_{ghi} &= \sigma_g \zeta_g^2 + \sigma_h \zeta_h^2 + \sigma_\epsilon \zeta_{ghi} & \text{if } i \text{ is even.} \end{aligned} \quad (20)$$

- $\zeta_g^1$  and  $\zeta_g^2$  are  $N(0, 1)$  random effects which apply to the odd-numbered and even-numbered observations within the  $g^{\text{th}}$  cluster.
- $\zeta_h^1$  and  $\zeta_h^2$  are  $N(0, 1)$  random effects which apply to the odd-numbered and even-numbered observations within the  $h^{\text{th}}$  cluster. The  $\zeta_{ghi}$  are independent standard normals.
- $\sigma_g$  and  $\sigma_h$  are specified as functions of correlations  $\rho_g$  and  $\rho_h$ , with  $\sigma_j = (\rho_j / (1 - \rho_j))^{1/2}$  for  $j = g, h$ .
- The value of  $\sigma_\epsilon$  is  $(1 - \sigma_g^2 - \sigma_h^2)^{1/2}$ , so that  $\text{Var}(z_{ghi}) = 1$ .

The cluster sizes in the  $G$  dimension are given by

$$N_g = \left[ N \frac{\exp(\gamma g/G)}{\sum_{j=1}^G \exp(\gamma j/G)} \right], \quad g = 1, \dots, G-1, \quad (21)$$

where  $[x]$  denotes the integer part of  $x$ . Then  $N_G = N - \sum_{g=1}^{G-1} N_g$ .

The cluster sizes in the  $G$  dimension are given by

$$N_g = \left[ N \frac{\exp(\gamma g/G)}{\sum_{j=1}^G \exp(\gamma j/G)} \right], \quad g = 1, \dots, G-1, \quad (21)$$

where  $[x]$  denotes the integer part of  $x$ . Then  $N_G = N - \sum_{g=1}^{G-1} N_g$ .  
(21), perhaps with a different  $\gamma$ , is also used in the  $H$  dimension.

The cluster sizes in the  $G$  dimension are given by

$$N_g = \left[ N \frac{\exp(\gamma g/G)}{\sum_{j=1}^G \exp(\gamma j/G)} \right], \quad g = 1, \dots, G-1, \quad (21)$$

where  $[x]$  denotes the integer part of  $x$ . Then  $N_G = N - \sum_{g=1}^{G-1} N_g$ .

(21), perhaps with a different  $\gamma$ , is also used in the  $H$  dimension.

Assuming that the distributions are independent,  $N_{gh} \approx N_g N_h / N$ .

The cluster sizes in the  $G$  dimension are given by

$$N_g = \left\lfloor N \frac{\exp(\gamma g/G)}{\sum_{j=1}^G \exp(\gamma j/G)} \right\rfloor, \quad g = 1, \dots, G-1, \quad (21)$$

where  $\lfloor x \rfloor$  denotes the integer part of  $x$ . Then  $N_G = N - \sum_{g=1}^{G-1} N_g$ .

(21), perhaps with a different  $\gamma$ , is also used in the  $H$  dimension.

Assuming that the distributions are independent,  $N_{gh} \approx N_g N_h / N$ .

In a final step, the cluster sizes are adjusted to ensure that they are all integers with  $N = \sum_{g=1}^G N_g = \sum_{h=1}^H N_h = \sum_{g=1}^G \sum_{h=1}^H N_{gh}$ .

The cluster sizes in the  $G$  dimension are given by

$$N_g = \left\lfloor N \frac{\exp(\gamma g/G)}{\sum_{j=1}^G \exp(\gamma j/G)} \right\rfloor, \quad g = 1, \dots, G-1, \quad (21)$$

where  $\lfloor x \rfloor$  denotes the integer part of  $x$ . Then  $N_G = N - \sum_{g=1}^{G-1} N_g$ .

(21), perhaps with a different  $\gamma$ , is also used in the  $H$  dimension.

Assuming that the distributions are independent,  $N_{gh} \approx N_g N_h / N$ .

In a final step, the cluster sizes are adjusted to ensure that they are all integers with  $N = \sum_{g=1}^G N_g = \sum_{h=1}^H N_h = \sum_{g=1}^G \sum_{h=1}^H N_{gh}$ .

The experiments currently involve **normally distributed regressors**, which follow the factor model (20).

The cluster sizes in the  $G$  dimension are given by

$$N_g = \left\lfloor N \frac{\exp(\gamma g/G)}{\sum_{j=1}^G \exp(\gamma j/G)} \right\rfloor, \quad g = 1, \dots, G-1, \quad (21)$$

where  $\lfloor x \rfloor$  denotes the integer part of  $x$ . Then  $N_G = N - \sum_{g=1}^{G-1} N_g$ .

(21), perhaps with a different  $\gamma$ , is also used in the  $H$  dimension.

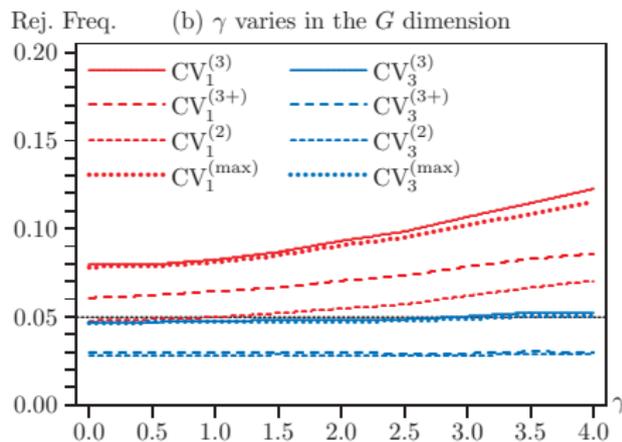
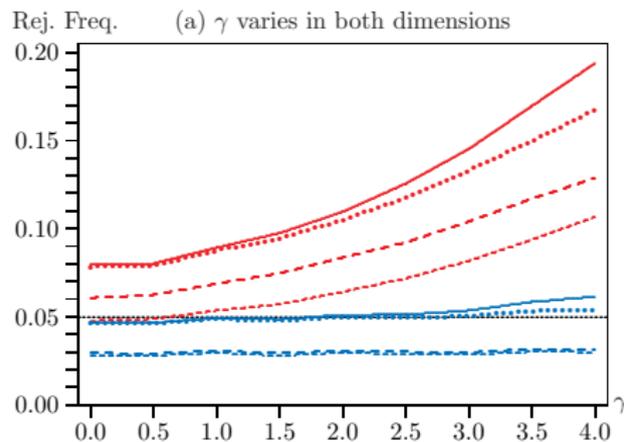
Assuming that the distributions are independent,  $N_{gh} \approx N_g N_h / N$ .

In a final step, the cluster sizes are adjusted to ensure that they are all integers with  $N = \sum_{g=1}^G N_g = \sum_{h=1}^H N_h = \sum_{g=1}^G \sum_{h=1}^H N_{gh}$ .

The experiments currently involve **normally distributed regressors**, which follow the factor model (20).

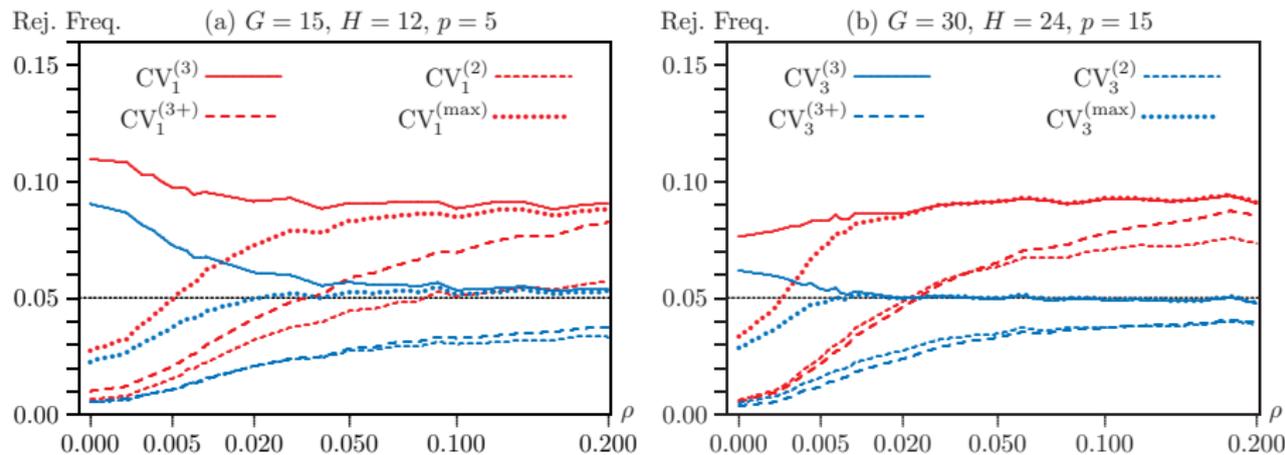
In most experiments, we set  $\rho_g^x = \rho_h^x = 0.2$  for the regressors and  $\rho_g = \rho_h = 0.1$  for the disturbances.

Figure 1. Rejection frequencies as functions of cluster size variation



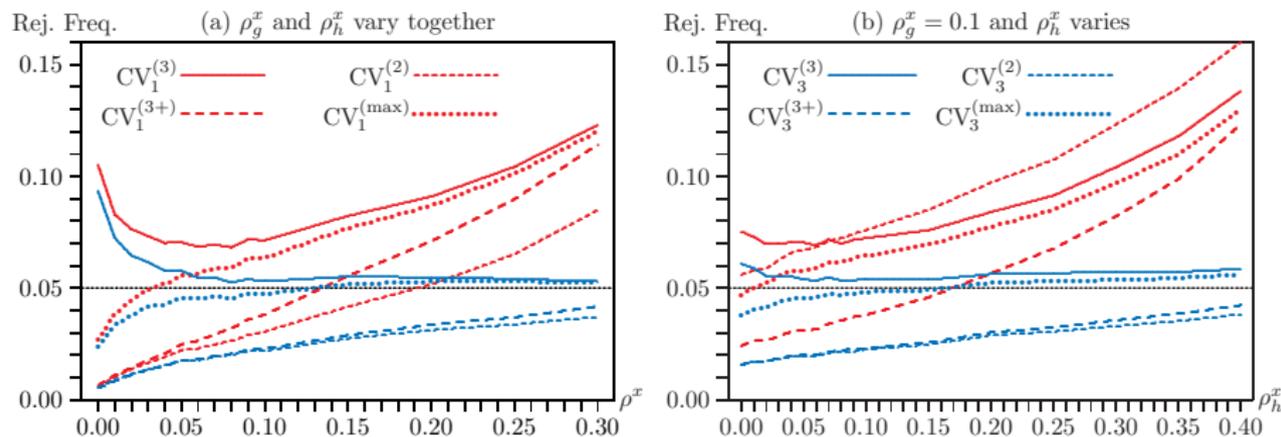
- $N = 10,000, G = 15, H = 12, I = 180, p = 10, k = 36$ .
- Regressors are from factor model (20), with  $\rho_g^x = \rho_h^x = 0.2$ .
- Disturbances are from factor model (20), with  $\rho_g = \rho_h = 0.1$ .
- Results are based on 100,000 replications.

Figure 2. Rejection frequencies as functions of disturbance correlations



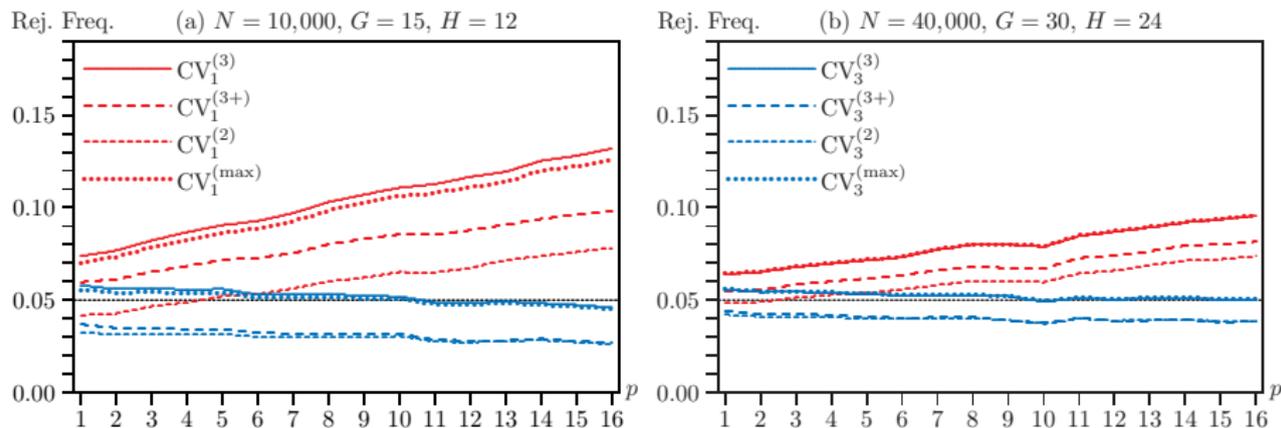
- (a)  $N = 10,000, G = 15, H = 12, I = 180, p = 5, \gamma = 2$ .
- (b)  $N = 40,000, G = 30, H = 24, I = 720, p = 15, \gamma = 2$ .
- Regressors are from factor model (20), with  $\rho_g^x = \rho_h^x = 0.2$ .
- Disturbances are from factor model (20), with  $\rho_g = \rho_h$  that vary.
- Results are based on 100,000 replications.

Figure 3. Rejection frequencies as functions of regressor correlations



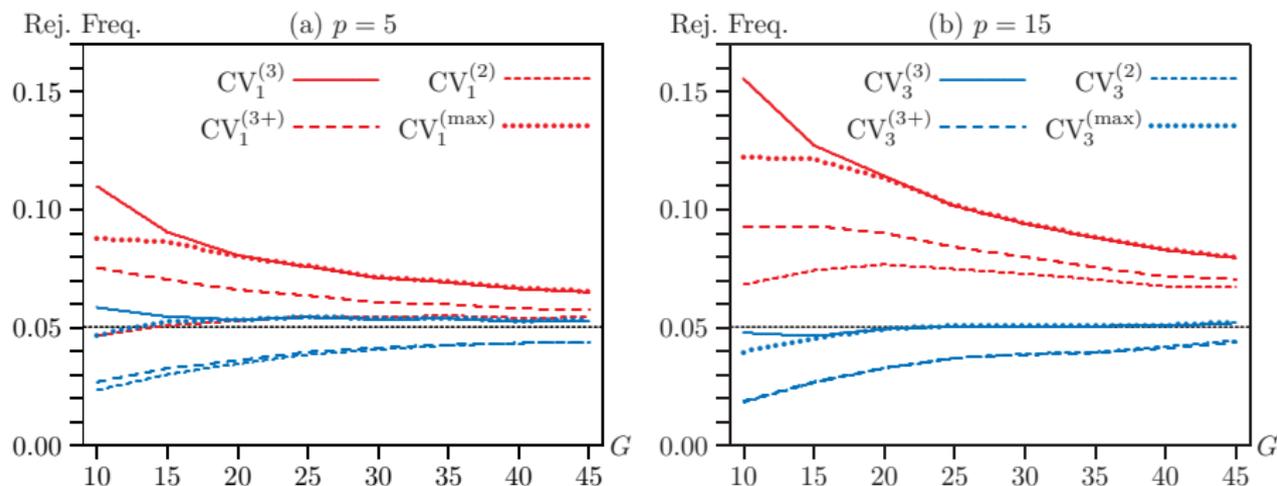
- $N = 10,000, G = 15, H = 12, I = 180, p = 5, \gamma = 2$ .
- Disturbances are from factor model (20), with  $\rho_g = \rho_h = 0.1$ .
- Regressors are from factor model; one or both values of  $\rho^x$  vary.
- Results are based on 100,000 replications.

Figure 4. Rejection frequencies as functions of number of regressors



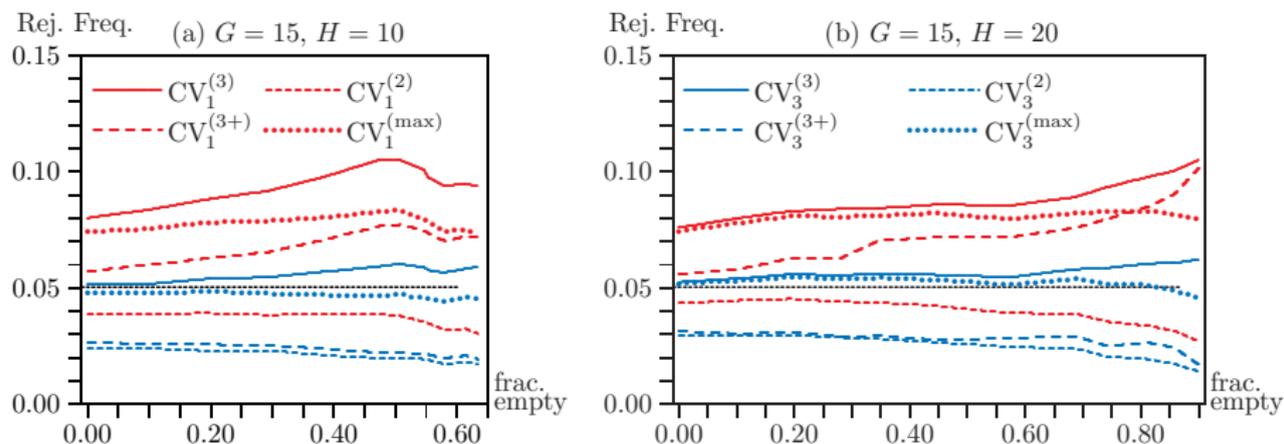
- (a)  $N = 10,000, G = 15, H = 12, I = 180, \gamma = 2$ .
- (b)  $N = 40,000, G = 30, H = 24, I = 720, \gamma = 2$ .
- Regressors are from factor model (20), with  $\rho_g^x = \rho_h^x = 0.2$ .
- Disturbances are from factor model (20), with  $\rho_g = \rho_h = 0.1$ .

Figure 5. Rejection frequencies as functions of numbers of clusters



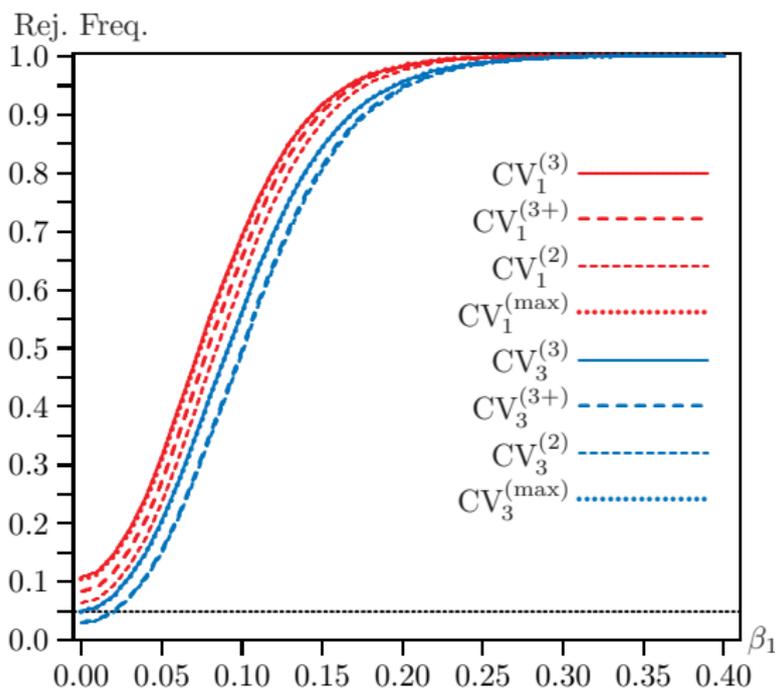
- The value of  $G$  varies from 5 to 45 by 5, with  $H = 4G/5$ .
- The value of  $N$  varies from 1,111 to 90,000.
- Regressors are from factor model (20), with  $\rho_g^x = \rho_h^x = 0.2$ .
- Disturbances are from factor model (20), with  $\rho_g = \rho_h = 0.1$ .

Figure 6. Rejection frequencies as functions of fraction of empty intersections



- $N = 6000$  in Panel (a) and  $N = 12000$  in Panel (b).
- The first 5 regressors are from the model (20), with  $\rho_g^x = \rho_h^x = 0.2$ .
- The extra 5 regressors are binary and equal 1 with probability 0.25.
- Disturbances are from factor model (20), with  $\rho_g = \rho_h = 0.1$ .

Figure 7. Power functions for eight tests



- $N = 10,000, G = 15, H = 12, I = 180, p = 5, \gamma = 2$ .
- Regressors and disturbances are from factor model (20).
- All coefficients except  $\beta_1$  equal 0.

# Empirical Examples

# Empirical Examples

The first example is based on Alsan (2015), which studies the impact of the tsetse fly on economic development in Africa.

# Empirical Examples

The first example is based on Alsan (2015), which studies the impact of the tsetse fly on economic development in Africa.

- The main explanatory variable is the TSI or “Tsetse Suitability Index.” It measures how suitable an area is to support the tsetse fly, which carries a parasite that affects humans and livestock.

# Empirical Examples

The first example is based on Alsan (2015), which studies the impact of the tsetse fly on economic development in Africa.

- The main explanatory variable is the TSI or “Tsetse Suitability Index.” It measures how suitable an area is to support the tsetse fly, which carries a parasite that affects humans and livestock.
- The paper tests the extent to which the tsetse fly inhibited political and agricultural development in parts of Africa.

# Empirical Examples

The first example is based on Alsan (2015), which studies the impact of the tsetse fly on economic development in Africa.

- The main explanatory variable is the TSI or “Tsetse Suitability Index.” It measures how suitable an area is to support the tsetse fly, which carries a parasite that affects humans and livestock.
- The paper tests the extent to which the tsetse fly inhibited political and agricultural development in parts of Africa.
- There are seven different outcome variables. These are regressed on TSI and various controls.

# Empirical Examples

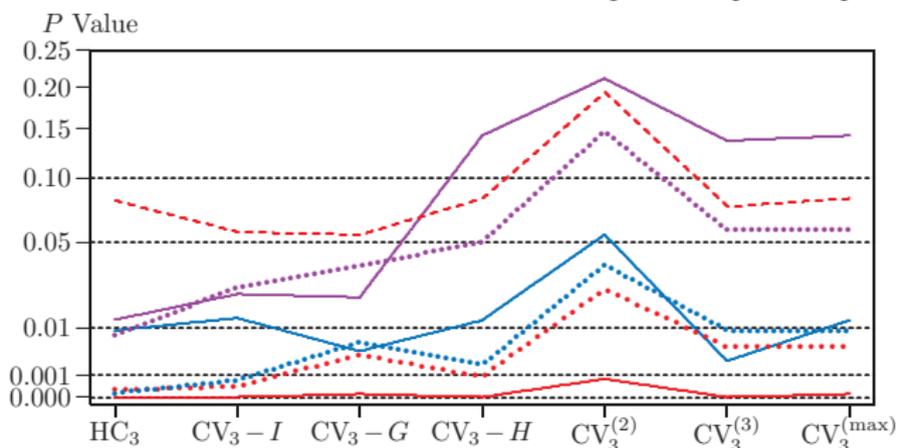
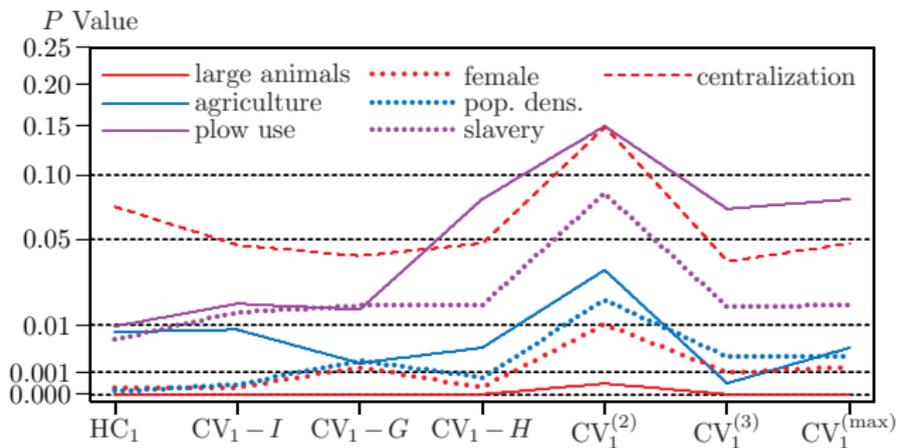
The first example is based on Alsan (2015), which studies the impact of the tsetse fly on economic development in Africa.

- The main explanatory variable is the TSI or “Tsetse Suitability Index.” It measures how suitable an area is to support the tsetse fly, which carries a parasite that affects humans and livestock.
- The paper tests the extent to which the tsetse fly inhibited political and agricultural development in parts of Africa.
- There are seven different outcome variables. These are regressed on TSI and various controls.
- Sample sizes vary from 315 to 485. There are two clustering dimensions, country and “cultural province.” Most results use one-way clustering by the latter.

# Empirical Examples

The first example is based on Alsan (2015), which studies the impact of the tsetse fly on economic development in Africa.

- The main explanatory variable is the TSI or “Tsetse Suitability Index.” It measures how suitable an area is to support the tsetse fly, which carries a parasite that affects humans and livestock.
- The paper tests the extent to which the tsetse fly inhibited political and agricultural development in parts of Africa.
- There are seven different outcome variables. These are regressed on TSI and various controls.
- Sample sizes vary from 315 to 485. There are two clustering dimensions, country and “cultural province.” Most results use one-way clustering by the latter.
- There are 44 countries, 43 or 44 provinces, and between 112 and 142 non-empty intersections. Since  $44^2 = 1936$ ,  $I \ll GH$ .



The second example studies the effects of minimum wages on earnings for young people in Canada, using Labour Force Survey data.

The second example studies the effects of minimum wages on earnings for young people in Canada, using Labour Force Survey data.

- Individuals aged 18–24 who have been in Canada less than ten years. 28,599 observations in 10 provinces for 2008 to 2019.

The second example studies the effects of minimum wages on earnings for young people in Canada, using Labour Force Survey data.

- Individuals aged 18–24 who have been in Canada less than ten years. 28,599 observations in 10 provinces for 2008 to 2019.
- The dependent variable is the log of weekly earnings, and the regressor of interest is the log of the minimum wage, which varies by province. There are 63 distinct values.

The second example studies the effects of minimum wages on earnings for young people in Canada, using Labour Force Survey data.

- Individuals aged 18–24 who have been in Canada less than ten years. 28,599 observations in 10 provinces for 2008 to 2019.
- The dependent variable is the log of weekly earnings, and the regressor of interest is the log of the minimum wage, which varies by province. There are 63 distinct values.

$$\begin{aligned} \log \text{earn}_{ipmt} = & \alpha + \beta \log \text{mw}_{pmt} + \gamma \text{big city}_{ipmt} + \delta \text{older}_{ipmt} \\ & + \text{year}_t + \text{month}_m + \text{prov}_p + \epsilon_{ipmt} \end{aligned}$$

The second example studies the effects of minimum wages on earnings for young people in Canada, using Labour Force Survey data.

- Individuals aged 18–24 who have been in Canada less than ten years. 28,599 observations in 10 provinces for 2008 to 2019.
- The dependent variable is the log of weekly earnings, and the regressor of interest is the log of the minimum wage, which varies by province. There are 63 distinct values.

$$\begin{aligned} \log \text{earn}_{ipmt} = & \alpha + \beta \log \text{mw}_{pmt} + \gamma \text{big city}_{ipmt} + \delta \text{older}_{ipmt} \\ & + \text{year}_t + \text{month}_m + \text{prov}_p + \epsilon_{ipmt} \end{aligned}$$

- Observations per province ( $H = 10$ ) vary from 163 to 6554.

The second example studies the effects of minimum wages on earnings for young people in Canada, using Labour Force Survey data.

- Individuals aged 18–24 who have been in Canada less than ten years. 28,599 observations in 10 provinces for 2008 to 2019.
- The dependent variable is the log of weekly earnings, and the regressor of interest is the log of the minimum wage, which varies by province. There are 63 distinct values.

$$\begin{aligned} \log \text{earn}_{ipmt} = & \alpha + \beta \log \text{mw}_{pmt} + \gamma \text{big city}_{ipmt} + \delta \text{older}_{ipmt} \\ & + \text{year}_t + \text{month}_m + \text{prov}_p + \epsilon_{ipmt} \end{aligned}$$

- Observations per province ( $H = 10$ ) vary from 163 to 6554.
- Observations per year ( $G = 12$ ) vary from 2051 to 2723.

The second example studies the effects of minimum wages on earnings for young people in Canada, using Labour Force Survey data.

- Individuals aged 18–24 who have been in Canada less than ten years. 28,599 observations in 10 provinces for 2008 to 2019.
- The dependent variable is the log of weekly earnings, and the regressor of interest is the log of the minimum wage, which varies by province. There are 63 distinct values.

$$\begin{aligned} \log \text{earn}_{ipmt} = & \alpha + \beta \log \text{mw}_{pmt} + \gamma \text{big city}_{ipmt} + \delta \text{older}_{ipmt} \\ & + \text{year}_t + \text{month}_m + \text{prov}_p + \epsilon_{ipmt} \end{aligned}$$

- Observations per province ( $H = 10$ ) vary from 163 to 6554.
- Observations per year ( $G = 12$ ) vary from 2051 to 2723.
- Observations per intersection ( $I = 120$ ) vary from 3 to 710.

The second example studies the effects of minimum wages on earnings for young people in Canada, using Labour Force Survey data.

- Individuals aged 18–24 who have been in Canada less than ten years. 28,599 observations in 10 provinces for 2008 to 2019.
- The dependent variable is the log of weekly earnings, and the regressor of interest is the log of the minimum wage, which varies by province. There are 63 distinct values.

$$\begin{aligned} \log \text{earn}_{ipmt} = & \alpha + \beta \log \text{mw}_{pmt} + \gamma \text{big city}_{ipmt} + \delta \text{older}_{ipmt} \\ & + \text{year}_t + \text{month}_m + \text{prov}_p + \epsilon_{ipmt} \end{aligned}$$

- Observations per province ( $H = 10$ ) vary from 163 to 6554.
- Observations per year ( $G = 12$ ) vary from 2051 to 2723.
- Observations per intersection ( $I = 120$ ) vary from 3 to 710.
- Coefficient estimate is 0.2934, with standard errors between 0.0254 for  $\text{HC}_1$  and 0.1663 for  $\text{CV}_3^{(2)}$ .

We also run **placebo regressions**, which generalize the idea of “placebo laws” proposed in Bertrand, Duflo, and Mullainathan (2004).

We also run **placebo regressions**, which generalize the idea of “placebo laws” proposed in Bertrand, Duflo, and Mullainathan (2004).

- The idea is to add a randomly generated regressor that looks similar to the minimum wage to the actual regression.

We also run **placebo regressions**, which generalize the idea of “placebo laws” proposed in Bertrand, Duflo, and Mullainathan (2004).

- The idea is to add a randomly generated regressor that looks similar to the minimum wage to the actual regression.
- This is done many times (100,000 in our simulations). Placebo regressor changes across simulations, but not the regressand.

We also run **placebo regressions**, which generalize the idea of “placebo laws” proposed in Bertrand, Duflo, and Mullainathan (2004).

- The idea is to add a randomly generated regressor that looks similar to the minimum wage to the actual regression.
- This is done many times (100,000 in our simulations). Placebo regressor changes across simulations, but not the regressand.
- Since the placebo regressors are random, they should have no explanatory power if the regression is specified correctly.

We also run **placebo regressions**, which generalize the idea of “placebo laws” proposed in Bertrand, Duflo, and Mullainathan (2004).

- The idea is to add a randomly generated regressor that looks similar to the minimum wage to the actual regression.
- This is done many times (100,000 in our simulations). Placebo regressor changes across simulations, but not the regressand.
- Since the placebo regressors are random, they should have no explanatory power if the regression is specified correctly.
- If a test at the .05 level rejects much more or less often than 5%, then we should not trust results of that test.

We also run **placebo regressions**, which generalize the idea of “placebo laws” proposed in Bertrand, Duflo, and Mullainathan (2004).

- The idea is to add a randomly generated regressor that looks similar to the minimum wage to the actual regression.
- This is done many times (100,000 in our simulations). Placebo regressor changes across simulations, but not the regressand.
- Since the placebo regressors are random, they should have no explanatory power if the regression is specified correctly.
- If a test at the .05 level rejects much more or less often than 5%, then we should not trust results of that test.
- We model the placebo minimum wage as a two-stage process at the province-year level.

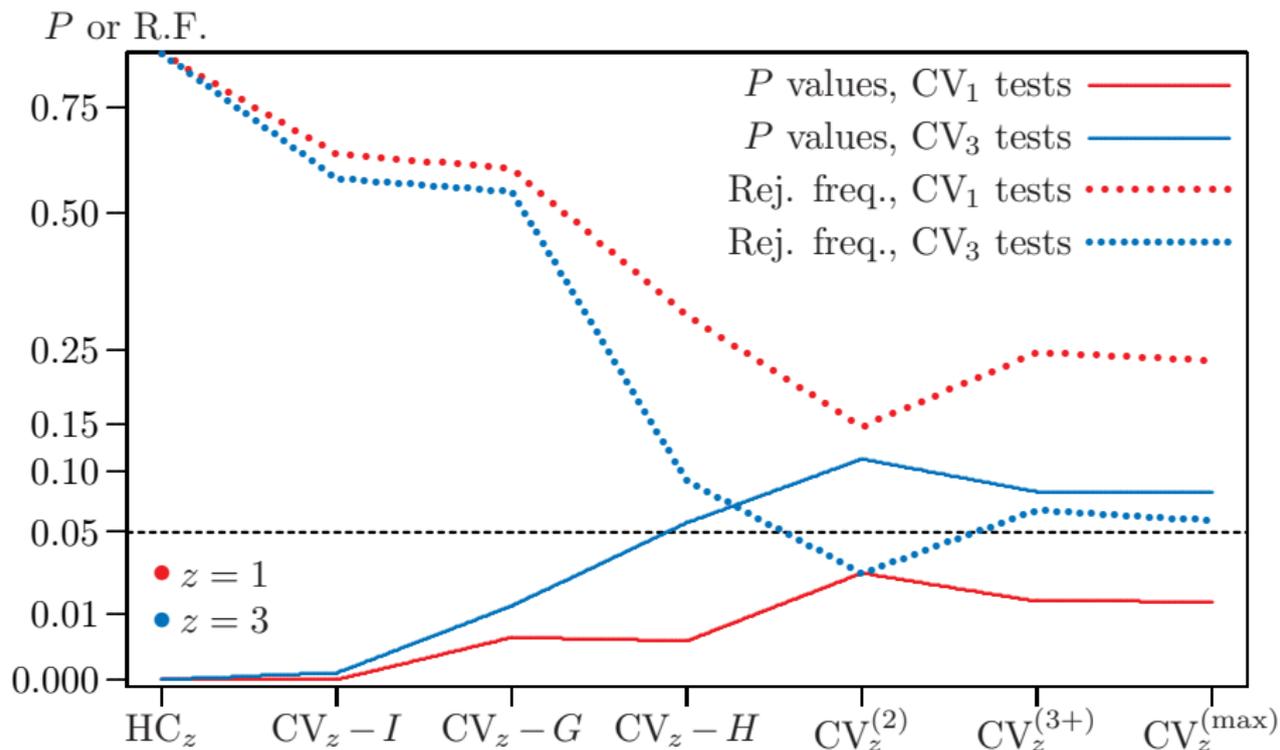
We also run **placebo regressions**, which generalize the idea of “placebo laws” proposed in Bertrand, Duflo, and Mullainathan (2004).

- The idea is to add a randomly generated regressor that looks similar to the minimum wage to the actual regression.
- This is done many times (100,000 in our simulations). Placebo regressor changes across simulations, but not the regressand.
- Since the placebo regressors are random, they should have no explanatory power if the regression is specified correctly.
- If a test at the .05 level rejects much more or less often than 5%, then we should not trust results of that test.
- We model the placebo minimum wage as a two-stage process at the province-year level.
- It increases if random variables  $\nu_p$  and  $\nu_y$  both exceed threshold values. The increase is a random amount of 0.25, 0.50, 0.75, or 1.00.

We also run **placebo regressions**, which generalize the idea of “placebo laws” proposed in Bertrand, Duflo, and Mullainathan (2004).

- The idea is to add a randomly generated regressor that looks similar to the minimum wage to the actual regression.
- This is done many times (100,000 in our simulations). Placebo regressor changes across simulations, but not the regressand.
- Since the placebo regressors are random, they should have no explanatory power if the regression is specified correctly.
- If a test at the .05 level rejects much more or less often than 5%, then we should not trust results of that test.
- We model the placebo minimum wage as a two-stage process at the province-year level.
- It increases if random variables  $\nu_p$  and  $\nu_y$  both exceed threshold values. The increase is a random amount of 0.25, 0.50, 0.75, or 1.00.

The details matter, and current results are preliminary.

Figure 9.  $P$  Values and Placebo Regression Rejection Frequencies

# Conclusions

# Conclusions

- 1 We propose two-way cluster jackknife CRVEs. New  $CV_3^{(3)}$  and  $CV_3^{(\max)}$  estimators often work very well indeed. They can over-reject or under-reject, but usually quite modestly.

# Conclusions

- 1 We propose two-way cluster jackknife CRVEs. New  $CV_3^{(3)}$  and  $CV_3^{(\max)}$  estimators often work very well indeed. They can over-reject or under-reject, but usually quite modestly.
- 2 In contrast,  $t$ -statistics based on the widely-used  $CV_1^{(3)}$  CRVE for OLS with two-way clustering often over-reject severely.

# Conclusions

- 1 We propose two-way cluster jackknife CRVEs. New  $CV_3^{(3)}$  and  $CV_3^{(\max)}$  estimators often work very well indeed. They can over-reject or under-reject, but usually quite modestly.
- 2 In contrast,  $t$ -statistics based on the widely-used  $CV_1^{(3)}$  CRVE for OLS with two-way clustering often over-reject severely.
- 3 Using an eigen-decomposition when  $V_1^{(3)}$  is not positive definite yields  $CV_1^{(3+)}$  (Stata default). It is parametrization-dependent!

# Conclusions

- 1 We propose two-way cluster jackknife CRVEs. New  $CV_3^{(3)}$  and  $CV_3^{(\max)}$  estimators often work very well indeed. They can over-reject or under-reject, but usually quite modestly.
- 2 In contrast,  $t$ -statistics based on the widely-used  $CV_1^{(3)}$  CRVE for OLS with two-way clustering often over-reject severely.
- 3 Using an eigen-decomposition when  $V_1^{(3)}$  is not positive definite yields  $CV_1^{(3+)}$  (Stata default). It is parametrization-dependent!
- 4 Fixed effects must be handled with care when computing cluster-jackknife ( $CV_3$ ) CRVEs for two-way clustering.

# Conclusions

- 1 We propose two-way cluster jackknife CRVEs. New  $CV_3^{(3)}$  and  $CV_3^{(\max)}$  estimators often work very well indeed. They can over-reject or under-reject, but usually quite modestly.
- 2 In contrast,  $t$ -statistics based on the widely-used  $CV_1^{(3)}$  CRVE for OLS with two-way clustering often over-reject severely.
- 3 Using an eigen-decomposition when  $V_1^{(3)}$  is not positive definite yields  $CV_1^{(3+)}$  (Stata default). It is parametrization-dependent!
- 4 Fixed effects must be handled with care when computing cluster-jackknife ( $CV_3$ ) CRVEs for two-way clustering.
- 5 The  $CV_3^{(2)}$  CRVE is cheaper but usually under-rejects, sometimes severely. So does  $CV_3^{(3+)}$  in some cases with fixed effects.

# Conclusions

- ① We propose two-way cluster jackknife CRVEs. New  $CV_3^{(3)}$  and  $CV_3^{(\max)}$  estimators often work very well indeed. They can over-reject or under-reject, but usually quite modestly.
- ② In contrast,  $t$ -statistics based on the widely-used  $CV_1^{(3)}$  CRVE for OLS with two-way clustering often over-reject severely.
- ③ Using an eigen-decomposition when  $V_1^{(3)}$  is not positive definite yields  $CV_1^{(3+)}$  (Stata default). It is parametrization-dependent!
- ④ Fixed effects must be handled with care when computing cluster-jackknife ( $CV_3$ ) CRVEs for two-way clustering.
- ⑤ The  $CV_3^{(2)}$  CRVE is cheaper but usually under-rejects, sometimes severely. So does  $CV_3^{(3+)}$  in some cases with fixed effects.
- ⑥ The number of regressors and their features matter!