

# Jackknife Standard Errors for Clustered Regression

Bruce E. Hansen\*

University of Wisconsin<sup>†</sup>

This version: December 2023

## Abstract

This paper presents a theoretical case for replacement of conventional heteroskedasticity-consistent and cluster-robust variance estimators with jackknife variance estimators, in the context of linear regression with heteroskedastic and/or cluster-dependent observations. We examine the bias of variance estimation and the coverage probabilities of confidence intervals. Concerning bias, we show that conventional variance estimators have full downward worst-case bias, while our jackknife variance estimator is never downward biased. Concerning confidence intervals, we show that intervals based on conventional standard errors have worst-case coverage equalling zero, while our jackknife-based confidence interval has coverage probability bounded by the Cauchy distribution. We also extend the Bell-McCaffrey (2002) student  $t$  approximation to our jackknife  $t$ -ratio, resulting in confidence intervals with improved coverage probabilities. Our theory holds under minimal assumptions, allowing arbitrary cluster sizes, regressor leverage, within-cluster correlation, heteroskedasticity, regression with a single treated cluster, fixed effects, and delete-cluster invertibility failures. Our theoretical findings are consistent with the extensive simulation literature investigating heteroskedasticity-consistent and cluster-robust variance estimation.

---

\*Research support from the NSF and the Phipps Chair are gratefully acknowledged. Over the course of this research, I have received helpful comments and suggestions from many individuals, including Andrew Chesher, Harold Chiang, Grant Hillier, Rustam Ibragimov, James MacKinnon, Ulrich Müller, Morten Nielsen, Marc Paoella, Peter Phillips, and Thilo Welz. Thanks to J.C. Lazzaro for helpful research assistance.

<sup>†</sup>Department of Economics, 1180 Observatory Drive, University of Wisconsin, Madison WI 53706.

# 1 Introduction

Heteroskedasticity-consistent (HC) and cluster-robust (CRVE) variance estimators and standard errors for regression coefficient estimators are foundational for applied economic analysis. They are used for measuring estimation precision, confidence interval construction, and tests of hypotheses.

Unfortunately, under standard conditions, conventional HC and CRVE variance estimators can be fully downward biased, conventional HC and CRVE t-tests can exhibit unbounded size distortions, and conventional HC and CRVE confidence intervals can have coverage rates equal to zero, even when the regression errors are normally distributed.

This situation – full downward bias, unbounded size, and zero coverage probability – can be corrected by replacing conventional variance estimators by an appropriate jackknife estimator. The latter is simple to calculate. This simple change – the use of jackknife instead of conventional standard errors – has the result that variance estimation is never downward biased and size distortion is bounded.

Our case for jackknife standard errors is based on the following new theoretical insights. First, under standard assumptions, the CRVE variance estimator has full worst-case downward estimation bias, while our jackknife variance estimator is conservative (its expectation is larger than the exact variance). Second, under the assumption of cluster-dependent normality, the worst-case coverage probability of the CRVE-based confidence interval equals zero, while the worst-case coverage probability of our jackknife-based confidence interval is uniformly bounded and controlled by the Cauchy distribution, where uniformity is across all possible regressor and covariance matrix configurations.

Of the above described results, the most important contribution of this paper is the demonstration that the finite sample coverage probability of our jackknife confidence interval is uniformly bounded by the Cauchy distribution. This is not an elementary result, but requires a combination of finite sample analysis, convex analysis, and numerical computation.

All of our theory holds under minimal assumptions, allowing arbitrary cluster sizes, regressor leverage, within-cluster correlation, heteroskedasticity, and delete-cluster invertibility failures. Concerning the latter – delete-cluster invertibility failures – our results allow for the context where a regressor is non-zero only for a single cluster; for example, regressions with included cluster-level fixed effects, regressions with cluster-level treatment indicators when only one cluster is treated, and saturated regressions with sparse cell proportions. In contrast, conventional variance estimators (including conventional jackknife variance estimators) can fail miserably in such contexts.

A new discovery we highlight is that for jackknife variance estimation to be robust to invertibility failures, it is critical to carefully modify the jackknife formula so not to discard clusters. Existing methods (which make *ad hoc* modifications) fail to be robust.

Our theoretical results are backed by the existing extensive simulation literature investigating robust variance estimation which has demonstrated that jackknife standard errors provide dramatically improved finite sample inference in a wide variety of settings.

Furthermore, we investigate the exact distribution of jackknife confidence intervals. We provide a new exact characterization of the finite sample distribution, and provide a practical approximation to this distribution based on a student  $t$  distribution.

The family of heteroskedasticity-consistent variance estimators are often written by the monikers  $HC_0$ ,  $HC_1$ ,  $HC_2$ , and  $HC_3$ . The  $HC_0$  version was introduced by Eicker (1963), Huber (1967), and White (1980). The degree-of-freedom correction known as  $HC_1$  was suggested by Hinkley (1977), and became ubiquitous in applied econometric practice by its designation as the default “r” robust option in Stata. Together,  $HC_0$  and  $HC_1$  are known as Eicker-Huber-White (EHW) variance estimators.  $HC_2$  and  $HC_3$  were introduced by MacKinnon and White (1985) as unbiased estimators under homoskedasticity and the jackknife principle, respectively.

The cluster-robust variance estimator (CRVE) was introduced by Liang and Zeger (1986) and Arellano (1987), is available in Stata through its ubiquitous `cluster` standard error option, and currently dominates applied econometric practice. Bell and McCaffrey (2002) introduced two generalizations for clustered regression similar to  $HC_2$  and  $HC_3$ . The review by MacKinnon, Nielsen, and Webb (2023a) use the monikers  $CV_1$ ,  $CV_2$ , and  $CV_3$  to denote these three estimators.

The recognition that finite-sample inference based on EHW confidence intervals can be severely distorted is a recurrent theme in econometrics. Some investigations include Chesher and Jewitt (1987), Chesher (1989), Chesher and Austin (1991), Long and Ervin (2000), and Young (2019).

There has been a substantial recent literature proposing improved standard errors over the EHW class under independent sampling. This includes Bera, Suprayitno, and Premaratne (2002), Cattaneo, Jansson, and Newey (2018), and Kline, Saggio, and Solvsten (2020). These new variance estimators have the advantage that they are (approximately) unbiased, but have as disadvantages that the variance estimators are computationally burdensome in large samples and are not necessarily positive semi-definite. These methods, while promising, have not been generalized to the clustered sampling setting and are not investigated in this paper.

Jackknife standard errors can be paired with conventional critical values (student  $t$  or normal) or with alternative methods. The latter include the distributional adjustments of Bell and McCaffrey (2002), Imbens and Kolesár (2016), and Pustejovsky and Tipton (2018), bootstrap percentile- $t$  methods (Cameron, Gelbach, and Miller (2008)), and the conditional critical values of Pötscher and Preinerstorfer (2023).

The recommendation to use jackknife/ $HC_3$  variance estimators is not new. Authors making this recommendation include Efron and Stein (1981), MacKinnon and White (1985), Andrews (1991), Chesher and Austin (1991), Long and Erwin (2000), and MacKinnon, Nielsen and Webb (2023abc).

Our analysis applies to all regression contexts. One where the inadequacy of conventional approximations has received particular attention is regression with a small number of clusters and/or a small number of treated clusters. This literature includes Conley and Taber (2011), Ibragimov and Müller (2016), Rokicki, Cohen, Fink, Salomon, and Landrum (2018), Ferman and Pinto (2019), Hagemann (2019, 2023), MacKinnon and Webb (2020), Canay, Santos, and Shaikh (2021), and Niccodemi and Wansbeek (2022). These applications will also benefit from jackknife standard errors and our adjusted degree of freedom approximation, as their application will reduce size distortions.

The organization of the paper is as follows. Section 2 introduces clustered regression and the cluster-robust variance estimator. Section 3 introduces jackknife variance estimation. Section 4 presents results on variance estimation bias. Section 5 presents the core results on the coverage rates of confidence

intervals. Section 6 presents results on the distribution of the jackknife  $t$ -statistic. Section 7 proposes an adjusted degree-of-freedom for a student  $t$  distributional approximation. Section 8 discusses models with cluster-level fixed effects. Section 9 presents simulation evidence. Section 10 presents an empirical application. A conclusion is presented in Section 11. The technical proofs are presented in Section 12.

The R code which generates the numerical calculations presented in the paper, as well as R and Stata code to implement the proposed variance estimator and degree-of-freedom correction, is posted on the author's webpage `users.ssc.wisc.edu/~bhansen/`.

## 2 Clustered Regression and the CRVE

The model is a standard clustered sampling regression. The observations are separated into  $G$  unbalanced mutually independent clusters. Notationally, we write the observations on the  $i$ th individual in the  $g$ th cluster as  $(Y_{ig}, X_{ig})$ , for  $i = 1, \dots, n_g$  and  $g = 1, \dots, G$ . The variable  $Y_{ig}$  is scalar while  $X_{ig}$  is a  $k$ -vector. It is useful to stack the observations by cluster, so that  $\mathbf{Y}_g = (Y_{1g}, \dots, Y_{n_gg})'$  is an  $n_g \times 1$  vector and  $\mathbf{X}_g = (X_{1g}, \dots, X_{n_gg})'$  is an  $n_g \times k$  matrix. The number of observations in the  $g$ th cluster is  $n_g$  and the total number of observations is  $n = \sum_{g=1}^G n_g$ . Stacking the observations conventionally, we obtain the full sample  $(\mathbf{Y}, \mathbf{X})$ .

The observations satisfy the standard linear regression model  $Y_{ig} = X'_{ig}\beta + e_{ig}$  where  $\beta$  is a  $k \times 1$  coefficient vector and  $e_{ig}$  is an error. Written at the level of the cluster, the model is

$$\mathbf{Y}_g = \mathbf{X}_g\beta + \mathbf{e}_g \tag{1}$$

$$\mathbb{E}[\mathbf{e}_g] = 0 \tag{2}$$

where  $\mathbf{e}_g = (e_{1g}, \dots, e_{n_gg})'$ . It is also sometimes convenient to use the full-sample notation  $\mathbf{Y} = \mathbf{X}\beta + \mathbf{e}$ . We will treat the regressors as fixed, but all results go through in the random regressor setting by conditioning. Define the cluster-level covariance matrices

$$\mathbb{E}[\mathbf{e}_g\mathbf{e}'_g] = \boldsymbol{\Sigma}_g. \tag{3}$$

The specification (3) allows the covariance matrices  $\boldsymbol{\Sigma}_g$  to be a function of the regressors (and hence conditionally heteroskedastic), and/or to be a function of the cluster  $g$  (and hence unconditionally heteroskedastic). We follow the clustering literature and impose no structure on  $\boldsymbol{\Sigma}_g$ . We also define the full-sample covariance matrix  $\boldsymbol{\Sigma} = \text{diag}\{\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_G\}$ .

The model (1)-(3) includes heteroskedastic regression as the special case where  $n_g = 1$  for all  $g$ . We call this the “no clustering” or “absence of clustering” case.

**Assumption 1** *Model (1)-(3) holds,  $\mathbf{X}$  is full rank, and  $\boldsymbol{\Sigma}$  has finite elements.*

We focus on the least squares estimator

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1} (\mathbf{X}'\mathbf{Y}) = \left( \sum_{g=1}^G \mathbf{X}'_g \mathbf{X}_g \right)^{-1} \left( \sum_{g=1}^G \mathbf{X}'_g \mathbf{Y}_g \right).$$

It is well known that in under Assumption 1,  $\hat{\beta}$  is unbiased for  $\beta$  with exact covariance matrix

$$\mathbf{V} = \text{var}[\hat{\beta}] = (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{g=1}^G \mathbf{X}'_g \boldsymbol{\Sigma}_g \mathbf{X}_g \right) (\mathbf{X}'\mathbf{X})^{-1}.$$

The nearly ubiquitous cluster-robust variance estimator (CRVE<sub>1</sub>) for  $\mathbf{V}$  is

$$\hat{\mathbf{V}}_1 = \frac{G(n-1)}{(G-1)(n-k)} (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{g=1}^G \mathbf{X}'_g \hat{\boldsymbol{\epsilon}}_g \hat{\boldsymbol{\epsilon}}'_g \mathbf{X}_g \right) (\mathbf{X}'\mathbf{X})^{-1}, \quad (4)$$

where  $\hat{\boldsymbol{\epsilon}}_g = \mathbf{Y}_g - \mathbf{X}_g \hat{\beta}$  denotes the least squares residual vector for the  $g$ th cluster. CRVE<sub>1</sub> was introduced by Liang and Zeger (1986) and Arellano (1987). The constant appearing in (4) is an *ad hoc* degree-of-freedom correction, apparently introduced by the Stata “`cluster`” covariance matrix option.

In the absence of clustering, (4) specializes to the HC<sub>1</sub> “heteroskedasticity-robust” variance estimator, which is the Eicker-Huber-White (EHW) estimator of Eicker (1963), Huber (1967), and White (1980), multiplied by an  $n/(n-k)$  degree-of-freedom correction as suggested by Hickley (1977). HC<sub>1</sub> dominates empirical practice due to its encoding as the “`r`” covariance matrix option in Stata.

As an alternative to CRVE<sub>1</sub>, Bell and McCaffrey (2002), Imbens and Kolesár (2016), and Kolesár (2023) recommended the CRVE<sub>2</sub> variance estimator

$$\hat{\mathbf{V}}_2 = (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{g=1}^G \mathbf{X}'_g \mathbf{M}_g^{+1/2} \hat{\boldsymbol{\epsilon}}_g \hat{\boldsymbol{\epsilon}}'_g \mathbf{M}_g^{+1/2} \mathbf{X}_g \right) (\mathbf{X}'\mathbf{X})^{-1} \quad (5)$$

where  $\mathbf{M}_g^{+1/2}$  is the Moore-Penrose generalized inverse of the symmetric square root of the partial projection matrix

$$\mathbf{M}_g = \mathbf{I}_{n_g} - \mathbf{X}_g (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'_g. \quad (6)$$

In the absence of clustering, and when all  $\mathbf{M}_g$  are invertible, CRVE<sub>2</sub> specializes to the HC<sub>2</sub> estimator of MacKinnon and White (1985).

The original definitions of HC<sub>2</sub> in MacKinnon and White (1985) and CRVE<sub>2</sub> in Bell and McCaffrey (2002) required that all  $\mathbf{M}_g$  are invertible. As discussed by Kolesár (2023), the generalized inverse in (5) allows CRVE<sub>2</sub> to be defined even when  $\mathbf{M}_g$  is non-invertible, and this is the implementation in Stata 18 through its `vce(hc2 clustvar)` option.

The motivation for the CRVE<sub>2</sub> estimator (5) that it is unbiased when the regression errors are i.i.d. (so that  $\boldsymbol{\Sigma}_g = \mathbf{I}_{n_g} \sigma^2$ ) and all  $\mathbf{M}_g$  are invertible.

### 3 Jackknife Variance Estimation

The jackknife estimator of variance of Tukey (1958) extended to clustered dependence is

$$\widehat{V}_3 = \frac{G-1}{G} \sum_{g=1}^G (\widehat{\beta}_{-g} - \bar{\beta})(\widehat{\beta}_{-g} - \bar{\beta})' \quad (7)$$

where

$$\begin{aligned} \widehat{\beta}_{-g} &= \left( \sum_{j \neq g} \mathbf{X}'_j \mathbf{X}_j \right)^{-1} \left( \sum_{j \neq g} \mathbf{X}'_j \mathbf{Y}_j \right) \\ &= \left( \mathbf{X}'\mathbf{X} - \mathbf{X}'_g \mathbf{X}_g \right)^{-1} \left( \mathbf{X}'\mathbf{Y} - \mathbf{X}'_g \mathbf{Y}_g \right) \end{aligned} \quad (8)$$

and

$$\bar{\beta} = \frac{1}{G} \sum_{g=1}^G \widehat{\beta}_{-g}. \quad (9)$$

The delete-one-cluster estimator  $\widehat{\beta}_{-g}$  in (8) is obtained by applying least squares to the sample after deleting the observations in cluster  $g$ . A variant of  $\widehat{V}_3$  is

$$\widehat{V}_4 = \frac{G-1}{G} \sum_{g=1}^G (\widehat{\beta}_{-g} - \widehat{\beta})(\widehat{\beta}_{-g} - \widehat{\beta})', \quad (10)$$

which centers at the full-sample estimator  $\widehat{\beta}$  rather than at  $\bar{\beta}$ . In Stata,  $\widehat{V}_3$  and  $\widehat{V}_4$  can be calculated by the `vce(jackknife)` and `vce(jackknife,mse)` options.

The estimators (7) and (10) as written are undefined if there is a cluster  $g$  for which  $\mathbf{X}'\mathbf{X} - \mathbf{X}'_g \mathbf{X}_g$  is noninvertible<sup>1</sup>. The implementation in Stata 18 circumvents this difficulty by excluding from the sums in (7)-(10) any cluster where (8) is undefined<sup>2</sup>. We follow this interpretation; henceforth, we assume that (7)-(10) are implemented with this modification. We describe (7)-(10) as the ‘‘conventional’’ jackknife variance estimators.

As we show in the next section, the conventional estimators  $\widehat{V}_3$  and  $\widehat{V}_4$  can exhibit downward bias. To eliminate this possibility, we define the following jackknife estimator:

$$\widehat{V}_5 = \sum_{g=1}^G (\widetilde{\beta}_{-g} - \widehat{\beta})(\widetilde{\beta}_{-g} - \widehat{\beta})' \quad (11)$$

where

$$\widetilde{\beta}_{-g} = \left( \mathbf{X}'\mathbf{X} - \mathbf{X}'_g \mathbf{X}_g \right)^+ \left( \mathbf{X}'\mathbf{Y} - \mathbf{X}'_g \mathbf{Y}_g \right) \quad (12)$$

and  $\mathbf{A}^+$  denotes the Moore-Penrose generalized inverse of  $\mathbf{A}$ .  $\widetilde{\beta}_{-g}$  is a generalized delete-one-cluster estimator. It is defined for all  $g$  and therefore (11) includes all clusters.  $\widetilde{\beta}_{-g}$  is a minimizer of the delete-one-cluster least squares criterion, and is therefore a valid delete-one-cluster version of the full-sample

<sup>1</sup>This is identical to the context where  $\mathbf{M}_g$  in (6) is noninvertible.

<sup>2</sup>This follows the recommendation in Shao and Tu (1995) for noninvertible bootstrap replications.

estimator  $\hat{\beta}$ .

Our estimator  $\hat{V}_5$  differs from  $\hat{V}_3$  and  $\hat{V}_4$  in three respects, all of which contribute to the inequality  $\hat{V}_5 > \hat{V}_4 > \hat{V}_3$ . First,  $\hat{V}_5$  does not drop noninvertible clusters. Second,  $\hat{V}_5$  is centered at the full-sample estimator  $\hat{\beta}$  rather than at  $\tilde{\beta}$ . Third,  $\hat{V}_5$  does not have the degree-of-freedom correction  $(G-1)/G$ . In some applications these differences will be negligible, but in others, as we show later, they can be substantial.

Jackknife estimation is ideally suited for the context where the delete-one-cluster estimators  $\tilde{\beta}_{-g}$  are well-defined for all clusters, which requires that the matrices  $\mathbf{X}'\mathbf{X} - \mathbf{X}'_g\mathbf{X}_g$  are invertible for all  $g$ . We call this context *clusterwise invertibility* and its failure *clusterwise noninvertibility*. Noninvertibility occurs when deletion of a single cluster renders the regressor design matrix singular. Most typically, this occurs when a regressor (or some linear combination of regressors) only takes non-zero values for a single cluster. Examples include regressions with included cluster-level fixed effects, regressions with cluster-level treatment indicators when only one cluster is treated, and saturated regressions with sparse cell proportions. Such regressions are commonplace in applications, so it is desirable for variance estimation to be sufficiently flexible to handle their occurrence, which is presumably the motivation for the “drop noninvertible clusters” modification described above. However, the properties of any such modification need to be investigated to preclude undesirable outcomes.

When the sample satisfies clusterwise invertibility then the estimators  $\hat{V}_4$  and  $\hat{V}_5$  only differ by the degree-of-freedom correction  $(G-1)/G$ , which is typically inconsequential. They can differ more substantially, however, under clusterwise noninvertibility.

The jackknife estimators  $\hat{V}_3$  and  $\hat{V}_4$  for clustered samples were developed by Cochran (1977), Rust and Rao (1996), and Bell and McCaffrey (2002). See MacKinnon, Nielsen, and Webb (2023abc) for detailed discussions. Stata has codified the modification to delete noninvertible clusters. MacKinnon, Nielsen, and Webb (2023c) propose and implement a generalized delete-one-cluster estimator similar to (12), though they do not investigate the statistical properties of the resulting jackknife variance estimator.

Alternative algebraic representations of the jackknife estimator are available; see MacKinnon, Nielsen, and Webb (2023b). One we will find useful is based on the delete-one-cluster prediction errors

$$\tilde{\mathbf{e}}_g = \mathbf{Y}_g - \mathbf{X}_g\tilde{\beta}_{-g}. \quad (13)$$

In the proof of Theorem 1 in Section 12 we show that

$$\tilde{\beta}_{-g} - \hat{\beta} = -(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'_g\tilde{\mathbf{e}}_g, \quad (14)$$

and interestingly, this equality holds even allowing for clusterwise noninvertibility. Given (14), we can write (11) as

$$\hat{V}_5 = (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{g=1}^G \mathbf{X}'_g\tilde{\mathbf{e}}_g\tilde{\mathbf{e}}'_g\mathbf{X}_g \right) (\mathbf{X}'\mathbf{X})^{-1}. \quad (15)$$

The jackknife estimators (7)-(11) simplify in the absence of clustering, though existing proposals and

analysis all assume that the  $\mathbf{M}_g$  are all invertible<sup>3</sup>. Under these conditions the estimator (7) corresponds to that proposed by MacKinnon and White (1985) and the estimator (11) corresponds to that proposed by Andrews (1991) and Davidson and MacKinnon (1993). The latter is known as HC<sub>3</sub>, and can be calculated in Stata by the `vce(hc3)` option. When clusterwise invertibility fails, the estimators HC<sub>2</sub> and HC<sub>3</sub> are undefined<sup>4</sup>. For more details, see the monographs of Efron (1982) and Shao and Tu (1995).

## 4 Biased vs Conservative Variance Estimation

In the classical regression model, the classical variance estimator is unbiased for the exact finite sample variance. Unsurprisingly, outside the classical model, robust variances estimators are not unbiased. We now examine the worst-case downward bias of the five variance estimators described in the previous sections. We focus on downward bias, as this is the issue which causes undercoverage of confidence intervals and oversized tests. The main contribution of this section is the following result.

**Theorem 1** *Under Assumption 1,*

$$\mathbb{E}[\widehat{\mathbf{V}}_5] \geq \mathbf{V}. \quad (16)$$

Theorem 1 shows that our recommended jackknife estimator  $\widehat{\mathbf{V}}_5$  is *never downward biased* in a positive semi-definite sense. A never-downward-biased variance estimator can be described as *conservative*. This means that in any regression context, and any sample size, we can be confident that the jackknife estimator is not downward biased. We will find that this bias property is important as it is directly connected to the coverage probabilities of confidence intervals.

Theorem 1 holds quite broadly, holding for all sample sizes, regressor matrices, variance matrices, and violations of clusterwise noninvertibility. In particular, the robustness to clusterwise noninvertibility is new and surprising.

Theorem 1 augments Theorem 2 of Bell and McCaffrey (2002), which established that  $\widehat{\mathbf{V}}_4$  is never downward biased when the regressors satisfy clusterwise invertibility and the errors  $e_{ig}$  are i.i.d. (that is, when  $\boldsymbol{\Sigma} = \mathbf{I}_n \sigma^2$ ).

Theorem 1 is also related to the seminal work of Efron and Stein (1981). Their results are typically described as stating that (16) holds for  $\widehat{\mathbf{V}}_3$  but this is incorrect. Instead, Efron and Stein's Theorem 2 states that  $\widehat{\mathbf{V}}_3$  is never-downward-biased as an estimator of  $\text{var}[\tilde{\boldsymbol{\beta}}]$ , not as an estimator of  $\text{var}[\widehat{\boldsymbol{\beta}}]$ . Furthermore, Theorem 2 below shows  $\widehat{\mathbf{V}}_3$  does not satisfy the never-downward-biased property (16).

We now explore the worst-case bias properties of the other CRVE variance estimators. For these results we focus on individual coefficient estimates and their variance estimators. For some non-zero  $k \times 1$  vector  $R$ , define the scalar parameter  $\theta = R' \boldsymbol{\beta}$ . This includes individual coefficients and linear combinations. Its estimator is  $\widehat{\theta} = R' \widehat{\boldsymbol{\beta}}$ . Under Assumption 1,  $\widehat{\theta}$  is unbiased for  $\theta$  and has exact variance

$$v^2 = \text{var}[\widehat{\theta}] = R' \mathbf{V} R.$$

<sup>3</sup>In the absence of clustering, this means that the leverage values  $h_{ii} = \mathbf{X}'_i (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}_i$  all satisfy  $h_{ii} < 1$ .

<sup>4</sup>A word of caution: Stata 18 reports HC<sub>2</sub> and HC<sub>3</sub> standard errors even when the regressor matrix is clusterwise noninvertible, and does not document how they are calculated in this situation.



The estimators of  $v^2$  are  $\hat{v}_j^2 = R' \hat{V}_j R$  for  $j = 1, \dots, 5$ . Standard errors for  $\hat{\theta}$  are their square roots  $\hat{v}_j = \sqrt{R' \hat{V}_j R}$ .

It will be convenient to define sets of models. For fixed  $k$  and  $G$  let  $\mathcal{F}$  be the class of all regressor and covariance matrices  $(\mathbf{X}, \mathbf{\Sigma})$  such that  $\mathbf{X}$  is full rank,  $\mathbf{\Sigma}$  has finite elements, and  $v^2 > 0$ . Let  $\mathcal{F}^* \subset \mathcal{F}$  be the subset where  $\mathbf{X}$  satisfies clusterwise invertibility. Let  $\mathcal{F}_0^* \subset \mathcal{F}^*$  and  $\mathcal{F}_0 \subset \mathcal{F}$  be the further subsets where  $\mathbf{\Sigma} = \mathbf{I}_n \sigma^2$ .

The following result shows that the variance estimators other than  $\hat{V}_5$  can be severely biased.

**Theorem 2** *Suppose Assumption 1 holds. For  $(\mathbf{X}, \mathbf{\Sigma}) \in \mathcal{F}_0^*$ ,*

$$\inf_{(\mathbf{X}, \mathbf{\Sigma}) \in \mathcal{F}_0^*} \frac{\mathbb{E}[\hat{v}_1^2]}{v^2} = 0. \quad (17)$$

For  $(\mathbf{X}, \mathbf{\Sigma}) \in \mathcal{F}^*$ ,

$$\inf_{(\mathbf{X}, \mathbf{\Sigma}) \in \mathcal{F}^*} \frac{\mathbb{E}[\hat{v}_2^2]}{v^2} = 0, \quad (18)$$

$$\inf_{(\mathbf{X}, \mathbf{\Sigma}) \in \mathcal{F}^*} \frac{\mathbb{E}[\hat{v}_3^2]}{v^2} = \left( \frac{G-1}{G} \right)^2 < 1, \quad (19)$$

and

$$\inf_{(\mathbf{X}, \mathbf{\Sigma}) \in \mathcal{F}^*} \frac{\mathbb{E}[\hat{v}_4^2]}{v^2} = \frac{G-1}{G} < 1. \quad (20)$$

For  $(\mathbf{X}, \mathbf{\Sigma}) \in \mathcal{F}_0$ ,

$$\inf_{(\mathbf{X}, \mathbf{\Sigma}) \in \mathcal{F}_0} \frac{\mathbb{E}[\hat{v}_3^2]}{v^2} = 0 \quad (21)$$

and

$$\inf_{(\mathbf{X}, \mathbf{\Sigma}) \in \mathcal{F}_0} \frac{\mathbb{E}[\hat{v}_4^2]}{v^2} = 0. \quad (22)$$

For  $(\mathbf{X}, \mathbf{\Sigma}) \in \mathcal{F}$ ,

$$\inf_{(\mathbf{X}, \mathbf{\Sigma}) \in \mathcal{F}} \frac{\mathbb{E}[\hat{v}_5^2]}{v^2} = 1. \quad (23)$$

Equation (17) shows that the expected value of the scaled CRVE<sub>1</sub> estimator can be arbitrarily close to zero. This means that there is some distribution under which variance estimation has arbitrarily large downward bias. In words, we say that CRVE<sub>1</sub> is *fully downwardly biased*. The set of models considered are those satisfying  $\mathbf{\Sigma} = \mathbf{I}_n \sigma^2$  and clusterwise invertibility, so the result (17) has nothing to do with heteroskedasticity, correlated errors, or invertibility failure. Rather, it is a consequence of extreme regressor leverage (unbalanced regressors and/or cluster sizes).

Equation (18) shows a similar result for CRVE<sub>2</sub> in the class of models which allows general covariance matrices (heteroskedasticity and correlation) but still requires clusterwise invertibility. The difference between (17) and (18) is that CRVE<sub>1</sub> can be fully downward biased due to regressor leverage alone, while the result for CRVE<sub>2</sub> requires non-i.i.d. errors.

Equations (19) and (20) show related results for the conventional jackknife estimators. They show that  $\hat{v}_3^2$  and  $\hat{v}_4^2$  violate the never-downward-biased property, even under clusterwise invertibility. The

magnitude of the violation, however, is small when  $G$  is large. Regardless, (19) and (20) show that the conventional interpretation<sup>5</sup> of the results of Efron and Stein (1981) is incorrect, and that the never-downward-biased result of Bell and McCaffrey (2002) is not robust to non-i.i.d. errors.

Equations (21) and (22) extend the analysis of the conventional jackknife estimators to the class of models allowing clusterwise noninvertibility. The results show that under clusterwise noninvertibility the conventional estimators  $\hat{v}_3^2$  and  $\hat{v}_4^2$  are fully downward biased. This arises even though (21)-(22) restrict the set of models to satisfy  $\Sigma = I_n \sigma^2$ . This means that the full downward bias of the conventional jackknife estimators is entirely a consequence of the deletion of noninvertible clusters, and does not require heteroskedasticity or correlated errors.

Equation (23) examines our recommended jackknife estimator  $\hat{v}_5^2$  in the broadest context, allowing for general heteroskedasticity, correlated errors, and clusterwise noninvertibility. As implied by Theorem 1,  $\hat{v}_5^2$  is never downward biased. Equation (23) extends this further and demonstrates that the infimum equals one, meaning that the inequality of Theorem 1 is sharp.

Theorem 2 draws a stark contrast between the five variance estimators. The full downward bias of  $CRVE_1$ ,  $CRVE_2$ , and the conventional jackknife estimators means that on average they can be “much too small” relative to the true variance. In contrast, the never-downward-biased property of  $\hat{v}_5^2$  means that there is no situation where it is expected to be “too small”, even slightly.

The model classes studied in Theorem 2 hold fixed the number of regressors  $k$  and number of clusters  $G$ , but allow the cluster sizes  $n_g$ , regressors  $X$ , and covariance matrices  $\Sigma$  to vary freely. It is important to understand that the statements of Theorem 2 hold for all  $G$ , from the very small to very large. Thus the bias of  $CRVE_1$ ,  $CRVE_2$ , and the conventional jackknife can be arbitrarily bad in both very small and very large samples.

The worstcase downward bias in (17)-(22) is calculated by studying models with extreme leverage, arising when the regressor of interest has variation which is dominated by a single cluster. Intuitively, when a small number of clusters dominate the sample, standard variance estimators are highly biased towards zero. The least squares estimator overfits the dominating clusters, shrinking the residuals for these clusters relative to the true errors. This leads to downward estimation bias. Conventional fixes, such as the degree-of-freedom adjustment of  $CRVE_1$ , are insufficient to counter the bias.

## 5 Unbounded vs Bounded Inference

Given a standard error  $\hat{v}_j$  and a pre-selected critical value  $c$ , a confidence interval for  $\theta$  is

$$\hat{C}_j(c) = \hat{\theta} \pm c \hat{v}_j.$$

For a nominal  $100(1 - \alpha)\%$  interval the conventional<sup>6</sup> choice for the critical value is  $c = t_{G-1}^{1-\alpha/2}$ , the  $1 - \alpha/2$  quantile of the student  $t$  distribution with  $G - 1$  degrees of freedom.

For the results of this section, we require that the errors are normally distributed.

<sup>5</sup>This distinction was recognized by Efron and Stein (1981). For a further discussion see Section 4.5 of Efron (1982).

<sup>6</sup>For example, that used in Stata.

**Assumption 2**  $e_g$  is distributed  $N(0, \Sigma_g)$ .

Assumption 2 states that the cluster error vectors  $e_g$  are normally distributed, but their within-cluster covariance structure is left unrestricted. As described after (3), this allows the covariance matrices to vary both with the regressors as well as the cluster, and thus allows both unconditional and conditional heteroskedasticity.

Normality is a strong assumption, and is not meant to be taken literally. Rather, by studying coverage rates under this assumption we can gain insight into their behavior in finite samples without relying on asymptotic approximations.

We now evaluate the worst-case coverage probabilities of CRVE confidence intervals.

Under i.i.d. normal errors, it is well known that confidence intervals with classical standard errors have exact coverage  $1 - \alpha$ . What may be less well known (or at least less emphasized in standard curricula) is that this result does not extend to confidence intervals constructed with HC and CRVE standard errors.

We now present the most important contribution of the paper, which provides a lower bound on the coverage rate of our recommended jackknife confidence interval. Let  $F(x; k_1, k_2)$  denote the  $F$  distribution function with degrees of freedom  $k_1$  and  $k_2$ .

**Theorem 3** Under Assumptions 1-2, for any  $1 \leq c < \infty$ ,

$$\inf_{(X, \Sigma) \in \mathcal{F}} \mathbb{P}[\theta \in \widehat{C}_5(c)] \geq F(c; 1, 1). \quad (24)$$

Equation (24) shows that the interval  $\widehat{C}_5(c)$  has coverage probability which is uniformly bounded<sup>7</sup> away from zero in the broad model class  $\mathcal{F}$ . The lower bound is the  $F$  distribution with (1,1) degrees of freedom, which is the square of the Cauchy distribution. An important implication is that the finite sample coverage probability of the  $\widehat{C}_5(c)$  jackknife confidence interval has bounded distortion from its nominal level. This result holds over the broadest model class  $\mathcal{F}$ . This includes all regression models, including those with the most extreme leverage, within-cluster correlation, heteroskedasticity, and clusterwise noninvertibility. Thus the jackknife interval  $\widehat{C}_5(c)$  is robust to these contexts.

We contrast (24) with the worst-case coverage of confidence intervals constructed with the other standard errors.

**Theorem 4** Under Assumptions 1-2, for any  $0 \leq c < \infty$ ,

$$\inf_{(X, \Sigma) \in \mathcal{F}_0^*} \mathbb{P}[\theta \in \widehat{C}_1(c)] = 0, \quad (25)$$

$$\inf_{(X, \Sigma) \in \mathcal{F}^*} \mathbb{P}[\theta \in \widehat{C}_2(c)] = 0, \quad (26)$$

$$\inf_{(X, \Sigma) \in \mathcal{F}_0} \mathbb{P}[\theta \in \widehat{C}_3(c)] = 0, \quad (27)$$

---

<sup>7</sup>The bound (28) requires  $c \geq 1$ . This is not an important restriction for inference as all conventional critical values exceed 1.

and

$$\inf_{(\mathbf{X}, \boldsymbol{\Sigma}) \in \mathcal{F}_0} \mathbb{P}[\theta \in \widehat{C}_4(c)] = 0. \quad (28)$$

Equation (25) shows that the worst-case coverage of the CRVE<sub>1</sub> confidence interval equals 0. This demonstrates that coverage can be arbitrarily distorted from the nominal level, and this holds for *any* critical value  $c$ . This means that it is impossible to uniformly<sup>8</sup> achieve any desired coverage probability. The set of models considered are those satisfying  $\boldsymbol{\Sigma} = \mathbf{I}_n \sigma^2$  and clusterwise invertibility, so the result (25) is not due to heteroskedasticity, correlated errors, or invertibility failure. Rather, it is a consequence of extreme regressor leverage (unbalanced regressors and/or cluster sizes).

Equation (26) shows a similar result for the CRVE<sub>2</sub> confidence interval in the class of models which allows general covariance matrices: that the worst-case coverage of the CRVE<sub>2</sub> confidence interval equals 0. Again, this demonstrates that coverage can be arbitrarily distorted from the nominal level. The difference with (25) is that (26) requires the model class to include non-i.i.d. errors.

Equations (27) and (28) show that the conventional jackknife intervals also have worst-case coverage of 0 if the model class is broadened to include clusterwise noninvertibility. This is due to invertibility failure. These results show that the conventional (e.g., Stata) modification, which is explicitly intended to allow regressions with clusterwise noninvertibility, is not actually robust to clusterwise noninvertibility.

The results (25)-(28) should not be surprising given Theorem 2, which showed that the variance estimators  $\widehat{v}_1^2 - \widehat{v}_4^2$  can be arbitrarily downward biased. Indeed, the proof of Theorem 4 is a simple manipulation of Theorem 2. What is important about these results is that they show that these confidence intervals have no *a priori* guarantee that they are in any sense a confidence interval. It is also important to understand that the zero coverage rates of (25)-(28) cannot be fixed by simply using a larger critical value  $c$ , as these results hold for any finite  $c$ .

Returning to the confidence interval  $\widehat{C}_5(c)$ , equation (24) bounds its smallest coverage probability in *any* regression with normal errors. For example, with the conventional  $c = 1.96$  critical value, the bound (24) is 0.70. Thus, the finite sample coverage of  $\widehat{C}_5(1.96)$  can never be less than 70%. Similarly, the finite sample size of a t-test using the standard error  $\widehat{v}_5$  and critical value  $c = 1.96$  can never be greater than 30%.

Another implication of (24) is that the Cauchy distribution can be used for finite sample inference (substituting the Cauchy for student  $t$  critical values). Doing so will produce inferential statements (hypothesis tests and confidence intervals) with uniform size control. This uniformity holds over all regression designs  $\mathbf{X}$  and error variances  $\boldsymbol{\Sigma}$ .

In practice, however, it is unlikely that researchers will use the Cauchy distribution for inference, as it is exceedingly conservative. For example, while the 5% normal critical value is 1.96, that for the Cauchy distribution is 12.7. It is difficult to imagine a user declaring a t-ratio equalling 10 to be “insignificant” simply because it is less than the Cauchy critical value.

Instead, the practical message of Theorems 3 and 4 is that conventional EHW and CRVE confidence intervals can have arbitrary coverage distortion, while the jackknife interval has bounded distortion.

---

<sup>8</sup>Here, “uniformly” means over all regression designs  $\mathbf{X}$ .

This is a strong motivation for replacement of the EHW and CRVE standard errors by simple-to-calculate jackknife standard errors.

While the bound (24) may appear simple, its derivation is not. Our derivation relies on a characterization of the exact distribution of the t-ratio, the variance bound (23), convex optimization, and numerical calculation. The use of numerical calculation in the proof is theoretically inelegant; it is only needed for the clusterwise noninvertible case. Under the added assumption of clusterwise invertibility, the proof of Theorem 3 does not require a numerical argument.

## 6 Distribution of Jackknife t statistic

In this section we present a characterization of the finite sample distribution of the squared jackknife t-statistic

$$T_5^2 = \frac{(\hat{\theta} - \theta)^2}{\hat{v}_5^2}. \quad (29)$$

Let  $F(x; k_1, k_2)$  denote the  $F$  distribution function with degrees of freedom  $k_1$  and  $k_2$ .

**Theorem 5** *Under Assumptions 1-2,  $v^2 > 0$ , and any  $0 \leq x < \infty$ ,*

$$\mathbb{P}[T_5^2 \leq x] \geq \mathbb{P}\left[\frac{v^2 \xi_0^2}{\sum_{j=1}^G \phi_j \xi_j^2} \leq x\right] \simeq F(a^2 x; 1, K) \quad (30)$$

where  $\xi_j^2$  are mutually independent  $\chi_1^2$  random variables,  $\phi_j$  are the eigenvalues of the  $G \times G$  matrix  $\mathbf{L}$  defined in (94),

$$a = \sqrt{\frac{\text{tr}[\mathbf{L}]}{v^2}}, \quad (31)$$

and

$$K = \frac{(\text{tr}[\mathbf{L}])^2}{\text{tr}[\mathbf{L}\mathbf{L}]} \quad (32)$$

The approximation in (30) is the Satterthwaite (1946) approximation to the weighted sum of chi-squares  $\sum_{j=1}^G \phi_j \xi_j^2$ .

Theorem 5 provides two distributional characterizations in (30). The first shows that the exact distribution of  $T_5^2$  is bounded below by that of a ratio of independent weighted sums of chi-squares. The second shows that the latter is approximately  $F$  distributed with a non-standard scale and degree-of-freedom. Equivalently, this shows that the distribution of the t-statistic  $T_5 = (\hat{\theta} - \theta) / \hat{v}_5$  is bounded by an approximate student  $t$  with scale  $a$  and degree-of-freedom  $K$ .

One new feature of Theorem 5 is the bound for the exact distribution by a ratio of *independent* weighted sums of chi-squares. It is well known that this representation holds for common t-statistics under i.i.d. errors, and is also widely known that a similar representation, but with correlated weighted sums of chi-squares, holds for non-i.i.d. errors. The fact that the representation (30) holds for independent  $\chi_1^2$  variables is new.

The constants  $a$  and  $K$  in (31)-(32) are written as functions of the eigenvalues of the  $G \times G$  matrix  $\mathbf{L}$  defined in the proof of Theorem 5. Of particular importance is the non-standard degree-of-freedom  $K$ . Examining (32) and observing that the matrix  $\mathbf{L}$  is  $G \times G$ , we can deduce that  $K$  takes values in  $[1, G]$ . The case  $K = G$  occurs when the matrix  $\mathbf{L}$  is perfectly balanced so that all its eigenvalues are equal. The case  $K = 1$  occurs when the matrix  $\mathbf{L}$  is highly unbalanced with only one non-zero eigenvalue. Most applications will lie between these two extremes. Note in contrast that conventional software uses  $G - 1$  degrees of freedom.

The constants  $a$  and  $K$  which determine the F distribution of Theorem 5 depend on the unknown variance matrices  $\Sigma_g$ . If the  $\Sigma_g$  are known then  $a$  and  $K$  can be calculated and this F distribution used. However, in practice these variance matrices are unknown. Bell and McCaffery suggested calculating  $a$  and  $K$  based on a reference model, in particular,  $\Sigma_g = \mathbf{I}_{n_g}$  (which holds under i.i.d. errors). Other plausible reference models are  $\Sigma_g = \mathbf{1}_{n_g} \mathbf{1}'_{n_g}$  (holds under perfect within-cluster correlation) and  $\Sigma_g = \mathbf{X}_g \mathbf{X}'_g$  (strong conditional heteroskedasticity). We follow Bell and McCaffery and recommend calculation of  $a$  and  $K$  based on the simple reference model  $\Sigma_g = \mathbf{I}_{n_g}$ . If desired, other reference models for  $\Sigma_g$  could be used.

We now present explicit expressions for the constants  $a$  and  $K$  under this reference model which do not rely on explicit computation of the matrix  $\mathbf{L}$  nor its eigenvalues.

**Theorem 6** *Computationally convenient expressions for the constants appearing in (31)-(32) for the case  $\Sigma_g = \mathbf{I}_{n_g}$  are*

$$v^2 = R' (\mathbf{X}' \mathbf{X})^{-1} R, \quad (33)$$

$$\text{tr}[\mathbf{L}] = \sum_{g=1}^G \mathbf{S}'_g \mathbf{S}_g - \frac{\mathbf{T}' \mathbf{T}}{\mathbf{Z}' \mathbf{Z}} + \text{tr}[\mathbf{U}' \mathbf{U} \hat{\mathbf{X}}' \hat{\mathbf{X}}] - 2 \text{tr}[\mathbf{U}' \mathbf{V}], \quad (34)$$

and

$$\begin{aligned} \text{tr}[\mathbf{L}\mathbf{L}] &= \sum_{g=1}^G (\mathbf{S}'_g \mathbf{S}_g)^2 - 2 \frac{\sum_{g=1}^G (\mathbf{S}'_g \mathbf{S}_g) T_g^2}{\mathbf{Z}' \mathbf{Z}} + \left( \frac{\mathbf{T}' \mathbf{T}}{\mathbf{Z}' \mathbf{Z}} \right)^2 - 2 \frac{\mathbf{T}' \mathbf{U} \hat{\mathbf{X}}' \hat{\mathbf{X}} \mathbf{U}' \mathbf{T}}{\mathbf{Z}' \mathbf{Z}} \\ &+ \text{tr}[\hat{\mathbf{X}}' \hat{\mathbf{X}} \mathbf{U}' \mathbf{U} \hat{\mathbf{X}}' \hat{\mathbf{X}} \mathbf{U}' \mathbf{U}] + 2 \text{tr}[\mathbf{V}' \mathbf{U} \mathbf{V}' \mathbf{U}] + 2 \text{tr}[\mathbf{U}' \mathbf{W} \hat{\mathbf{X}}' \hat{\mathbf{X}}] \\ &- 4 \text{tr}[\mathbf{V}' \mathbf{W}] + 4 \frac{\mathbf{T}' \mathbf{V} \mathbf{U}' \mathbf{T}}{\mathbf{Z}' \mathbf{Z}} - 4 \text{tr}[\mathbf{U}' \mathbf{U} \hat{\mathbf{X}}' \hat{\mathbf{X}} \mathbf{U}' \mathbf{V}] + 2 \text{tr}[\mathbf{U}' \mathbf{U} \mathbf{V}' \mathbf{V}], \end{aligned} \quad (35)$$

where

$$\begin{aligned}
\mathbf{Z} &= \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}R \\
\widehat{\mathbf{X}} &= \mathbf{X} - \mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{X} \\
\mathbf{U}_g &= (\mathbf{X}'\mathbf{X} - \mathbf{X}'_g\mathbf{X}_g)^+ \mathbf{X}'_g\mathbf{Z}_g \\
\mathbf{S}_g &= \mathbf{Z}_g + \mathbf{X}_g\mathbf{U}_g \\
\mathbf{V}_g &= \widehat{\mathbf{X}}'_g\mathbf{S}_g \\
\mathbf{W}_g &= \mathbf{U}_g\mathbf{S}'_g\mathbf{S}_g \\
\mathbf{T}_g &= \mathbf{S}'_g\mathbf{Z}_g
\end{aligned}$$

$$\mathbf{U} = \begin{bmatrix} \mathbf{U}'_1 \\ \vdots \\ \mathbf{U}'_G \end{bmatrix}, \quad \mathbf{V} = \begin{bmatrix} \mathbf{V}'_1 \\ \vdots \\ \mathbf{V}'_G \end{bmatrix}, \quad \mathbf{W} = \begin{bmatrix} \mathbf{W}'_1 \\ \vdots \\ \mathbf{W}'_G \end{bmatrix}, \quad \mathbf{S} = \begin{bmatrix} \mathbf{S}_1 \\ \vdots \\ \mathbf{S}_G \end{bmatrix}, \quad \mathbf{T} = \begin{bmatrix} T_1 \\ \vdots \\ T_G \end{bmatrix}.$$

The expression (34) and (35) in Theorem 6 look imposing, but they are computationally simple to implement. All of the matrices for which the trace operator are applied are  $k \times k$ . The expressions (34)-(35) are functions of (and only of) the regressors  $\mathbf{X}$  and vector  $R$ .

## 7 Adjusted Confidence Intervals

Simulation evidence shows that while our jackknife confidence interval using conventional critical values has improved coverage rates relative to the  $\text{CRVE}_1$  and  $\text{CRVE}_2$  intervals, it can still under-cover. The source of the problem is the conventional student  $t$  distributional approximation.

A promising class of distributional adjustments was introduced by Bell and McCaffery (2002) for  $\text{HC}_2$  and  $\text{CRVE}_2$   $t$ -ratios. They are simple to calculate and produce greatly improved coverage rates in finite samples. These methods have been endorsed by other authors, including Imbens and Kolesár (2016). In this section, we extend the Bell-McCaffery adjustment to jackknife  $t$  tests and confidence intervals. Our extension involves both scale and degree-of-freedom adjustments.

Our Bell-McCaffery adjustment is based on the approximation of Theorem 5, which states that the finite sample distribution of the squared  $t$ -statistic is approximately  $F$  distributed (or the non-squared  $t$ -statistic is approximately student  $t$  distributed) with a scale adjustment and a non-standard degree-of-freedom. The original Bell-McCaffery adjustment (for  $\text{HC}_2$  and  $\text{CRVE}_2$   $t$ -ratios) does not have a scale adjustment, because the  $\text{CRVE}_2$  variance estimator is unbiased under i.i.d. errors. In contrast, our adjustment for the jackknife statistic has a scale adjustment to account for jackknife variance estimation bias.

As discussed in the previous section, these approximations need to be calculated under a reference model for the unknown variances, and following Bell and McCaffery we recommend  $\boldsymbol{\Sigma}_g = \mathbf{I}_{n_g}$ . In this context, efficient formula for the inputs for the constants  $a$  and  $K$  are given in expressions (33)-(35).

The adjusted  $100(1 - \alpha)\%$  confidence interval for  $\theta$  is

$$\tilde{C}_5 = \hat{\theta} \pm \frac{t_K^{1-\alpha/2} \hat{v}_5}{a} \quad (36)$$

where  $t_K^{1-\alpha/2}$  is the  $1 - \alpha/2$  quantile of the student  $t$  distribution with  $K$  degrees of freedom and  $\hat{v}_5$  is our recommended jackknife standard error. Equivalently,  $t_K^{1-\alpha/2} = \sqrt{q_{1,K}^{1-\alpha}}$ , where  $q_{1,K}^{1-\alpha}$  is the  $1 - \alpha$  quantile of the  $F(x; 1, K)$  distribution.

Specifically, the confidence interval is formed by the following steps. Take any linear coefficient  $\theta = R' \beta$ .

1. Calculate the jackknife variance  $\hat{V}_5$  from (11) and standard error  $\hat{v}_5 = \sqrt{R' \hat{V}_5 R}$ .
2. Calculate the adjustment parameters  $K$  and  $a$  from (34) and (35).
3. Calculate the critical value  $t_K^{1-\alpha/2}$  from the student  $t$  distribution with the degree-of-freedom  $K$ .
4. Set  $\hat{\theta} = R' \hat{\beta}$  and  $\tilde{C}_5$  as in (36).

The adjustment parameters  $K$  and  $a$  are coefficient-specific, so the above steps need to be repeated for each element of  $\beta$  in order to calculate confidence intervals for all individual coefficients.

The adjusted degree-of-freedom  $K$  will typically be smaller than the conventional degree-of-freedom, thereby increasing the critical value and the width of the confidence interval. This adjustment accounts for the non-standard distribution of the variance estimator. In contrast, the adjustment factor  $a$  is typically greater than 1, so the scale adjustment decreases the width of the confidence interval. This counters the bias of the jackknife variance estimator.

Similarly, the adjusted p-value for a test of  $\theta = \theta_0$  is

$$p = 1 - F\left(a^2 \left(\frac{\hat{\theta} - \theta_0}{\hat{v}_5}\right)^2; 1, K\right) \quad (37)$$

where  $F(x; 1, K)$  is the F distribution with degrees of freedom  $(1, K)$ . As for the confidence intervals, this makes both a degree-of-freedom adjustment (through  $K$ ) and a scale adjustment (through  $a$ ).

The theoretical justification for the confidence interval (36) is as follows. Using the equalities  $\tilde{C}_5 = \hat{C}_5(t_K^{1-\alpha/2}/a)$ ,  $t_K^{1-\alpha/2} = \sqrt{q_{1,K}^{1-\alpha}}$ , then Theorem 5, we find

$$\begin{aligned} \mathbb{P}[\theta \in \tilde{C}_5] &= \mathbb{P}[\theta \in \hat{C}_5(t_K^{1-\alpha/2}/a)] \\ &= \mathbb{P}\left[T_5^2 \leq q_{1,K}^{1-\alpha}/a^2\right] \\ &\gtrsim F\left(q_{1,K}^{1-\alpha}; 1, K\right) \\ &= 1 - \alpha. \end{aligned} \quad (38)$$

The inequality (38) combines the inequality and approximation of Theorem 5, and holds under the reference model  $\Sigma = I_n \sigma^2$ .



This shows that, approximately,  $\tilde{C}_5$  will have coverage close to or exceeding the desired level. There are three sources of approximation error: (1) the inequality bound in (30); (2) the Satterthwaite approximation in (30); (3) the reference model approximation. Each may contribute to a distortion of actual from nominal coverage.

## 8 Fixed Effects

It is common in panel/clustering contexts to include cluster-specific fixed effects<sup>9</sup>. This possibility is allowed in our framework either by explicit inclusion of fixed effect dummy variables in the regressor matrix  $\mathbf{X}$ , or by specifying  $\mathbf{Y}_g$  and  $\mathbf{X}_g$  as within-transformed. In the latter case, equation (1) represents the model *after* the within transformation has been applied. Note that in this case, the covariance matrix (3) is that of the within-transformed errors, not the original equation errors. As (3) is unstructured and allows  $\Sigma_g$  to be singular, this is without loss of generality.

When fixed effect dummy variables are present in the regressor matrix  $\mathbf{X}$ , the traditional jackknife variances estimators  $\hat{V}_3$  and  $\hat{V}_4$  are undefined. In contrast, our jackknife variance estimator  $\hat{V}_5$  is well-defined, and conservative by Theorem 1 for all the coefficients including the estimated fixed effects.

Regardless, our recommendation<sup>10</sup> is to first apply the within transformation and then apply least squares estimation to the within transformed variables, to obtain the coefficient estimates and jackknife standard errors. This is both computationally and theoretically preferred. It is theoretically preferred because the model after the within transformation is clusterwise invertible, so has reduced distributional distortions. In contrast, the model with the included fixed effect dummy variables is clusterwise noninvertible, and while our finite sample distribution theory applies, the inequalities suggest that the confidence intervals will be conservative.

## 9 Simulation Evidence

We present a simulation experiment to investigate the performance of the methods. The experiment concerns the clustered linear regression (1)-(3). The goal is 95% confidence intervals for the slope coefficients.

Our baseline model is the simple regression  $\mathbf{Y}_g = \alpha + \mathbf{X}_g\beta + \mathbf{e}_g$  with  $\mathbf{X}_g$  a single  $n_g \times 1$  stochastic regressor. Within this model we consider six designs which vary the distributions of the regressors  $\mathbf{X}_g$ , the heteroskedasticity of the equation errors, and the cluster size heterogeneity. Let  $\mathbf{I}_g$  denote the  $n_g \times n_g$  identity matrix,  $\mathbf{1}_g = (1, 1, \dots)$  the  $n_g \times 1$  vector of ones, and  $\mathbf{h}_g = (1, -1, 1, -1, \dots)$  the  $n_g \times 1$  vector containing alternating  $\pm 1$ .

The designs for the regressors are:

1. Normal with Clustered Dependence:  $\mathbf{X}_g \sim N(\mathbf{1}_g, \mathbf{I}_g + \mathbf{1}_g\mathbf{1}_g')$ .

<sup>9</sup>We focus on the case where the fixed effects are included for the same clusters as used

<sup>10</sup>These comments apply to the case where the fixed effects are applied at the same level as clustering, and not when fixed effects are applied at a different or more aggregate level.

2. LogNormal with Clustered Dependence:  $\mathbf{X}_g \sim \exp\left(\mathbf{N}\left(0, \mathbf{I}_g + \mathbf{1}_g \mathbf{1}_g'\right)\right)$ , recentered and rescaled so the elements have a mean of 1 and variance of 2.

In the first design, the regressors are normally distributed with a cluster-level random effects covariance matrix. In the second design, the regressors are log-normally distributed with the same cluster-level random effects covariance matrix. In both models the regressors are centered and scaled so to have the same first and second moments.

The designs for the errors are:

1. Homoskedastic with Clustered Dependence:  $\mathbf{e}_g \sim \mathbf{N}\left(0, \mathbf{I}_g + \mathbf{1}_g \mathbf{1}_g' + \mathbf{h}_g \mathbf{h}_g'\right)$ .
2. Heteroskedastic:  $\mathbf{e}_g \sim \mathbf{N}\left(0, \mathbf{I}_g + (\mathbf{X}_g - \mathbf{1}_g)(\mathbf{X}_g - \mathbf{1}_g)'\right)$ .

In the first design, the errors are normally distributed with a cluster-level factor covariance matrix. In the second design, they are conditionally heteroskedastic. The scalings are set so that the unconditional variances are the same in both designs.

The designs for the cluster sizes are

1. Homogeneous:  $n_g = 10$  for  $g = 1, \dots, G$ .
2. Heterogeneous:  $n_1 = n_2 = 5G/2 + 5$ , and  $n_g = 5$  for  $g = 3, \dots, G$ .

The homogeneous cluster size design sets all clusters equal to 10. The heterogeneous cluster size design puts between 1/4 and 1/3 of the observations in each of the first and second clusters, with the remaining spread among the other clusters. There are an average of 10 observations per cluster to match the homogeneous cluster size design.

The six designs combine these elements as described in Table 1. For each design, the number of clusters  $G$  is varied among  $\{6, 12, 40, 100\}$ . As there are an average of 10 observations per cluster, the associated total sample sizes are  $\{60, 120, 400, 1000\}$ .

Table 1: Simulation Designs

	Regressor	Error	Cluster Size
Design 1	Normal	Homoskedastic	Homogeneous
Design 2	LogNormal	Homoskedastic	Homogeneous
Design 3	LogNormal	Homoskedastic	Heterogeneous
Design 4	Normal	Heteroskedastic	Homogeneous
Design 5	LogNormal	Heteroskedastic	Homogeneous
Design 6	LogNormal	Heteroskedastic	Heterogeneous

For each simulation replication we estimate the coefficients by least squares. We use six methods to calculate standard errors: the five variances estimators described in the text:  $\text{CRVE}_1(\hat{\mathbf{v}}_1)$ ,  $\text{CRVE}_2(\hat{\mathbf{v}}_2)$ ,

the two conventional jackknife estimators  $\hat{v}_3$  and  $\hat{v}_4$ , our recommended jackknife estimator  $\hat{v}_5$ , and the nonparametric pairs cluster bootstrap<sup>11</sup> which we denote as  $\hat{v}_6$ .

We calculate confidence intervals for  $\beta$  using twelve methods.

The first six are conventional, based on the six standard errors and conventional student  $t$  critical values. Thus, given each standard error  $\hat{v}_j$ , we form the confidence interval  $\hat{\beta} \pm t_{G-1}^{0.975} \hat{v}_j$  where  $t_{G-1}^{0.975}$  is the 0.975 quantile of the  $t_{G-1}$  distribution. We use the  $t_{G-1}^{0.975}$  critical value as this is the current implementation in Stata for cluster-robust inference.

The next two intervals are adjusted  $t$  intervals. The first is that recommended by Bell and McCaffrey (2002) and equals  $\hat{\beta} \pm t_K^{0.975} \hat{v}_2$  where  $K$  is calculated<sup>12</sup> similar to (32) under the reference model  $\Sigma_g = I_{n_g}$ . The second interval is our adjusted interval (36) using the jackknife standard error.

The next two intervals are nonparametric pairs cluster bootstrap symmetric percentile- $t$  intervals, using the standard errors  $\hat{v}_1$  and  $\hat{v}_5$ , and 1000 bootstrap replications. The bootstrap samples are constructed by nonparametric pairs resampling, as described earlier. On each bootstrap sample we calculate the least squares estimate  $\hat{\beta}^*$  and the CRVE<sub>1</sub> and jackknife standard errors  $\hat{v}_1^*$  and  $\hat{v}_5^*$ . From the 1000 bootstrap samples we calculate the 95% quantiles  $\hat{c}_1^*$  and  $\hat{c}_5^*$  of the statistics  $|\hat{\beta}^* - \hat{\beta}|/\hat{v}_1^*$  and  $|\hat{\beta}^* - \hat{\beta}|/\hat{v}_5^*$ . The confidence intervals are  $\hat{\beta} \pm \hat{c}_1^* \hat{v}_1$  and  $\hat{\beta} \pm \hat{c}_5^* \hat{v}_5$ .

The final two intervals are wild cluster bootstrap symmetric percentile- $t$  intervals, using the standard errors  $\hat{v}_1$  and  $\hat{v}_5$ , and 1000 bootstrap replications. This is the method proposed by Cameron, Gelbach, and Miller (2008), Djogbenou, MacKinnon, and Nielsen (2019), and Canay, Santos, and Shaikh (2021) for hypothesis testing, and can be used to construct a confidence interval by test inversion<sup>13</sup>. First, the coefficients are re-estimated imposing known  $\beta$  to obtain restricted estimates  $(\tilde{\alpha}, \tilde{\beta} = \beta)$  and residuals  $\tilde{e}_g$ . Next, the clusters, regressors  $X_g$ , and restricted residuals  $\tilde{e}_g$  are held fixed. The bootstrap errors are generated as  $e_g^* = \zeta_g \tilde{e}_g$  where  $\zeta_g$  is an independent Rademacher variable (equals +1 and -1 each with probability 1/2), and  $Y_g^* = \tilde{\alpha} + X_g \tilde{\beta} + e_g^*$ . The bootstrap sample then consists of the observations  $(Y_g^*, X_g)$ . On each bootstrap sample we calculate the least squares estimate  $\hat{\beta}^*$  and the standard errors  $\hat{v}_1^*$  and  $\hat{v}_5^*$ . From the 1000 bootstrap samples we calculate the 95% quantiles  $\hat{c}_1^*(\beta)$  and  $\hat{c}_5^*(\beta)$  of the the statistics  $|\hat{\beta}^* - \beta|/\hat{v}_1^*$  and  $|\hat{\beta}^* - \beta|/\hat{v}_5^*$ . The confidence intervals<sup>14</sup> consist of all  $\beta$  such that  $|\hat{\beta} - \beta|/\hat{v}_1 \leq \hat{c}_1^*(\beta)$  and  $|\hat{\beta} - \beta|/\hat{v}_5 \leq \hat{c}_5^*(\beta)$ , respectively.

We compute the actual coverage probability of these nominal 95% intervals by simulation with 20,000 replications. These estimates are precise, as their standard errors are all less than 0.003.

We report the results for the baseline regression model in Table 2. The top block reports the results for  $G = 6$ . The first six columns are for the conventional confidence intervals. We can see that the conventional CRVE<sub>1</sub> confidence interval has substantial under-coverage in most designs. The worst-case is

<sup>11</sup>This is the standard implementation of the bootstrap for clustered observations. Each bootstrap sample is constructed by resampling  $G$  clusters  $(Y_g, X_g)$  with replacement from the original sample of clusters. Least squares estimation is applied to the bootstrap sample. The bootstrap standard error is the empirical standard deviation of the bootstrap least squares estimates. 1000 bootstrap replications were made in each simulation replication. For details see Cameron, Gelbach, and Miller (2008).

<sup>12</sup>The CRVE<sub>2</sub> standard error  $\hat{v}_2$  and Bell-McCaffrey adjusted interval are calculated using the `dfadjust` R package of Kolesár (2023), which is identical to the implementation in Stata 18.

<sup>13</sup>MacKinnon, Nielsen and Webb (2023b) review several variants of the wild cluster bootstrap. Our implementation corresponds to their WCR-V method.

<sup>14</sup>To assess the coverage rate, it is sufficient to do the calculation for the true value of  $\beta$ .

Table 2: Baseline Regression Model. Coverage of Nominal 95% Confidence Intervals for  $\beta$

Cr. Value	Conventional $t_{G-1}$						Adjusted $t_K$		Pairs Boot		Wild Boot	
St. Error	$\hat{v}_1$	$\hat{v}_2$	$\hat{v}_3$	$\hat{v}_4$	$\hat{v}_5$	$\hat{v}_6$	$\hat{v}_2$	$\hat{v}_5$	$\hat{v}_1$	$\hat{v}_5$	$\hat{v}_1$	$\hat{v}_5$
$G = 6$												
Design 1	0.91	0.93	0.95	0.95	0.96	0.94	0.95	0.96	0.96	0.96	0.93	0.94
Design 2	0.85	0.91	0.95	0.95	0.96	0.98	0.99	0.99	0.93	0.96	0.96	0.95
Design 3	0.83	0.90	0.95	0.95	0.96	0.99	0.99	0.99	0.93	0.97	0.95	0.95
Design 4	0.89	0.91	0.93	0.93	0.95	0.91	0.94	0.95	0.95	0.96	0.93	0.94
Design 5	0.64	0.74	0.88	0.89	0.90	0.92	0.93	0.95	0.80	0.92	0.91	0.94
Design 6	0.61	0.72	0.88	0.88	0.90	0.92	0.92	0.94	0.81	0.93	0.91	0.94
$G = 12$												
Design 1	0.92	0.93	0.95	0.95	0.95	0.93	0.95	0.95	0.95	0.95	0.95	0.95
Design 2	0.83	0.89	0.94	0.94	0.94	0.97	0.98	0.99	0.89	0.93	0.96	0.96
Design 3	0.80	0.88	0.94	0.94	0.95	0.98	0.99	0.99	0.90	0.95	0.95	0.95
Design 4	0.91	0.92	0.93	0.93	0.94	0.91	0.94	0.94	0.95	0.95	0.94	0.95
Design 5	0.63	0.74	0.87	0.87	0.88	0.86	0.93	0.95	0.82	0.92	0.92	0.93
Design 6	0.59	0.71	0.87	0.87	0.88	0.86	0.93	0.95	0.81	0.92	0.92	0.94
$G = 40$												
Design 1	0.94	0.94	0.95	0.95	0.95	0.94	0.95	0.95	0.95	0.95	0.95	0.95
Design 2	0.86	0.90	0.93	0.93	0.93	0.95	0.97	0.98	0.89	0.92	0.95	0.96
Design 3	0.81	0.88	0.94	0.94	0.94	0.95	0.98	0.99	0.90	0.95	0.94	0.96
Design 4	0.93	0.94	0.94	0.94	0.95	0.93	0.94	0.94	0.95	0.95	0.95	0.95
Design 5	0.70	0.79	0.88	0.88	0.89	0.84	0.93	0.95	0.90	0.94	0.94	0.94
Design 6	0.64	0.75	0.87	0.87	0.88	0.83	0.93	0.95	0.87	0.94	0.93	0.94
$G = 100$												
Design 1	0.95	0.95	0.95	0.95	0.95	0.94	0.95	0.95	0.95	0.95	0.95	0.95
Design 2	0.89	0.91	0.93	0.93	0.93	0.95	0.97	0.97	0.91	0.93	0.95	0.95
Design 3	0.81	0.88	0.94	0.94	0.95	0.92	0.97	0.98	0.92	0.96	0.94	0.96
Design 4	0.94	0.94	0.95	0.95	0.95	0.94	0.95	0.95	0.95	0.95	0.95	0.95
Design 5	0.77	0.84	0.90	0.90	0.90	0.85	0.93	0.95	0.92	0.95	0.95	0.95
Design 6	0.69	0.78	0.89	0.89	0.89	0.83	0.93	0.95	0.89	0.94	0.94	0.95

Design 6, where the interval has only 61% coverage. The  $CRVE_2$  confidence interval has slightly better coverage, but still substantially under-covers in most designs. The jackknife and bootstrap confidence intervals have better coverage, but undercover in Designs 5 & 6.

The next two columns are for the confidence intervals using adjusted critical values. The adjusted  $CRVE_2$  interval has improved coverage over the conventional  $CRVE_2$  interval, but slightly undercovers (92%) for Designs 5 and 6. The adjusted jackknife interval has excellent coverage, with the coverage probability exceeding 94% in all designs.

The results for the pairs bootstrap- $t$  methods are reported in the next two columns. The intervals based on the  $CRVE_1$  standard error generally undercover (and severely for Designs 5 & 6). The intervals based on the jackknife standard errors have much better coverage, with coverage exceeding 92% in all designs.

Table 3: Regression with Treatment Dummy. Coverage of Nominal 95% Confidence Intervals for  $\beta$

Cr. Value	Conventional $t_{G-1}$						Adjusted $t_K$		Pairs Boot		Wild Boot	
St. Error	$\hat{v}_1$	$\hat{v}_2$	$\hat{v}_3$	$\hat{v}_4$	$\hat{v}_5$	$\hat{v}_6$	$\hat{v}_2$	$\hat{v}_5$	$\hat{v}_1$	$\hat{v}_5$	$\hat{v}_1$	$\hat{v}_5$
$G = 6$												
Design 1	0.91	0.93	0.94	0.94	0.96	0.94	0.95	0.95	0.95	0.96	0.93	0.94
Design 2	0.85	0.91	0.92	0.93	0.97	0.93	0.99	0.99	0.92	0.96	0.96	0.96
Design 3	0.83	0.90	0.87	0.89	0.96	0.87	0.99	0.98	0.90	0.95	0.95	0.95
Design 4	0.89	0.91	0.91	0.92	0.95	0.90	0.94	0.93	0.93	0.94	0.93	0.94
Design 5	0.64	0.74	0.81	0.83	0.90	0.84	0.93	0.94	0.81	0.92	0.92	0.94
Design 6	0.61	0.72	0.74	0.76	0.89	0.74	0.92	0.93	0.79	0.91	0.92	0.94
$G = 12$												
Design 1	0.92	0.93	0.94	0.94	0.96	0.93	0.95	0.95	0.95	0.95	0.95	0.95
Design 2	0.83	0.89	0.92	0.93	0.94	0.95	0.99	0.99	0.89	0.93	0.96	0.96
Design 3	0.82	0.89	0.89	0.90	0.95	0.89	0.99	0.99	0.88	0.93	0.95	0.96
Design 4	0.91	0.92	0.93	0.93	0.94	0.91	0.94	0.94	0.94	0.95	0.94	0.94
Design 5	0.63	0.74	0.84	0.85	0.88	0.83	0.93	0.95	0.81	0.91	0.92	0.93
Design 6	0.59	0.71	0.76	0.77	0.87	0.73	0.93	0.94	0.79	0.90	0.92	0.94
$G = 40$												
Design 1	0.94	0.94	0.95	0.95	0.95	0.94	0.95	0.95	0.95	0.95	0.95	0.95
Design 2	0.86	0.90	0.93	0.93	0.93	0.95	0.97	0.98	0.89	0.92	0.95	0.96
Design 3	0.82	0.89	0.91	0.91	0.94	0.89	0.98	0.98	0.88	0.93	0.94	0.96
Design 4	0.93	0.94	0.94	0.94	0.95	0.93	0.94	0.94	0.95	0.95	0.95	0.95
Design 5	0.70	0.79	0.88	0.88	0.89	0.83	0.93	0.95	0.90	0.94	0.94	0.94
Design 6	0.64	0.75	0.79	0.79	0.88	0.73	0.93	0.95	0.84	0.92	0.93	0.94
$G = 100$												
Design 1	0.95	0.95	0.95	0.95	0.95	0.94	0.95	0.95	0.95	0.95	0.95	0.95
Design 2	0.89	0.91	0.93	0.93	0.94	0.95	0.97	0.97	0.91	0.93	0.95	0.95
Design 3	0.83	0.89	0.92	0.92	0.94	0.89	0.97	0.98	0.89	0.94	0.93	0.96
Design 4	0.94	0.94	0.95	0.95	0.95	0.94	0.95	0.95	0.95	0.95	0.95	0.95
Design 5	0.77	0.84	0.90	0.90	0.90	0.85	0.93	0.95	0.92	0.95	0.95	0.95
Design 6	0.69	0.78	0.80	0.80	0.89	0.74	0.93	0.95	0.86	0.92	0.94	0.95

The results for the wild bootstrap- $t$  are presented in the final two columns. The coverage rates are excellent when the jackknife standard errors are used.

The following blocks are for  $G = 12$ ,  $G = 40$ , and  $G = 100$ . Qualitatively, the results are similar to the  $G = 6$  case. For some designs and methods the coverage probabilities improve slightly as  $G$  increases. A notable exception is the conventional interval with bootstrap standard errors whose coverage probabilities worsen as  $G$  increases in Designs 3, 5, and 6. In general, the best performance is obtained by our adjusted jackknife interval and the wild bootstrap using jackknife standard errors. A close competitor is the pairs bootstrap percentile- $t$  using jackknife standard errors. No other method has reliable coverage across the designs and the four sample sizes.

If we compare performance by standard error method, within each type of inference method (conventional, adjusted, pairs bootstrap- $t$ , and wild bootstrap- $t$ ), we systematically see that methods based

Table 4: Regression with Treatment Dummy. Coverage of Nominal 95% Confidence Intervals for  $\gamma$

Cr. Value	Conventional $t_{G-1}$						Adjusted $t_K$		Pairs Boot		Wild Boot	
St. Error	$\hat{v}_1$	$\hat{v}_2$	$\hat{v}_3$	$\hat{v}_4$	$\hat{v}_5$	$\hat{v}_6$	$\hat{v}_2$	$\hat{v}_5$	$\hat{v}_1$	$\hat{v}_5$	$\hat{v}_1$	$\hat{v}_5$
<b>G = 6</b>												
Design 1	0.64	0.66	0.68	0.68	1.00	0.66	0.73	1.00	0.83	1.00	1.00	1.00
Design 2	0.64	0.66	0.66	0.67	1.00	0.64	0.72	1.00	0.85	1.00	1.00	1.00
Design 3	0.59	0.64	0.67	0.68	1.00	0.66	0.76	1.00	0.82	1.00	0.98	0.98
Design 4	0.71	0.75	0.77	0.77	1.00	0.75	0.82	1.00	0.83	0.99	1.00	1.00
Design 5	0.69	0.72	0.74	0.75	1.00	0.72	0.76	1.00	0.87	1.00	1.00	1.00
Design 6	0.72	0.77	0.79	0.80	1.00	0.79	0.86	1.00	0.87	1.00	0.99	0.99
<b>G = 12</b>												
Design 1	0.49	0.50	0.51	0.51	1.00	0.49	0.52	1.00	0.57	1.00	1.00	1.00
Design 2	0.47	0.48	0.48	0.49	1.00	0.47	0.51	1.00	0.53	1.00	1.00	1.00
Design 3	0.47	0.51	0.55	0.56	1.00	0.50	0.59	1.00	0.59	1.00	0.99	0.98
Design 4	0.61	0.64	0.66	0.66	1.00	0.61	0.66	1.00	0.73	1.00	1.00	1.00
Design 5	0.54	0.58	0.62	0.62	1.00	0.56	0.60	1.00	0.66	1.00	1.00	1.00
Design 6	0.64	0.70	0.74	0.75	1.00	0.70	0.77	1.00	0.76	1.00	0.99	0.99
<b>G = 40</b>												
Design 1	0.28	0.28	0.29	0.29	1.00	0.28	0.29	1.00	0.29	1.00	1.00	1.00
Design 2	0.26	0.27	0.27	0.27	1.00	0.26	0.27	1.00	0.27	1.00	1.00	1.00
Design 3	0.33	0.38	0.43	0.43	1.00	0.35	0.43	1.00	0.48	1.00	0.99	0.99
Design 4	0.44	0.45	0.45	0.45	1.00	0.43	0.45	1.00	0.48	1.00	1.00	1.00
Design 5	0.36	0.40	0.44	0.44	1.00	0.39	0.41	1.00	0.47	1.00	1.00	1.00
Design 6	0.56	0.62	0.67	0.68	1.00	0.60	0.68	1.00	0.69	1.00	1.00	1.00
<b>G = 100</b>												
Design 1	0.17	0.17	0.17	0.17	1.00	0.17	0.17	1.00	0.18	1.00	1.00	1.00
Design 2	0.16	0.16	0.16	0.16	1.00	0.16	0.16	1.00	0.16	1.00	1.00	1.00
Design 3	0.28	0.33	0.38	0.38	1.00	0.29	0.37	1.00	0.47	1.00	0.99	0.98
Design 4	0.30	0.31	0.31	0.31	1.00	0.30	0.31	1.00	0.32	1.00	1.00	1.00
Design 5	0.29	0.32	0.36	0.36	1.00	0.31	0.32	1.00	0.39	1.00	1.00	1.00
Design 6	0.54	0.60	0.64	0.64	1.00	0.54	0.65	1.00	0.69	1.00	1.00	1.00

on jackknife standard errors perform better than any other standard error method. Specifically, this holds whether inference uses conventional student  $t$  critical values, adjusted critical values, pairs bootstrap- $t$  critical values, or wild bootstrap- $t$  critical values. This is strong evidence favoring jackknife standard errors, regardless of the inference method.

We expand the analysis by examining a model with clusterwise noninvertibility. This is  $Y_g = \alpha + X_g\beta + D_g\gamma + e_g$  with  $D_g$  a dummy indicator for the first cluster. This model is cluster-level treatment with a single treated cluster. The regression is clusterwise noninvertible as the least squares estimator is undefined when the first cluster is omitted. In this model we examine confidence intervals for both  $\beta$  and  $\gamma$ . For this model we consider the same six designs as in the baseline model for the distributions of the regressors, errors, and cluster sizes.

At this point we need to discuss our implementation of the pairs bootstrap. As  $D_g$  is non-zero only

Table 5: Regression with Strongly Skewed Errors. Coverage of Nominal 95% Confidence Intervals for  $\beta$

Cr. Value	Conventional $t_{G-1}$						Adjusted $t_K$		Pairs Boot		Wild Boot	
St. Error	$\hat{v}_1$	$\hat{v}_2$	$\hat{v}_3$	$\hat{v}_4$	$\hat{v}_5$	$\hat{v}_6$	$\hat{v}_2$	$\hat{v}_5$	$\hat{v}_1$	$\hat{v}_5$	$\hat{v}_1$	$\hat{v}_5$
$G = 6$												
Design 1	0.92	0.94	0.96	0.96	0.97	0.95	0.96	0.97	0.96	0.97	0.94	0.94
Design 2	0.86	0.92	0.96	0.96	0.97	0.99	0.99	0.99	0.93	0.97	0.96	0.95
Design 3	0.84	0.90	0.96	0.96	0.97	0.99	0.99	0.99	0.93	0.97	0.94	0.95
Design 4	0.86	0.88	0.90	0.90	0.92	0.88	0.91	0.92	0.93	0.94	0.91	0.91
Design 5	0.64	0.74	0.87	0.87	0.89	0.91	0.92	0.94	0.79	0.91	0.91	0.93
Design 6	0.61	0.72	0.87	0.88	0.89	0.93	0.92	0.93	0.80	0.92	0.91	0.94
$G = 12$												
Design 1	0.93	0.94	0.95	0.95	0.96	0.94	0.95	0.96	0.96	0.96	0.95	0.95
Design 2	0.85	0.91	0.94	0.94	0.95	0.97	0.99	0.99	0.90	0.94	0.96	0.96
Design 3	0.81	0.89	0.94	0.95	0.95	0.98	0.99	0.99	0.90	0.95	0.95	0.95
Design 4	0.88	0.90	0.91	0.91	0.92	0.89	0.91	0.92	0.93	0.93	0.92	0.92
Design 5	0.62	0.72	0.85	0.85	0.86	0.85	0.92	0.94	0.80	0.90	0.91	0.93
Design 6	0.58	0.69	0.84	0.85	0.86	0.85	0.91	0.94	0.79	0.90	0.91	0.93
$G = 40$												
Design 1	0.94	0.95	0.95	0.95	0.95	0.94	0.95	0.95	0.96	0.96	0.95	0.95
Design 2	0.87	0.91	0.94	0.94	0.94	0.96	0.98	0.98	0.90	0.93	0.96	0.96
Design 3	0.82	0.89	0.94	0.94	0.94	0.95	0.98	0.99	0.91	0.95	0.95	0.96
Design 4	0.92	0.92	0.93	0.93	0.93	0.91	0.93	0.93	0.94	0.94	0.93	0.93
Design 5	0.68	0.77	0.85	0.85	0.86	0.81	0.91	0.93	0.86	0.91	0.92	0.92
Design 6	0.62	0.72	0.84	0.84	0.85	0.80	0.90	0.93	0.84	0.91	0.91	0.92
$G = 100$												
Design 1	0.94	0.95	0.95	0.95	0.95	0.94	0.95	0.95	0.95	0.95	0.95	0.95
Design 2	0.89	0.92	0.93	0.93	0.94	0.95	0.97	0.97	0.91	0.93	0.95	0.95
Design 3	0.82	0.89	0.94	0.94	0.94	0.93	0.97	0.98	0.91	0.95	0.94	0.95
Design 4	0.94	0.94	0.94	0.94	0.94	0.93	0.94	0.94	0.95	0.95	0.94	0.94
Design 5	0.74	0.81	0.87	0.87	0.87	0.83	0.91	0.93	0.89	0.93	0.92	0.92
Design 6	0.66	0.75	0.86	0.86	0.86	0.80	0.90	0.93	0.86	0.92	0.91	0.92

for one cluster, a high percentage (about 37%) of bootstrap samples have a singular regressor matrix. In this context it is not clear how the bootstrap should treat these sample draws. We follow the recommendation of Shao and Tu (1995) and implementation in Stata which discards these bootstrap samples. Thus bootstrap standard errors and critical values are calculated from the bootstrap samples which have nonsingular regressor matrices.

Table 3 presents the coverage rates for confidence intervals for  $\beta$ . In general, the results are similar to those of Table 2 but with two important differences. First, in Table 3 we see a meaningful divergence in performance between the intervals based on conventional jackknife standard errors  $\hat{v}_3$  and  $\hat{v}_4$  and those based on our recommended jackknife standard errors  $\hat{v}_5$ . In Table 2, these three methods were nearly identical; in Table 3 we can see that the two conventional intervals exhibit substantial undercoverage in most designs (which diminish as  $G$  increases). This difference is due to the treatment of noninvert-

Table 6: Baseline Regression. Average Interval Length

Cr. Value	Conventional $t_{G-1}$						Adjusted $t_K$		Pairs Boot	
St. Error	$\hat{v}_1$	$\hat{v}_2$	$\hat{v}_3$	$\hat{v}_4$	$\hat{v}_5$	$\hat{v}_6$	$\hat{v}_2$	$\hat{v}_5$	$\hat{v}_1$	$\hat{v}_5$
$G = 6$										
Design 1	1.00	1.12	1.28	1.29	1.41	1.05	1.26	1.43	1.85	2.08
Design 2	1.69	2.09	3.17	3.28	3.60	3.32	3.96	5.30	3.15	5.65
Design 3	1.71	2.18	3.51	3.66	4.01	3.92	4.46	5.87	3.77	7.24
Design 4	0.94	1.04	1.17	1.17	1.29	0.97	1.18	1.31	1.91	2.25
Design 5	0.83	1.08	1.71	1.78	1.95	1.46	2.15	2.96	2.94	6.26
Design 6	0.83	1.10	1.83	1.91	2.09	1.65	2.35	3.15	3.50	7.60
$G = 12$										
Design 1	0.67	0.71	0.76	0.76	0.80	0.67	0.76	0.80	0.85	0.90
Design 2	0.86	1.02	1.37	1.39	1.45	1.40	1.84	2.39	1.22	1.91
Design 3	0.93	1.15	1.67	1.70	1.77	1.67	2.26	2.97	1.58	2.70
Design 4	0.66	0.70	0.76	0.76	0.79	0.66	0.76	0.79	0.92	1.01
Design 5	0.63	0.82	1.23	1.25	1.31	0.93	1.64	2.30	2.29	4.65
Design 6	0.63	0.85	1.33	1.35	1.42	1.01	1.81	2.50	2.70	5.62
$G = 40$										
Design 1	0.36	0.37	0.38	0.38	0.38	0.36	0.38	0.38	0.39	0.39
Design 2	0.37	0.41	0.49	0.49	0.50	0.48	0.61	0.75	0.47	0.61
Design 3	0.45	0.55	0.73	0.73	0.74	0.64	0.91	1.19	0.75	1.12
Design 4	0.38	0.39	0.40	0.40	0.41	0.38	0.40	0.41	0.43	0.44
Design 5	0.55	0.69	0.93	0.93	0.94	0.69	1.15	1.56	1.69	2.80
Design 6	0.56	0.72	1.04	1.05	1.06	0.76	1.32	1.85	1.98	3.66
$G = 100$										
Design 1	0.23	0.23	0.24	0.24	0.24	0.23	0.24	0.24	0.24	0.24
Design 2	0.22	0.23	0.26	0.26	0.26	0.26	0.30	0.35	0.26	0.30
Design 3	0.31	0.37	0.47	0.47	0.48	0.39	0.54	0.71	0.55	0.80
Design 4	0.25	0.25	0.25	0.25	0.26	0.25	0.25	0.26	0.26	0.26
Design 5	0.51	0.60	0.75	0.75	0.76	0.58	0.87	1.12	1.23	1.85
Design 6	0.51	0.65	0.89	0.89	0.90	0.65	1.07	1.48	1.57	2.79

ible clusters. Second, the conventional interval using bootstrap standard errors exhibits a substantial deterioration in coverage relative to the clusterwise invertible case of Table 2.

Next, we examine Table 4, which presents coverage rates for the dummy variable coefficient  $\gamma$  in the regression model  $\mathbf{Y}_g = \alpha + \mathbf{X}_g\beta + \mathbf{D}_g\gamma + \mathbf{e}_g$ . This is a treacherous context, as we are essentially making inference for a coefficient based on a single cluster. The results in Table 4 reveal that the conventional inference methods fail<sup>15</sup>, and worsen as the sample size increases. The Bell-McCaffrey adjusted interval, conventional bootstrap interval, and pairs bootstrap percentile- $t$  using CRVE<sub>1</sub> standard errors similarly fail. The exceptions (which generally produce coverage rates of 100%) are any intervals (conventional, adjusted, bootstrap percentile- $t$ , and wild bootstrap- $t$ ) using the jackknife  $\hat{v}_5$  standard error, and both wild bootstrap percentile- $t$  intervals. The reason why the intervals based on the jackknife standard error

<sup>15</sup>The failure of conventional methods in this context is well known.



$\hat{v}_5$  have coverage rates of 100% is because in this context  $\hat{v}_5$  tends to approximately equal the coefficient estimate  $\hat{\gamma}$ , so the confidence interval always covers the true value of 0. Coverage is conservative, but at least there is no tendency towards false significance. The message of Table 4 is that our recommended jackknife standard errors and confidence intervals will not produce misleading significance, even in the most extreme context of inference on a dummy indicator for a single treated cluster.

Our inference theory is developed under the assumption of normally distributed errors, which raises the question if the coverage rates are sensitive to departures to normality. To explore this question, we repeat the analysis of the baseline model, but with the errors drawn from a skewed heavy-tailed distribution. Specifically, we use the “strongly skewed” distribution displayed in Figure 3.7(b) of Hansen (2022) which is a 9-component normal mixture distribution with a skew of 1.34 and kurtosis of 6.7. We sample the regression errors from this skewed distribution, with the same cluster covariances as in the baseline model. The results are displayed in Table 5. Comparing Tables 2 and 5, we can see that the results are qualitatively similar, with reduced coverage accuracy of several methods in Table 5. The relative rankings of the methods, however, are unchanged.

We finally address the accuracy of the confidence intervals by calculating their average length in the baseline model. We do so for all methods except for the wild bootstrap. This is because calculation of the full wild bootstrap confidence interval must be done by numerical test statistic inversion, which is computationally demanding<sup>16</sup>. We report the results in Table 6. The results are a bit difficult to compare across methods, as most methods have substantial undercoverage. However, we can make the following observations. First, if we compare the conventional methods (the first six), we can see that while our proposed jackknife-based intervals have the longest length, the six are of rough similar length, and this holds despite the fact that the other intervals have substantial under-coverage. Second, our proposed adjusted confidence interval has longer average length than the conventional methods, but again of rough similarity. Third, the relative differences decrease as the sample size increases. Fourth, the pairs bootstrap- $t$  intervals can be substantially wider than our proposed intervals, despite having similar coverage probabilities. In general, the results of Table 6 suggest that the proposed methods have reasonable accuracy.

In summary, we are able to draw the following conclusions from the simulation evidence. First, the standard error which produces confidence intervals with the best coverage rates is  $\hat{v}_5$ . Second, the simple confidence interval  $\hat{\beta} \pm t_{G-1}^{1-\alpha/2} \hat{v}_5$  has good coverage rates in many contexts, but can under-cover in extreme designs. Third, the adjusted confidence interval  $\hat{\beta} \pm t_K^{1-\alpha/2} \hat{v}_5 / a$  has excellent (but conservative) coverage in all contexts examined. Fourth, excellent coverage is also attained by the pairs and wild bootstrap- $t$  intervals, but only if combined with the jackknife standard error  $\hat{v}_5$ . Fifth, issues such as clusterwise invertibility should not be handled by *ad hoc* computational implementations, but rather by methods justified by theoretical insight. In particular, the jackknife should be implemented without the discarding of iterations with noninvertible design matrices.

---

<sup>16</sup>Our calculation of Table 6 using parallel processing took 3 days. Adding the wild bootstrap would increase the computation time by an order of magnitude.

## 10 Empirical Illustration

We illustrate the applicability of the methods with an empirical example. We follow Canay, Santos, and Shaikh (2021) by revisiting an application by Meng, Qian, and Yared (2015) into the causes of the Chinese Great Famine between 1958 and 1960. Their regressions (Table 2 of Meng-Qian-Yared) take the form  $Y = Z_1\beta_1 + Z_2\beta_2 + W'\gamma + e$  for  $G = 19$  provinces between 1953 and 1982, where  $Y$  equals the log of deaths in the province,  $Z_1$  equals the log of predicted grain production,  $Z_2$  equals the product of  $Z_1$  and an indicator for a famine year, and  $W$  are other controls. The focus is on the coefficient sum  $\beta_1 + \beta_2$ . The authors report six specifications which vary the sample period (1953-1982 vs 1953-1965), the provinces (19 vs 23 provinces), and replacing predicted with reported grain production.

Canay, Santos, and Shaikh (2021) use this application to illustrate hypothesis testing using the cluster wild bootstrap. In contrast, we are interested in standard error calculation and confidence interval construction, in addition to hypothesis testing. Following these authors, we cluster by province.

We estimate the same six regression specifications as Meng, Qian, and Yared (2015) and focus on the coefficient sum  $\beta_1 + \beta_2$ . In Table 7 we report the least squares estimates  $\hat{\beta}_1 + \hat{\beta}_2$  plus three standard errors:  $CRVE_1$ ,  $CRVE_2$ , and our recommended jackknife standard errors. What you can see from the table is that in some of the specifications there are considerable differences between the three standard errors, and in particular between the jackknife and the other two. The discrepancies between the  $CRVE_1$  and jackknife standard errors range from 10% (in specification #1) to 66% (in specification #3). These are large and substantial differences.

We next construct 95% adjusted confidence intervals for the coefficient sum  $\beta_1 + \beta_2$ . We start by calculating the adjustment coefficients  $K$  and  $a$  for each of the six specifications, and report these coefficients in Table 7. The values for the adjusted degree-of-freedom  $K$  range between 4 and 6, which are all small. The values for the scale adjustment  $a$  range between 1.17 and 1.26. Together, these are used to construct the confidence intervals, which are reported in the Table.

Take, for example, specification #1, where  $\hat{\beta}_1 + \hat{\beta}_2 = 0.141$ ,  $\hat{v}_5 = 0.066$ ,  $K = 4.18$ , and  $a = 1.21$ . The 95% critical value from the  $t$  distribution with  $K = 4.18$  degrees of freedom is 2.73. The 95% confidence interval is therefore  $0.141 \pm 2.73 \times 0.066/1.21 = [-.01, .29]$ , as reported.

The confidence intervals are wide, indicating uncertainty about the value of the coefficient sum. The intervals do not vary greatly across the six specifications, indicating that the result is reasonably robust to the specification.

We also construct and report adjusted p-values for t-tests of the hypothesis  $\beta_1 + \beta_2 = 0$ . The t-statistic using the jackknife standard error is  $0.141/0.066 = 2.16$ . The adjusted p-value is  $1 - F(1.21^2 \times 2.16^2; 1, 4.18) = 0.058$ , as reported. None of the six p-values are significant at the 5% level. This contrasts with the p-values reported by Meng, Qian, and Yared (2015), which were all statistically significant, some greatly so. Several of our p-values are similar to those calculated by the wild cluster bootstrap as reported by Canay, Santos, and Shaikh (2021). For example, for the baseline specification #1, our p-value of 0.058 is nearly identical to their wild studentized p-value of 0.061.

Table 7: China's Great Famine, 1959-1961

	Dependent variable: log deaths in year $t + 1$					
	Constructed grain production				Reported grain production	
	19 provinces		23 provinces		19 provinces	
	1953-1982	1953-1965	1953-1982	1953-1965	1953-1982	1953-1965
	(1)	(2)	(3)	(4)	(5)	(6)
$\hat{\beta}_1 + \hat{\beta}_2$	0.141	0.098	0.115	0.094	0.113	0.089
CRVE <sub>1</sub>	(0.060)	(0.053)	(0.037)	(0.037)	(0.063)	(0.059)
CRVE <sub>2</sub>	[0.061]	[0.056]	[0.039]	[0.037]	[0.068]	[0.063]
Jackknife	<0.066>	<0.066>	<0.061>	<0.052>	<0.079>	<0.073>
$K$	4.18	3.95	5.34	5.03	5.18	5.97
$a$	1.21	1.26	1.21	1.23	1.18	1.17
Interval	[-0.01, 0.29]	[-0.05, 0.24]	[-0.01, 0.24]	[-0.01, 0.20]	[-0.06, 0.28]	[-0.06, 0.24]
p-value	0.058	0.135	0.069	0.076	0.151	0.200

Notes: All regressions include log total population, log urban population, and year fixed effects. CRVE<sub>1</sub>, CRVE<sub>2</sub>, and jackknife standard errors for  $\hat{\beta}_1 + \hat{\beta}_2$ , clustered by province, in parenthesis, square brackets, and angle brackets, respectively.  $K$  and  $a$  are the adjustment parameters for the confidence interval for  $\hat{\beta}_1 + \hat{\beta}_2$ . The p-value is for the t-test of the hypothesis  $\beta_1 + \beta_2 = 0$ .

## 11 Conclusion

Heteroskedasticity-consistent and cluster-robust standard errors are routinely reported in applied econometric practice. It is prudent for the profession to coalesce on simple yet well-behaved methods which produce reliable inference across reasonable estimation settings. It is our contention that jackknife variance estimators are superior to conventional (EHW and CRVE<sub>1</sub>) estimators, based on our analysis of worst-case downward bias and confidence interval coverage rates. They are also computationally simple to implement.

## 12 Technical Proofs

**Proof of Theorem 1:** We first show that (14) holds for definitions (12) and (13). The estimator (12) with any generalized inverse is a minimizer of the least-squares criterion, so solves the first order condition

$$\left( \mathbf{X}'\mathbf{X} - \mathbf{X}'_g\mathbf{X}_g \right) \tilde{\beta}_{-g} = \left( \mathbf{X}'\mathbf{Y} - \mathbf{X}'_g\mathbf{Y}_g \right).$$

Pre-multiplying by  $(\mathbf{X}'\mathbf{X})^{-1}$ , rearranging, and using (13), we obtain

$$\begin{aligned} \tilde{\beta}_{-g} &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y} - (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'_g\mathbf{Y}_g + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'_g\mathbf{X}_g\tilde{\beta}_{-g} \\ &= \hat{\beta} - (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'_g(\mathbf{Y}_g - \mathbf{X}_g\tilde{\beta}_{-g}) \\ &= \hat{\beta} - (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'_g\tilde{\mathbf{e}}_g, \end{aligned}$$

which is (14) as claimed.

By definition (13) and model (1), the prediction errors equal

$$\tilde{\mathbf{e}}_g = \mathbf{Y}_g - \mathbf{X}_g \tilde{\boldsymbol{\beta}}_{-g} = \mathbf{e}_g - \mathbf{X}_g (\tilde{\boldsymbol{\beta}}_{-g} - \boldsymbol{\beta}). \quad (39)$$

Squaring and expanding,

$$\tilde{\mathbf{e}}_g \tilde{\mathbf{e}}_g' = \mathbf{e}_g \mathbf{e}_g' - \mathbf{e}_g (\tilde{\boldsymbol{\beta}}_{-g} - \boldsymbol{\beta})' \mathbf{X}_g' - \mathbf{X}_g (\tilde{\boldsymbol{\beta}}_{-g} - \boldsymbol{\beta}) \mathbf{e}_g' + \mathbf{X}_g (\tilde{\boldsymbol{\beta}}_{-g} - \boldsymbol{\beta}) (\tilde{\boldsymbol{\beta}}_{-g} - \boldsymbol{\beta})' \mathbf{X}_g' \quad (40)$$

$$\geq \mathbf{e}_g \mathbf{e}_g' - \mathbf{e}_g (\tilde{\boldsymbol{\beta}}_{-g} - \boldsymbol{\beta})' \mathbf{X}_g' - \mathbf{X}_g (\tilde{\boldsymbol{\beta}}_{-g} - \boldsymbol{\beta}) \mathbf{e}_g'. \quad (41)$$

The inequality holds because the final term in (40) takes the form  $\mathbf{A}\mathbf{A}'$  and is thus positive semi-definite.

The first term in (41) has expectation  $\boldsymbol{\Sigma}_g$ . Observe that  $\mathbf{e}_g$  is independent of  $\tilde{\boldsymbol{\beta}}_{-g} - \boldsymbol{\beta}$  and mean zero, so the expectation of the second and third terms in (41) equals zero. We deduce that

$$\mathbb{E} [\tilde{\mathbf{e}}_g \tilde{\mathbf{e}}_g'] \geq \mathbb{E} [\mathbf{e}_g \mathbf{e}_g'] = \boldsymbol{\Sigma}_g. \quad (42)$$

Using expression (15) and inequality (42),

$$\begin{aligned} \mathbb{E} [\widehat{\mathbf{V}}_5] &= (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{g=1}^G \mathbf{X}_g' \mathbb{E} [\tilde{\mathbf{e}}_g \tilde{\mathbf{e}}_g'] \mathbf{X}_g \right) (\mathbf{X}'\mathbf{X})^{-1} \\ &\geq (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{g=1}^G \mathbf{X}_g' \boldsymbol{\Sigma}_g \mathbf{X}_g \right) (\mathbf{X}'\mathbf{X})^{-1} \\ &= \mathbf{V}. \end{aligned}$$

This is (16). ■

**Proof of Theorem 2, equation (17):** Without loss of generality, normalize  $\mathbf{X}'\mathbf{X} = \mathbf{I}_k$  and  $R'R = 1$ . Set  $\mathbf{Z}_g = \mathbf{X}_g R$  and  $c_g = \mathbf{Z}_g' \mathbf{Z}_g$ , and observe that  $\sum_{g=1}^G c_g = 1$ .

The CRVE<sub>1</sub> estimator can be written as

$$\widehat{v}_1^2 = d \sum_{g=1}^G \mathbf{Z}_g' \widehat{\mathbf{e}}_g \widehat{\mathbf{e}}_g' \mathbf{Z}_g \quad (43)$$

where  $d = G(n-1)/(G-1)(n-k)$ . The restriction  $(\mathbf{X}, \boldsymbol{\Sigma}) \in \mathcal{F}_0^*$  imposes  $\boldsymbol{\Sigma}_g = \sigma^2 \mathbf{I}_{n_g}$ , under which we can calculate that  $v^2 = \sigma^2$  and

$$\mathbb{E} [\widehat{\mathbf{e}}_g \widehat{\mathbf{e}}_g'] = \sigma^2 \mathbf{M}_g \quad (44)$$

where

$$\mathbf{M}_g = \mathbf{I}_{n_g} - \mathbf{X}_g \mathbf{X}_g' \leq \mathbf{I}_{n_g} - \mathbf{Z}_g \mathbf{Z}_g'. \quad (45)$$

The inequality in (45) can be shown by the following argument. Let  $\mathbf{R}_\perp$  be  $k \times (k-1)$  such that  $[\mathbf{R}, \mathbf{R}_\perp]$  is

orthonormal. Set  $\mathbf{Z}_{\perp g} = \mathbf{X}_g \mathbf{R}_{\perp}$ . Then

$$\mathbf{X}_g \mathbf{X}'_g = \mathbf{X}_g [\mathbf{R}, \mathbf{R}_{\perp}] \begin{bmatrix} \mathbf{R}' \\ \mathbf{R}'_{\perp} \end{bmatrix} \mathbf{X}'_g = [\mathbf{Z}_g, \mathbf{Z}_{\perp g}] \begin{bmatrix} \mathbf{Z}'_g \\ \mathbf{Z}'_{\perp g} \end{bmatrix} = \mathbf{Z}_g \mathbf{Z}'_g + \mathbf{Z}_{\perp g} \mathbf{Z}'_{\perp g} \geq \mathbf{Z}_g \mathbf{Z}'_g,$$

which implies (45). Taking expectations of (43) using (44), (45),  $\mathbf{Z}'_g \mathbf{Z}_g = c_g$ , and  $\sum_{g=1}^G c_g = 1$ , we find that

$$\begin{aligned} \mathbb{E}[\hat{v}_1^2] &= d \sum_{g=1}^G \mathbf{Z}'_g \mathbb{E}[\hat{\mathbf{e}}_g \hat{\mathbf{e}}'_g] \mathbf{Z}_g \\ &= d \sigma^2 \sum_{g=1}^G \mathbf{Z}'_g \mathbf{M}_g \mathbf{Z}_g \\ &\leq d \sigma^2 \sum_{g=1}^G \mathbf{Z}'_g (\mathbf{I}_{n_g} - \mathbf{Z}_g \mathbf{Z}'_g) \mathbf{Z}_g \\ &= d \sigma^2 \left( 1 - \sum_{g=1}^G c_g^2 \right). \end{aligned}$$

Together with  $v^2 = \sigma^2$ , we find that

$$\frac{\mathbb{E}[\hat{v}_1^2]}{v^2} \leq d \left( 1 - \sum_{g=1}^G c_g^2 \right). \quad (46)$$

Inequality (46) implies that the left side of (17) is weakly smaller than the infimum of the right-hand side of (46) over  $c_g$ . Thus

$$\inf_{(\mathbf{X}, \boldsymbol{\Sigma}) \in \mathcal{F}_0^*} \frac{\mathbb{E}[\hat{v}_1^2]}{v^2} \leq d \inf_{\sum_{g=1}^G c_g = 1} \left( 1 - \sum_{g=1}^G c_g^2 \right) = 0. \quad (47)$$

The rightmost equality in (47) is attained as  $(c_1, \dots, c_G) \rightarrow (1, 0, \dots, 0)$ . Since the left side of (47) is non-negative, the equation hold as an equality. This verifies (17). ■

**Proof of Theorem 2, equation (18):** We use the same normalizations and notation as in the proof of (17), add the assumption that  $\lambda_{\min}(\mathbf{M}_g) \geq \delta > 0$  for all  $g > 1$ , which means that all clusters other than  $g = 1$  are uniformly clusterwise invertible, and assume that  $\boldsymbol{\Sigma}_1 = \mathbf{I}_{n_1}$  and  $\boldsymbol{\Sigma}_g = \mathbf{0}$ , which is extreme heteroskedasticity. Under these conditions, you can calculate that  $\mathbf{V} = \mathbf{X}'_1 \mathbf{X}_1$  and  $v^2 = c_1$ .

The CRVE<sub>2</sub> estimator for  $v^2$  can be written as

$$\hat{v}_2^2 = \sum_{g=1}^G \mathbf{Z}'_g \mathbf{M}_g^{-1/2} \hat{\mathbf{e}}_g \hat{\mathbf{e}}'_g \mathbf{M}_g^{-1/2} \mathbf{Z}_g. \quad (48)$$

We can calculate that for  $g = 1$ ,

$$\mathbb{E}[\hat{\mathbf{e}}_1 \hat{\mathbf{e}}'_1] = \mathbf{I}_{n_1} - 2\mathbf{X}_1 \mathbf{X}'_1 + \mathbf{X}_1 \mathbf{X}'_1 \mathbf{X}_1 \mathbf{X}'_1 = \mathbf{M}_1 \mathbf{M}_1, \quad (49)$$

and for  $g > 1$ ,

$$\mathbb{E}[\hat{\mathbf{e}}_g \hat{\mathbf{e}}'_g] = \mathbf{X}_g \mathbf{X}'_1 \mathbf{X}_1 \mathbf{X}'_g. \quad (50)$$

Taking expectations of (48) and using (49)-(50)

$$\begin{aligned}\mathbb{E}[\widehat{v}_2^2] &= \sum_{g=1}^G \mathbf{Z}'_g \mathbf{M}_g^{-1/2} \mathbb{E}[\widehat{\boldsymbol{\epsilon}}_g \widehat{\boldsymbol{\epsilon}}'_g] \mathbf{M}_g^{-1/2} \mathbf{Z}_g \\ &= \mathbf{Z}'_1 \mathbf{M}_1^{-1/2} \mathbf{M}_1 \mathbf{M}_1 \mathbf{M}_1^{-1/2} \mathbf{Z}_1 + \sum_{g=2}^G \mathbf{Z}'_g \mathbf{M}_g^{-1/2} \mathbf{X}_g \mathbf{X}'_1 \mathbf{X}_1 \mathbf{X}'_g \mathbf{M}_g^{-1/2} \mathbf{Z}_g.\end{aligned}\quad (51)$$

The first term on the right side of (51) equals

$$\mathbf{Z}'_1 \mathbf{M}_1 \mathbf{Z}_1 \leq \mathbf{Z}'_1 (\mathbf{I}_{n_1} - \mathbf{Z}_1 \mathbf{Z}'_1) \mathbf{Z}_1 = c_1 - c_1^2,$$

where the inequality is (45). Since

$$\lambda_{\max}(\mathbf{M}_g^{-1/2} \mathbf{X}_g \mathbf{X}'_1 \mathbf{X}_1 \mathbf{X}'_g \mathbf{M}_g^{-1/2}) \leq \lambda_{\max}(\mathbf{M}_g^{-1}) \lambda_{\max}(\mathbf{X}'_1 \mathbf{X}_1) \lambda_{\max}(\mathbf{X}'_g \mathbf{X}_g) \leq \frac{1}{\delta}$$

(using the assumption  $\lambda_{\min}(\mathbf{M}_g) \geq \delta$  for any  $g \geq 2$  and the fact  $\mathbf{X}'_g \mathbf{X}_g \leq \mathbf{X}' \mathbf{X} = \mathbf{I}_k$  for any  $g \geq 1$ ), the second term on the right side of (51) satisfies

$$\sum_{g=2}^G \mathbf{Z}'_g \mathbf{M}_g^{-1/2} \mathbf{X}_g \mathbf{X}'_1 \mathbf{X}_1 \mathbf{X}'_g \mathbf{M}_g^{-1/2} \mathbf{Z}_g \leq \sum_{g=2}^G \frac{\mathbf{Z}'_g \mathbf{Z}_g}{\delta} = \frac{1 - c_1}{\delta}.\quad (52)$$

Together with  $v^2 = 1$ , we find that

$$\frac{\mathbb{E}[\widehat{v}_2^2]}{v^2} \leq 1 - c_1 + \frac{1 - c_1}{\delta c_1}.\quad (53)$$

The assumptions we have made are a special case of the model class  $\mathcal{F}^*$ . Therefore, the left side of (18) is weakly smaller than the infimum of (53) over  $c_1$ . Hence

$$\inf_{(\mathbf{X}, \boldsymbol{\Sigma}) \in \mathcal{F}^*} \frac{\mathbb{E}[\widehat{v}_2^2]}{v^2} \leq \inf_{0 < c_1 < 1} \left(1 - c_1 + \frac{1 - c_1}{\delta c_1}\right) = 0.\quad (54)$$

The rightmost equality in (54) is attained as  $c_1 \rightarrow 1$ . Equation (54) implies (18), as claimed.  $\blacksquare$

**Proof of Theorem 2, equation (19):** We first show that when  $\mathbf{X}$  is clusterwise invertible,

$$\mathbb{E}[\widehat{\mathbf{V}}_3] \geq \left(\frac{G-1}{G}\right)^2 \mathbf{V}.\quad (55)$$

The proof of (55) is analogous to that of Theorem 1. Using (9) under clusterwise invertibility, collecting

terms, (14), and then (39), we calculate that

$$\begin{aligned}
\tilde{\beta}_{-g} - \bar{\beta} &= (\tilde{\beta}_{-g} - \hat{\beta}) - \frac{1}{G} \sum_{h=1}^G (\tilde{\beta}_{-h} - \hat{\beta}) \\
&= \left( \frac{G-1}{G} \right) (\tilde{\beta}_{-g} - \hat{\beta}) - \frac{1}{G} \sum_{h \neq g} (\tilde{\beta}_{-h} - \hat{\beta}) \\
&= - \left( \frac{G-1}{G} \right) (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'_g \tilde{\mathbf{e}}_g + \frac{1}{G} (\mathbf{X}'\mathbf{X})^{-1} \sum_{h \neq g} \mathbf{X}'_h \tilde{\mathbf{e}}_h \\
&= - \left( \frac{G-1}{G} \right) (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'_g \mathbf{e}_g - \mathbf{S}_g
\end{aligned} \tag{56}$$

where

$$\mathbf{S}_g = - \left( \frac{G-1}{G} \right) (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'_g \mathbf{X}_g (\tilde{\beta}_{-g} - \beta) - \frac{1}{G} (\mathbf{X}'\mathbf{X})^{-1} \sum_{h \neq g} \mathbf{X}'_h \mathbf{e}_h + \frac{1}{G} (\mathbf{X}'\mathbf{X})^{-1} \sum_{h \neq g} \mathbf{X}'_h \mathbf{X}_h (\tilde{\beta}_{-h} - \beta).$$

Notice that the first two components of  $\mathbf{S}_g$  are uncorrelated with  $\mathbf{e}_g$ . Thus

$$\begin{aligned}
\mathbb{E} \left[ \left( \frac{G-1}{G} \right) (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'_g \mathbf{e}_g \mathbf{S}'_g \right] &= \frac{G-1}{G^2} (\mathbf{X}'\mathbf{X})^{-1} \sum_{h \neq g} \mathbf{X}'_g \mathbb{E} \left[ \mathbf{e}_g (\tilde{\beta}_{-h} - \beta)' \right] \mathbf{X}'_h \mathbf{X}_h (\mathbf{X}'\mathbf{X})^{-1} \\
&= \frac{G-1}{G^2} (\mathbf{X}'\mathbf{X})^{-1} \sum_{h \neq g} \mathbf{X}'_g \Sigma_g \mathbf{X}_g (\mathbf{X}'\mathbf{X} - \mathbf{X}'_h \mathbf{X}_h)^{-1} \mathbf{X}'_h \mathbf{X}_h (\mathbf{X}'\mathbf{X})^{-1} \\
&= \mathbf{A}_g,
\end{aligned}$$

say, where the second equality uses the relationship (under clusterwise invertibility)

$$\tilde{\beta}_{-h} - \beta = (\mathbf{X}'\mathbf{X} - \mathbf{X}'_h \mathbf{X}_h)^{-1} \left( \sum_{\ell \neq h} \mathbf{X}'_\ell \mathbf{e}_\ell \right).$$

Using (56), it follows that

$$\begin{aligned}
\mathbb{E} [\widehat{\mathbf{V}}_3] &= \left( \frac{G-1}{G} \right) \sum_{g=1}^G \mathbb{E} \left[ (\tilde{\beta}_{-g} - \bar{\beta}) (\tilde{\beta}_{-g} - \bar{\beta})' \right] \\
&= \left( \frac{G-1}{G} \right)^2 (\mathbf{X}'\mathbf{X})^{-1} \sum_{g=1}^G \mathbf{X}'_g \mathbb{E} \left[ \mathbf{e}_g \mathbf{e}'_g \right] \mathbf{X}_g (\mathbf{X}'\mathbf{X})^{-1} + \sum_{g=1}^G \left( \mathbb{E} \left[ \mathbf{S}_g \mathbf{S}'_g \right] + \mathbf{A}_g + \mathbf{A}'_g \right) \\
&\geq \left( \frac{G-1}{G} \right)^2 \mathbf{V}.
\end{aligned} \tag{57}$$

The positive semi-definite (PSD) inequality (57) holds because  $\mathbb{E} \left[ \mathbf{S}_g \mathbf{S}'_g \right]$  and  $\mathbf{A}_g + \mathbf{A}'_g$  are PSD. The sum  $\mathbf{A}_g + \mathbf{A}'_g$  is PSD because  $\mathbf{A}_g$  is the matrix product of PSD matrices, and thus has non-negative eigenvalues (see Zhang and Zhang (2006, Corollary 11)). This establishes (55).

Equation (55) implies

$$\inf_{(\mathbf{X}, \boldsymbol{\Sigma}) \in \mathcal{F}^*} \frac{\mathbb{E}[\hat{\mathbf{v}}_3^2]}{\mathbf{v}^2} = \inf_{(\mathbf{X}, \boldsymbol{\Sigma}) \in \mathcal{F}^*} \frac{R' \mathbb{E}[\hat{\mathbf{V}}_3] R}{R' \mathbf{V} R} \geq \left( \frac{G-1}{G} \right)^2. \quad (58)$$

To show that (58) holds as an equality we calculate an upper bound for the left side of (58) in the context of the example from the proof of (18). We adopt the assumptions made therein. The model class  $\mathcal{F}^*$  imposes clusterwise invertibility. As shown by equation (19) of MacKinnon, Nielsen, and Webb (2023b), clusterwise invertibility and  $\mathbf{X}' \mathbf{X} = \mathbf{I}_k$  implies

$$R'(\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}}_{-g}) = R' \mathbf{X}'_g \mathbf{M}_g^{-1} \hat{\mathbf{e}}_g = \mathbf{Z}'_g \mathbf{M}_g^{-1} \hat{\mathbf{e}}_g.$$

Consequently,

$$\begin{aligned} \hat{\mathbf{v}}_3^2 &= \left( \frac{G-1}{G} \right) \sum_{g=1}^G R' (\tilde{\boldsymbol{\beta}}_{-g} - \bar{\boldsymbol{\beta}}) (\tilde{\boldsymbol{\beta}}_{-g} - \bar{\boldsymbol{\beta}})' R \\ &= \left( \frac{G-1}{G} \right) \sum_{g=1}^G R' (\tilde{\boldsymbol{\beta}}_{-g} - \hat{\boldsymbol{\beta}}) (\tilde{\boldsymbol{\beta}}_{-g} - \hat{\boldsymbol{\beta}})' R - (G-1) R' (\hat{\boldsymbol{\beta}} - \bar{\boldsymbol{\beta}}) (\hat{\boldsymbol{\beta}} - \bar{\boldsymbol{\beta}})' R \\ &= \left( \frac{G-1}{G} \right) \sum_{g=1}^G \mathbf{Z}'_g \mathbf{M}_g^{-1} \hat{\mathbf{e}}_g \hat{\mathbf{e}}_g' \mathbf{M}_g^{-1} \mathbf{Z}'_g - \left( \frac{G-1}{G^2} \right) \left( \sum_{g=1}^G \mathbf{Z}'_g \mathbf{M}_g^{-1} \hat{\mathbf{e}}_g \right) \left( \sum_{g=1}^G \hat{\mathbf{e}}_g' \mathbf{M}_g^{-1} \mathbf{Z}'_g \right) \\ &= \left( \frac{G-1}{G} \right)^2 \mathbf{Z}'_1 \mathbf{M}_1^{-1} \hat{\mathbf{e}}_1 \hat{\mathbf{e}}_1' \mathbf{M}_1^{-1} \mathbf{Z}_1 \end{aligned} \quad (59)$$

$$+ \left( \frac{G-1}{G} \right)^2 \sum_{g=2}^G \mathbf{Z}'_g \mathbf{M}_g^{-1} \hat{\mathbf{e}}_g \hat{\mathbf{e}}_g' \mathbf{M}_g^{-1} \mathbf{Z}_g \quad (60)$$

$$- 2 \left( \frac{G-1}{G^2} \right) \sum_{g=2}^G \mathbf{Z}'_g \mathbf{M}_g^{-1} \hat{\mathbf{e}}_g \hat{\mathbf{e}}_1' \mathbf{M}_1^{-1} \mathbf{Z}_1 \quad (61)$$

$$- \left( \frac{G-1}{G^2} \right) \sum_{h \neq g \neq 1} \mathbf{Z}'_g \mathbf{M}_g^{-1} \hat{\mathbf{e}}_g \hat{\mathbf{e}}_h' \mathbf{M}_h^{-1} \mathbf{Z}_h. \quad (62)$$

Using (49), the expectation of (59) equals  $((G-1)/G)^2$  times

$$\mathbf{Z}'_1 \mathbf{M}_1^{-1} \mathbb{E}[\hat{\mathbf{e}}_1 \hat{\mathbf{e}}_1'] \mathbf{M}_1^{-1} \mathbf{Z}_1 = \mathbf{Z}'_1 \mathbf{M}_1^{-1} \mathbf{M}_1 \mathbf{M}_1 \mathbf{M}_1^{-1} \mathbf{Z}_1 = \mathbf{Z}'_1 \mathbf{Z}_1 = c_1.$$

Using (50) the expectation of (60) equals  $((G-1)/G)^2$  times

$$\sum_{g=2}^G \mathbf{Z}'_g \mathbf{M}_g^{-1} \mathbb{E}[\hat{\mathbf{e}}_g \hat{\mathbf{e}}_g'] \mathbf{M}_g^{-1} \mathbf{Z}_g = \sum_{g=2}^G \mathbf{Z}'_g \mathbf{M}_g^{-1} \mathbf{X}_g \mathbf{X}'_1 \mathbf{X}_1 \mathbf{X}'_g \mathbf{M}_g^{-1} \mathbf{Z}_g \leq \frac{1-c_1}{\delta^2},$$

where the inequality follows by similar same steps as for (52).

We calculate that for  $g \neq 1$

$$\mathbb{E}[\hat{\mathbf{e}}_g \hat{\mathbf{e}}_1'] = -\mathbf{X}_g (\mathbf{I}_k - \mathbf{X}'_1 \mathbf{X}_1) \mathbf{X}'_1. \quad (63)$$



The assumption  $\lambda_{\min}(\mathbf{M}_g) \geq \delta$  implies that

$$\sum_{g=2}^G \mathbf{X}'_g \mathbf{M}_g^{-1} \mathbf{X}_g \geq \frac{1}{\delta} \sum_{g=2}^G \mathbf{X}'_g \mathbf{X}_g = \frac{1}{\delta} (\mathbf{I}_k - \mathbf{X}'_1 \mathbf{X}_1). \quad (64)$$

Using the Woodbury identity and the fact  $\mathbf{X}'_1 \mathbf{X}_1 \leq \mathbf{X}' \mathbf{X} = \mathbf{I}_k$ ,

$$\begin{aligned} \mathbf{X}'_1 \mathbf{M}_1^{-1} \mathbf{X}_1 &= \mathbf{X}'_1 \mathbf{X}_1 + \mathbf{X}'_1 \mathbf{X}_1 (\mathbf{I}_k - \mathbf{X}'_1 \mathbf{X}_1)^{-1} \mathbf{X}'_1 \mathbf{X}_1 \\ &\leq \mathbf{I}_k + (\mathbf{I}_k - \mathbf{X}'_1 \mathbf{X}_1)^{-1} \\ &\leq 2(\mathbf{I}_k - \mathbf{X}'_1 \mathbf{X}_1)^{-1}. \end{aligned} \quad (65)$$

Combining (63), (64), and (65), the expectation of (61) equals  $2(G-1)/G^2$  times

$$\begin{aligned} R' \sum_{g=2}^G \mathbf{X}'_g \mathbf{M}_g^{-1} \mathbf{X}_g (\mathbf{I}_k - \mathbf{X}'_1 \mathbf{X}_1) \mathbf{X}'_1 \mathbf{M}_1^{-1} \mathbf{X}_1 R &\leq \frac{2}{\delta} R' (\mathbf{I}_k - \mathbf{X}'_1 \mathbf{X}_1) (\mathbf{I}_k - \mathbf{X}'_1 \mathbf{X}_1) (\mathbf{I}_k - \mathbf{X}'_1 \mathbf{X}_1)^{-1} R \\ &= \frac{2}{\delta} R' (\mathbf{I}_k - \mathbf{X}'_1 \mathbf{X}_1) R \\ &= \frac{2(1-c_1)}{\delta}. \end{aligned}$$

We calculate that for  $h \neq g \neq 1$

$$\mathbb{E}[\widehat{\mathbf{e}}_g \widehat{\mathbf{e}}'_h] = \mathbf{X}_g \mathbf{X}'_1 \mathbf{X}_1 \mathbf{X}'_h,$$

so the expectation of (62) equals  $(G-1)/G^2$  times

$$-R' \sum_{h \neq g \neq 1} \mathbf{X}'_g \mathbf{M}_g^{-1} \mathbf{X}_g \mathbf{X}'_1 \mathbf{X}_1 \mathbf{X}'_h \mathbf{M}_h^{-1} \mathbf{X}_h R \leq 0,$$

where the inequality holds since it is a quadratic form of a sum of positive semi-definite matrices.

Together, with  $v^2 = c_1$ , we have

$$\frac{\mathbb{E}[\widehat{v}_3^2]}{v^2} \leq \left(\frac{G-1}{G}\right)^2 + \left(\frac{G-1}{G}\right)^2 \frac{1-c_1}{\delta^2 c_1} + 4 \left(\frac{G-1}{G^2}\right) \frac{1-c_1}{\delta}. \quad (66)$$

As in the proof of (18), the left side of (19) is weakly smaller than the infimum of (66) over  $c_1$ . Hence

$$\inf_{(\mathbf{X}, \boldsymbol{\Sigma}) \in \mathcal{F}^*} \frac{\mathbb{E}[\widehat{v}_3^2]}{v^2} \leq \inf_{0 < c_1 < 1} \left( \left(\frac{G-1}{G}\right)^2 + \left(\frac{G-1}{G}\right)^2 \frac{1-c_1}{\delta^2 c_1} + 4 \left(\frac{G-1}{G^2}\right) \frac{1-c_1}{\delta} \right) = \left(\frac{G-1}{G}\right)^2.$$

The rightmost equality in (54) is attained as  $c_1 \rightarrow 1$ . Combined with (58) this establishes that (58) holds as an equality, which is (19), as claimed. ■

**Proof of Theorem 2, equation (20):** Since the model is clusterwise invertible,

$$\widehat{v}_4^2 = \left(\frac{G-1}{G}\right) \widehat{v}_5^2.$$

The result follows from (23), which we establish below. ■

**Proof of Theorem 2, equations (21) and (22):** We calculate an upper bound for the left side of (21)-(22) assuming one noninvertible cluster. Without loss of generality, assume  $\mathbf{X}'\mathbf{X} = \mathbf{I}_k$  and assume the noninvertible cluster is  $g = 1$ . This means that there is a linear combination  $\mathbf{X}_g\gamma$  of the regressors which is identically zero for  $g \neq 1$ . Without loss of generality, assume  $\gamma = R$ . Since the regressor matrix is singular when cluster 1 is omitted, this cluster is discarded from the calculation of  $\hat{v}_3^2$  and  $\hat{v}_4^2$ . Therefore the latter equals

$$\begin{aligned}\hat{v}_4^2 &= \left(\frac{G-2}{G-1}\right) \sum_{g=2}^G R' (\tilde{\beta}_{-g} - \hat{\beta}) (\tilde{\beta}_{-g} - \hat{\beta})' R \\ &= \left(\frac{G-2}{G-1}\right) \sum_{g=2}^G R' \mathbf{X}'_g \mathbf{M}_g^{-1} \hat{\mathbf{e}}_g \hat{\mathbf{e}}'_g \mathbf{M}_g^{-1} \mathbf{X}_g R \\ &= 0\end{aligned}$$

since  $\mathbf{X}_g R = 0$  for  $g \neq 1$ . Thus  $\hat{v}_4^2$  is identically zero. Since  $\hat{v}_3^2 \leq \hat{v}_4^2$  we find that  $\hat{v}_3^2 = 0$  as well. This implies (21)-(22) as stated. ■

**Proof of Theorem 2, equation (23):** Equation (16) implies

$$\inf_{(\mathbf{X}, \boldsymbol{\Sigma}) \in \mathcal{F}} \frac{\mathbb{E}[\hat{v}_5^2]}{v^2} = \inf_{(\mathbf{X}, \boldsymbol{\Sigma}) \in \mathcal{F}} \frac{R' \mathbb{E}[\hat{\mathbf{V}}_5] R}{R' \mathbf{V} R} \geq \inf_{(\mathbf{X}, \boldsymbol{\Sigma}) \in \mathcal{F}} \frac{R' \mathbf{V} R}{R' \mathbf{V} R} = 1. \quad (67)$$

To show that this is a strict equality we calculate an upper bound for the left side of (67) in the context of the example from the proof of (19). By the calculations from that proof,

$$\hat{v}_5^2 = \sum_{g=1}^G \mathbf{Z}'_g \mathbf{M}_g^{-1} \hat{\mathbf{e}}_g \hat{\mathbf{e}}'_g \mathbf{M}_g^{-1} \mathbf{Z}_g$$

which satisfies

$$\frac{\mathbb{E}[\hat{v}_5^2]}{v^2} \leq 1 + \frac{1 - c_1}{\delta^2 c_1}. \quad (68)$$

This implies

$$\inf_{(\mathbf{X}, \boldsymbol{\Sigma}) \in \mathcal{F}} \frac{\mathbb{E}[\hat{v}_5^2]}{v^2} \leq \inf_{0 < c_1 < 1} \left[ 1 + \frac{1 - c_1}{\delta^2 c_1} \right] = 1. \quad (69)$$

Combined with (67) this yields (23) as stated. ■

For the proof of Theorem 3 we use the following distributional lower bound.

**Theorem 7** For any  $n \geq 1$  let  $\boldsymbol{\xi} = (\xi_0, \xi_1, \dots, \xi_n) \sim N(\mathbf{0}, \boldsymbol{\Omega})$  where

$$\boldsymbol{\Omega} = \begin{bmatrix} 1 & \rho_1 & \cdots & \rho_n \\ \rho_1 & 1 & \cdots & 0 \\ \vdots & & \ddots & \vdots \\ \rho_n & 0 & \cdots & 1 \end{bmatrix}.$$

Then for any  $w_j \geq 0$  satisfying  $\sum_{j=1}^n w_j \geq 1$ , any  $\delta_j$ , and any  $1 \leq x < \infty$ ,

$$\mathbb{P} \left[ \frac{\xi_0^2}{\sum_{j=1}^n w_j (\xi_j + \delta_j)^2} \leq x \right] \geq F(x; 1, 1), \quad (70)$$

where  $F(x; k_1, k_2)$  denotes the  $F$  distribution function with degrees of freedom  $k_1$  and  $k_2$ .

**Proof of Theorem 7:** Our proof examines the inequality in four cases: (a)  $\delta_j = 0$  and  $\rho_j = 0$  for all  $j$ ; (b)  $\rho_j = 0$  for all  $j$  but  $\delta_j \neq 0$ ; (c)  $\delta_j = 0$  for all  $j$  but  $\rho_j \neq 0$ ; (d) general  $\rho_j$  and  $\delta_j$ .

Case (a):  $\delta_j = 0$  and  $\rho_j = 0$  for all  $j$ . The assumption that  $\rho_j = 0$  implies that the  $\xi_j$  are mutually independent. Set  $\bar{w} = \sum_{j=1}^n w_j \geq 1$ ,  $\alpha_j = w_j / \bar{w}$  for  $j \leq n-1$ ,  $\alpha_n = 1 - \sum_{j=1}^{n-1} \alpha_j$ ,  $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_{n-1})'$ ,  $\mathbf{U} = (\xi_1^2 - \xi_n^2, \dots, \xi_{n-1}^2 - \xi_n^2)$ , and let  $H(x) = \mathbb{P}[\chi_1^2 \leq x]$  denote the  $\chi_1^2$  cumulative distribution function. Using the law of iterated expectations, the monotonicity of  $H(x)$ , and  $\bar{w} \geq 1$ , the left side of (70) equals

$$\begin{aligned} \mathbb{E} \left[ \mathbb{P} \left[ \xi_1^2 \leq x \sum_{j=1}^n w_j \xi_j^2 \mid \xi_1^2, \dots, \xi_n^2 \right] \right] &= \mathbb{E} \left[ H \left( \bar{w} x \sum_{j=1}^n \alpha_j \xi_j^2 \right) \right] \\ &\geq \mathbb{E} \left[ H \left( x \sum_{j=1}^n \alpha_j \xi_j^2 \right) \right] \\ &= \mathbb{E} [H(x \boldsymbol{\alpha}' \mathbf{U} + x \xi_n^2)] \\ &\equiv f(\boldsymbol{\alpha}). \end{aligned} \quad (71)$$

We calculate that

$$\frac{\partial^2}{\partial \boldsymbol{\alpha} \partial \boldsymbol{\alpha}'} f(\boldsymbol{\alpha}) = x^2 \mathbb{E} [\mathbf{U} \mathbf{U}' H''(x \boldsymbol{\alpha}' \mathbf{U} + x \xi_n^2)].$$

This matrix is negative semi-definite since the  $\chi_1^2$  distribution function  $H(t)$  is globally concave. Hence, the function  $f(\boldsymbol{\alpha})$  is concave in  $\boldsymbol{\alpha}$ . It follows that the minimum of  $f(\boldsymbol{\alpha})$  over  $\boldsymbol{\alpha}$  is obtained at a corner. By symmetry we can take any corner, e.g.  $\boldsymbol{\alpha}_{\min} = (1, 0, \dots, 0)$ . Hence, the minimum of (71) over  $\boldsymbol{\alpha}$  equals

$$\mathbb{E} [H(x \boldsymbol{\alpha}'_{\min} \mathbf{U} + x \xi_n^2)] = \mathbb{E} [H(x \xi_1^2)] = \mathbb{P} \left[ \frac{\xi_0^2}{\xi_1^2} \leq x \right] = F(x; 1, 1),$$

since  $\xi_0^2 / \xi_1^2$  is distributed  $F$  with degrees of freedom (1, 1). This establishes (70).

Case (b):  $\rho_j = 0$  for all  $j$  but  $\delta_j \neq 0$ . In this case,  $\sum_{j=1}^n w_j (\xi_j + \delta_j)^2$  is independent of  $\xi_0^2$ . By Theorem 3 of Mathew and Nordström (1997),  $\sum_{j=1}^n w_j \xi_j^2$  is stochastically dominated by  $\sum_{j=1}^n w_j (\xi_j + \delta_j)^2$ . This

implies

$$\mathbb{P} \left[ \frac{\xi_0^2}{\sum_{j=1}^n w_j (\xi_j + \delta_j)^2} \leq x \right] \geq \mathbb{P} \left[ \frac{\xi_0^2}{\sum_{j=1}^n w_j \xi_j^2} \leq x \right]. \quad (72)$$

The right-side corresponds to Case (a), which we have already shown is bounded below by  $F(x; 1, 1)$ . This establishes (70).

Case (c):  $\delta_j = 0$  for all  $j$  but  $\rho_j \neq 0$ . Define the  $(n+1) \times (n+1)$  positive semi-definite matrices  $\mathbf{A} = \text{diag}\{1, 0, \dots, 0\}$ ,  $\mathbf{B} = \text{diag}\{0, w_1, \dots, w_n\}$ , and  $\mathbf{C} = \mathbf{\Omega}^{1/2} (\mathbf{A} - x\mathbf{B}) \mathbf{\Omega}^{1/2}$ . By the spectral decomposition,  $\mathbf{C} = \mathbf{H}\mathbf{\Lambda}\mathbf{H}'$ , where  $\mathbf{H}'\mathbf{H} = \mathbf{I}_{n+1}$ ,  $\mathbf{\Lambda} = \text{diag}\{\lambda_0, \dots, \lambda_n\}$ , and  $\lambda_0 \geq \lambda_1 \geq \dots \geq \lambda_n$  are the ordered eigenvalues of  $\mathbf{C}$ . Set  $\mathbf{u} = \mathbf{\Omega}^{-1/2} \boldsymbol{\xi} \sim \mathbf{N}(0, \mathbf{I}_{n+1})$  and  $\boldsymbol{\zeta} = \mathbf{H}'\mathbf{u} \sim \mathbf{N}(0, \mathbf{I}_{n+1})$ . Partition  $\boldsymbol{\zeta} = (\zeta_0, \dots, \zeta_n)'$  in conformity with  $\mathbf{\Lambda}$ .

The left side of (70) equals

$$\begin{aligned} \mathbb{P} \left[ \frac{\xi_0^2}{\sum_{j=1}^n w_j \xi_j^2} \leq x \right] &= \mathbb{P} \left[ \xi_0^2 \leq x \sum_{j=1}^n w_j \xi_j^2 \right] \\ &= \mathbb{P} [\boldsymbol{\xi}' (\mathbf{A} - x\mathbf{B}) \boldsymbol{\xi} \leq 0] \\ &= \mathbb{P} [\mathbf{u}' \mathbf{C} \mathbf{u} \leq 0] \\ &= \mathbb{P} [\boldsymbol{\zeta}' \mathbf{\Lambda} \boldsymbol{\zeta} \leq 0] \\ &= \mathbb{P} \left[ \sum_{j=0}^n \lambda_j \zeta_j^2 \leq 0 \right]. \end{aligned} \quad (73)$$

We next establish two bounds on the eigenvalues  $\lambda_j$ . First,

$$\begin{aligned} \lambda_0 &= \lambda_{\max}(\mathbf{C}) \\ &\leq \lambda_{\max}(\mathbf{\Omega}^{1/2} \mathbf{A} \mathbf{\Omega}^{1/2}) + \lambda_{\max}(-x\mathbf{\Omega}^{1/2} \mathbf{B} \mathbf{\Omega}^{1/2}) \\ &= \lambda_{\max}(\mathbf{A}\mathbf{\Omega}) - x\lambda_{\min}(\mathbf{B}\mathbf{\Omega}) \\ &= 1, \end{aligned} \quad (74)$$

since  $\lambda_{\max}(\mathbf{A}\mathbf{\Omega}) = 1$  by direct calculation, and  $\lambda_{\min}(\mathbf{B}\mathbf{\Omega}) = 0$  since  $\mathbf{B}$  has deficient rank.

Second, by a corollary of the Weyl eigenvalue inequality for Hermitian matrices (Corollary 4.3.15 of Horn and Johnson (2013)), for each  $j \geq 1$ ,

$$\begin{aligned} \lambda_j &= \lambda_{j+1}(\mathbf{C}) \\ &\leq \lambda_{j+1}(\mathbf{\Omega}^{1/2} \mathbf{A} \mathbf{\Omega}^{1/2}) + \lambda_{\max}(-x\mathbf{\Omega}^{1/2} \mathbf{B} \mathbf{\Omega}^{1/2}) \\ &= \lambda_{j+1}(\mathbf{A}\mathbf{\Omega}) - x\lambda_{\min}(\mathbf{B}\mathbf{\Omega}) \\ &= 0, \end{aligned} \quad (75)$$

the final equality since both  $\mathbf{A}$  and  $\mathbf{B}$  have deficient rank. Together, (74) and (75) show that the largest eigenvalue  $\lambda_0$  of  $\mathbf{C}$  is less than one, and the remaining are non-positive.

If  $\lambda_0 \leq 0$  then  $\lambda_j \leq 0$  for all  $j$  by (75), so (73) equals 1. This implies (70). We thus assume  $\lambda_0 > 0$  for the

remainder of the proof. For  $j = 1, \dots, n$  set  $\alpha_j = -\lambda_j / (\lambda_0 x) \geq 0$ . We see that (73) equals

$$\mathbb{P} \left[ \zeta_0^2 \leq x \sum_{j=1}^n \alpha_j \zeta_j^2 \right] = \mathbb{P} \left[ \frac{\zeta_0^2}{\sum_{j=1}^n \alpha_j \zeta_j^2} \leq x \right]. \quad (76)$$

We next establish

$$\sum_{j=1}^n \alpha_j \geq 1. \quad (77)$$

Using the property  $\text{tr}(\mathbf{C}) = \sum_{j=0}^n \lambda_j$  and the direct calculations  $\text{tr}(\mathbf{A}\mathbf{\Omega}) = 1$  and  $\text{tr}(\mathbf{B}\mathbf{\Omega}) = \sum_{j=1}^n w_j$ ,

$$\begin{aligned} \sum_{j=1}^n \alpha_j &= \frac{-\sum_{j=1}^n \lambda_j}{\lambda_0 x} \\ &= \frac{\lambda_0 - \text{tr}(\mathbf{C})}{\lambda_0 x} \\ &= \frac{\lambda_0 - \text{tr}(\mathbf{A}\mathbf{\Omega}) + x \text{tr}(\mathbf{B}\mathbf{\Omega})}{\lambda_0 x} \\ &= \frac{\lambda_0 - 1 + x \sum_{j=1}^n w_j}{\lambda_0 x} \\ &\geq \frac{\lambda_0 - 1 + x}{\lambda_0 x} \\ &\geq \frac{\lambda_0 + (x-1)\lambda_0}{\lambda_0 x} \\ &= 1 \end{aligned}$$

the first inequality using  $\sum_{j=1}^n w_j \geq 1$  and the second using  $x \geq 1$  and (74). This establishes (77). We have shown that the left side of (70) equals (76) where  $\zeta_j \sim \mathcal{N}(0, 1)$  are mutually independent and  $\sum_{j=1}^n \alpha_j \geq 1$ . This corresponds to Case (a), which we have already shown is bounded below by  $F(x; 1, 1)$ . This establishes (70).

Case (d): general  $\rho_j$  and  $\delta_j$ . Use the same notation as in the proof of Case (c), plus  $\boldsymbol{\delta} = (0, \delta_1, \dots, \delta_n)'$  and  $\boldsymbol{\mu} = \mathbf{H}'\mathbf{\Omega}^{-1/2}\boldsymbol{\delta}$ . By similar manipulations, the left side of (70) equals

$$\mathbb{P} \left[ \sum_{j=0}^n \lambda_j (\zeta_j + \mu_j)^2 \leq 0 \right] = \mathbb{P} \left[ \sum_{j=0}^n \lambda_j \chi_1^2(\mu_j^2) \leq 0 \right] \quad (78)$$

where  $\chi_1^2(\mu_j^2)$  are mutually independent non-central chi-square random variables with one degree of freedom and non-centrality parameter  $\mu_j^2$ , and  $\lambda_j$  are the ordered eigenvalues of  $\mathbf{C}$ . For these calcula-

tions we use the following explicit formulae. Standard manipulations reveal that

$$\mathbf{C} = \begin{bmatrix} \sigma^2 & \sigma\rho_1 & \cdots & \sigma\rho_n \\ \sigma\rho_1 & \rho_1^2 - xw_1 & & \rho_n\rho_1 \\ \vdots & & \ddots & \vdots \\ \sigma\rho_n & \rho_n\rho_1 & \cdots & \rho_n^2 - xw_2 \end{bmatrix}, \quad \boldsymbol{\mu} = \mathbf{H}' \begin{bmatrix} -\sum_{j=1}^n \rho_j \delta_j / \sigma \\ \delta_1 \\ \vdots \\ \delta_n \end{bmatrix}$$

where  $\sigma^2 = 1 - \sum_{j=1}^n \rho_j^2$ . Given the inputs  $(x, \rho_1, \dots, \rho_n, \delta_1, \dots, \delta_n, w_1, \dots, w_n)$ , the matrix  $\mathbf{C}$  can be numerically calculated, from  $\mathbf{C}$  the eigenvalues  $\lambda_j$  and eigenvector matrix  $\mathbf{H}$  can be calculated, and from these the non-centrality coefficients  $\mu_j$ . Given the inputs  $(\lambda_1, \dots, \lambda_n, \mu_1, \dots, \mu_n)$  the probability (78) can be calculated using the method of Imhof (1961). We use R package `CompQuadForm` of Duchesne and Lafaye de Micheaux (2010). The results can be compared to verify that (78) exceeds  $F(x; 1, 1)$ .

The most important cases to investigate are  $n = 1$  and  $n = 2$ , so for these our calculations were most thorough. For  $n = 1$  we calculated (78) for each  $\rho_1 \in (-1, 1)$  and  $\mu_1 \in [-3, 3]$  on a double grid with increments of 0.01, and for each  $x \geq 1$  corresponding to increments of  $F(x; 1, 1)$  of 0.01. Uniformly, the inequality (70) held.

For  $n = 2$  we generated 100,000,000 draws of parameters, calculated (78) for all  $x$  as for the case  $n = 1$ , and compared the result with  $F(x; 1, 1)$ . Uniformly across the 100,000,000 parameter draws and all values of  $x$ , the inequality (70) held. The parameters were drawn as follows. The correlations  $(\rho_1, \rho_2)$  were drawn from a uniform distribution on the unit disk. The coefficients  $(\delta_1, \delta_2)$  were drawn from the  $U[-2, 2]^2$  distribution. The weights  $(w_1, w_2)$  were drawn from a flat (uniform) Dirichlet distribution.

For the cases  $3 \leq n \leq 10$  our calculations were similar to the  $n = 2$  case, but using 10,000,000 draws of parameters for each  $n$ . The correlations  $(\rho_1, \dots, \rho_n)$  were drawn from a uniform distribution on the unit disk. The coefficients  $(\delta_1, \dots, \delta_n)$  were drawn from the  $U[-2, 2]^n$  distribution. The weights  $(w_1, \dots, w_n)$  were drawn from a flat (uniform) Dirichlet distribution. For each parameter draw and all values of  $x$ , the inequality (70) held.

We conclude that the inequality (70) holds as stated. ■

**Proof of Theorem 3:** Take any  $(\mathbf{X}, \boldsymbol{\Sigma}) \in \mathcal{F}$  and  $1 \leq c < \infty$ . Define

$$\xi_0 = \frac{\hat{\theta} - \theta}{v} \sim N(0, 1). \quad (79)$$

Define the delete-one-cluster operator  $\widetilde{\mathbf{M}}$ , which is the  $n \times n$  matrix with  $gj$ th block

$$\widetilde{\mathbf{M}}_{gj} = \begin{cases} \mathbf{I}_{n_g} & g = j \\ -\mathbf{X}_g (\mathbf{X}'\mathbf{X} - \mathbf{X}'_g \mathbf{X}_g)^+ \mathbf{X}'_j & g \neq j. \end{cases} \quad (80)$$

This operator has the algebraic property that it creates delete-one-cluster prediction errors,  $\widetilde{\mathbf{M}}\mathbf{Y} = \tilde{\mathbf{e}}$ , and is the jackknife analog of the least squares annihilation matrix. It has the property that  $\widetilde{\mathbf{M}}\mathbf{X} = \mathbf{0}$  under

clusterwise invertibility, but not under clusterwise noninvertibility. Define

$$\mathbf{Z} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}R, \quad (81)$$

partition by cluster as  $\mathbf{Z} = (\mathbf{Z}'_1, \dots, \mathbf{Z}'_G)'$ , and define

$$\tilde{\mathbf{Z}} = \text{diag}\{\mathbf{Z}_1, \dots, \mathbf{Z}_G\}. \quad (82)$$

Using (15), these definitions, and the equations  $\tilde{\mathbf{e}} = \tilde{\mathbf{M}}\mathbf{Y}$  and  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$  we find

$$\begin{aligned} \hat{v}_5^2 &= R'(\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{g=1}^G \mathbf{X}'_g \tilde{\mathbf{e}}_g \tilde{\mathbf{e}}'_g \mathbf{X}_g \right) (\mathbf{X}'\mathbf{X})^{-1} R \\ &= \sum_{g=1}^G \mathbf{Z}'_g \tilde{\mathbf{e}}_g \tilde{\mathbf{e}}'_g \mathbf{Z}_g \\ &= \tilde{\mathbf{e}}' \tilde{\mathbf{Z}} \tilde{\mathbf{Z}}' \tilde{\mathbf{e}} \\ &= \mathbf{Y}' \tilde{\mathbf{M}} \tilde{\mathbf{Z}} \tilde{\mathbf{Z}}' \tilde{\mathbf{M}} \mathbf{Y} \\ &= (\mathbf{X}\boldsymbol{\beta} + \mathbf{e})' \tilde{\mathbf{M}} \tilde{\mathbf{Z}} \tilde{\mathbf{Z}}' \tilde{\mathbf{M}} (\mathbf{X}\boldsymbol{\beta} + \mathbf{e}) \end{aligned} \quad (83)$$

$$= (\boldsymbol{\psi} + \mathbf{U})' (\boldsymbol{\psi} + \mathbf{U}) \quad (84)$$

where  $\boldsymbol{\psi} = \tilde{\mathbf{Z}}' \tilde{\mathbf{M}} \mathbf{X} \boldsymbol{\beta}$  and  $\mathbf{U} = \tilde{\mathbf{Z}}' \tilde{\mathbf{M}} \mathbf{e}$ .

Note that  $\mathbf{U} \sim N(0, \mathbf{A})$  where  $\mathbf{A} = \tilde{\mathbf{Z}}' \tilde{\mathbf{M}} \boldsymbol{\Sigma} \tilde{\mathbf{M}}' \tilde{\mathbf{Z}}$ . Let  $r = \text{rank}(\mathbf{A})$ . By the spectral decomposition,

$$\mathbf{A} = \mathbf{H} \begin{bmatrix} \boldsymbol{\Lambda} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{H}'$$

where  $\boldsymbol{\Lambda}$  has the  $r$  non-zero eigenvalues of  $\mathbf{A}$  on the diagonal and  $\mathbf{H}\mathbf{H}' = \mathbf{I}_n$ . Partition  $\mathbf{H} = [\mathbf{H}_1, \mathbf{H}_2]$  conformably and define  $\boldsymbol{\delta}_1 = \boldsymbol{\Lambda}^{-1/2} \mathbf{H}'_1 \boldsymbol{\psi}$ ,  $\boldsymbol{\delta}_2 = \mathbf{H}'_2 \boldsymbol{\psi}$ , and  $\boldsymbol{\xi}_1 = \boldsymbol{\Lambda}^{-1/2} \mathbf{H}'_1 \mathbf{U} \sim N(0, \mathbf{I}_r)$ . Observe that  $\mathbf{H}'_2 \tilde{\mathbf{Z}}' \tilde{\mathbf{M}} \mathbf{e} = \mathbf{0}$  almost surely. Since  $\mathbf{I}_n = \mathbf{H}\mathbf{H}'$ , (84) equals

$$\begin{aligned} (\mathbf{H}'\boldsymbol{\psi} + \mathbf{H}'\mathbf{U})' (\mathbf{H}'\boldsymbol{\psi} + \mathbf{H}'\mathbf{U}) &= \left( \begin{pmatrix} \boldsymbol{\Lambda}^{1/2} \boldsymbol{\delta}_1 \\ \boldsymbol{\delta}_2 \end{pmatrix} + \begin{pmatrix} \boldsymbol{\Lambda}^{1/2} \boldsymbol{\xi}_1 \\ \mathbf{0} \end{pmatrix} \right)' \left( \begin{pmatrix} \boldsymbol{\Lambda}^{1/2} \boldsymbol{\delta}_1 \\ \boldsymbol{\delta}_2 \end{pmatrix} + \begin{pmatrix} \boldsymbol{\Lambda}^{1/2} \boldsymbol{\xi}_1 \\ \mathbf{0} \end{pmatrix} \right) \\ &= \begin{pmatrix} \boldsymbol{\delta}_1 + \boldsymbol{\xi}_1 \\ \boldsymbol{\delta}_2 \end{pmatrix}' \begin{bmatrix} \boldsymbol{\Lambda} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{n-r} \end{bmatrix} \begin{pmatrix} \boldsymbol{\delta}_1 + \boldsymbol{\xi}_1 \\ \boldsymbol{\delta}_2 \end{pmatrix} \\ &= (\boldsymbol{\delta}_1 + \boldsymbol{\xi}_1)' \boldsymbol{\Lambda} (\boldsymbol{\delta}_1 + \boldsymbol{\xi}_1) + \boldsymbol{\delta}'_2 \boldsymbol{\delta}_2 \\ &\geq (\boldsymbol{\delta}_1 + \boldsymbol{\xi}_1)' \boldsymbol{\Lambda} (\boldsymbol{\delta}_1 + \boldsymbol{\xi}_1) \\ &= \sum_{j=1}^r \lambda_j (\xi_j + \delta_j)^2 \end{aligned} \quad (85)$$

where  $\boldsymbol{\xi}_1 = (\xi_1, \dots, \xi_r)$  and  $\boldsymbol{\delta}_1 = (\delta_1, \dots, \delta_r)$ .

We can write  $\tilde{\mathbf{M}} = \mathbf{I}_n - \tilde{\mathbf{P}}$ , where  $\tilde{\mathbf{P}}$  is a delete-one-cluster projection operator which has the property

that its diagonal blocks are zeros. The eigenvalues  $\lambda_j$  satisfy

$$\begin{aligned}
\sum_{j=1}^r \lambda_j &= \text{tr}(\mathbf{A}) \\
&= \text{tr}\left(\bar{\mathbf{Z}}' \widetilde{\mathbf{M}} \boldsymbol{\Sigma} \widetilde{\mathbf{M}}' \bar{\mathbf{Z}}\right) \\
&= \text{tr}\left(\bar{\mathbf{Z}}' \boldsymbol{\Sigma} \bar{\mathbf{Z}}\right) - 2 \text{tr}\left(\bar{\mathbf{Z}}' \widetilde{\mathbf{P}} \boldsymbol{\Sigma} \bar{\mathbf{Z}}\right) + \text{tr}\left(\bar{\mathbf{Z}}' \widetilde{\mathbf{P}} \boldsymbol{\Sigma} \widetilde{\mathbf{P}}' \bar{\mathbf{Z}}\right) \\
&\geq v^2
\end{aligned} \tag{86}$$

because

$$\text{tr}\left(\bar{\mathbf{Z}}' \boldsymbol{\Sigma} \bar{\mathbf{Z}}\right) = \sum_{g=1}^G \mathbf{Z}'_g \boldsymbol{\Sigma}_g \mathbf{Z}_g = v^2$$

and  $\text{tr}\left(\bar{\mathbf{Z}}' \widetilde{\mathbf{P}} \boldsymbol{\Sigma} \bar{\mathbf{Z}}\right) = 0$ , the latter because the matrices  $\bar{\mathbf{Z}}'$  and  $\boldsymbol{\Sigma} \bar{\mathbf{Z}}$  are block diagonal, while  $\widetilde{\mathbf{P}}$  has diagonal blocks equalling zero.

Using  $\widehat{C}_5(c) = \widehat{\theta} \pm \widehat{v}_5 c$ , the definition  $\xi_0 = (\widehat{\theta} - \theta) / v$ , and the derivations (84)-(85), we find that

$$\mathbb{P}\left[\theta \in \widehat{C}_5(c)\right] = \mathbb{P}\left[\frac{(\widehat{\theta} - \theta)^2}{\widehat{v}_5^2} \leq c^2\right] = \mathbb{P}\left[\frac{\xi_0^2}{\sum_{j=1}^r w_j (\xi_j + \delta_j)^2} \leq c^2\right], \tag{87}$$

where  $w_j = \lambda_j / v^2$  satisfy  $\sum_{j=1}^r w_j \geq 1$  by (86), and the  $\xi_j$  are normal random variables satisfying the conditions of Theorem 7 (they are mean zero, unit variance, and  $\xi_j$  are mutually uncorrelated for  $j \geq 1$ ). Hence by Theorem 7, (87) is bounded below by  $F(c; 1, 1)$ .

If the model is assumed to satisfy clusterwise invertibility, then  $\widetilde{\mathbf{M}} \mathbf{X} = 0$ , which implies  $\boldsymbol{\psi} = 0$ ,  $\boldsymbol{\delta}_1 = 0$ , and  $\delta_j = 1$  for  $j = 1, \dots, r$ . In this case, the application of Theorem 7 to bound (87) below by  $F(c; 1, 1)$  does not rely on the numerical argument, justifying the claim in the text.

This holds for all models in  $\mathcal{F}$  and thus establishes (24), completing the proof.  $\blacksquare$

For the proof of Theorem 4 we use an intermediate result which connects full downward bias with zero coverage.

**Theorem 8** *Under Assumptions 1-2, if for some variance estimator  $\widehat{v}^2$  and model class  $\mathcal{F}_a \subset \mathcal{F}$ ,*

$$\inf_{(\mathbf{X}, \boldsymbol{\Sigma}) \in \mathcal{F}_a} \frac{\mathbb{E}[\widehat{v}^2]}{v^2} = 0, \tag{88}$$

*then for  $\widehat{C}(c) = \widehat{\theta} \pm c \widehat{v}$  and any  $0 \leq c < \infty$ ,*

$$\inf_{(\mathbf{X}, \boldsymbol{\Sigma}) \in \mathcal{F}_a} \mathbb{P}\left[\theta \in \widehat{C}(c)\right] = 0.$$

**Proof of Theorem 8:** Set  $\varepsilon > 0$ . Let  $Q = (\widehat{\theta} - \theta)^2 / v^2 \sim \chi_1^2$ . Let  $q$  be the  $\varepsilon$ th quantile of the  $\chi_1^2$  distribution and set  $\eta = q / c^2$ . Define the events  $A = \{\widehat{v}^2 / v^2 \leq \eta\}$  and  $B = \{(\widehat{\theta} - \theta)^2 / \widehat{v}^2 \leq c^2\}$ . They jointly imply the



event  $\{Q \leq c^2\eta\}$ . Thus

$$\mathbb{P}[B \cap A] \leq \mathbb{P}[Q \leq c^2\eta] = \varepsilon. \quad (89)$$

Pick  $(\mathbf{X}, \Sigma) \in \mathcal{F}_a$  so that

$$\frac{\mathbb{E}[\hat{v}^2]}{v^2} \leq \eta\varepsilon, \quad (90)$$

which is feasible by (88). By Markov's inequality and (90),

$$\mathbb{P}[B \cap A^c] \leq \mathbb{P}[A^c] = \mathbb{P}\left[\frac{\hat{v}^2}{v^2} > \eta\right] \leq \frac{\mathbb{E}[\hat{v}^2]}{\eta v^2} \leq \varepsilon. \quad (91)$$

Equations (89) and (91) imply that

$$\mathbb{P}[\theta \in \hat{C}(c)] = \mathbb{P}[B] = \mathbb{P}[B \cap A] + \mathbb{P}[B \cap A^c] \leq 2\varepsilon.$$

As  $\varepsilon$  is arbitrary this establishes the stated result. ■

**Proof of Theorem 4:** Results (25), (26), (27), and (28) follow from Theorem 2, equations (17), (18), (21), and (22), combined with Theorem 8. ■

**Proof of Theorem 5:** Define  $\xi_0$ ,  $\mathbf{Z}$ , and  $\bar{\mathbf{Z}}$  as in (79), (81), and (82). Consider the multivariate regression of  $\mathbf{e}$  on  $\xi_0$ :

$$\mathbf{e} = v\Sigma\mathbf{Z}(\mathbf{Z}'\Sigma\mathbf{Z})^{-1}\xi_0 + \boldsymbol{\varepsilon} \quad (92)$$

where

$$\boldsymbol{\varepsilon} = \left(\mathbf{I}_n - \Sigma\mathbf{Z}(\mathbf{Z}'\Sigma\mathbf{Z})^{-1}\mathbf{Z}'\right)\mathbf{e}.$$

The variables  $\boldsymbol{\varepsilon}$  and  $\xi_0$  are jointly normal, uncorrelated, and therefore independent. The variance matrix of  $\boldsymbol{\varepsilon}$  is

$$\boldsymbol{\Omega} = \Sigma - \Sigma\mathbf{Z}(\mathbf{Z}'\Sigma\mathbf{Z})^{-1}\mathbf{Z}'\Sigma. \quad (93)$$

Combining equation (83) from the proof of Theorem 4 with (92), we find

$$\hat{v}_5^2 = (\mathbf{q} + \boldsymbol{\varepsilon})' \widetilde{\mathbf{M}} \bar{\mathbf{Z}} \bar{\mathbf{Z}}' \widetilde{\mathbf{M}} (\mathbf{q} + \boldsymbol{\varepsilon})$$

where

$$\mathbf{q} = \mathbf{X}\beta + v\Sigma\mathbf{Z}(\mathbf{Z}'\Sigma\mathbf{Z})^{-1}\xi_0.$$

Following similar steps as in the proof of Theorem 4, we find that the squared t-statistic satisfies

$$T_5^2 = \frac{(\hat{\theta} - \theta)^2}{\hat{v}_5^2} = \frac{v^2 \xi_0^2}{(\mathbf{t}_1 + \boldsymbol{\xi}_1)' \boldsymbol{\Phi} (\mathbf{t}_1 + \boldsymbol{\xi}_1) + \mathbf{t}_2' \mathbf{t}_2}$$

where  $\boldsymbol{\Phi}$  are the non-zero eigenvalues of

$$\mathbf{L} = \bar{\mathbf{Z}}' \widetilde{\mathbf{M}} \Sigma_{\boldsymbol{\varepsilon}} \widetilde{\mathbf{M}}' \bar{\mathbf{Z}}, \quad (94)$$

$\mathbf{H}^* = [\mathbf{H}_1^*, \mathbf{H}_2^*]$  are the associated eigenvectors,  $\mathbf{t}_1 = \Phi^{-1/2} \mathbf{H}_1^{*'} \tilde{\mathbf{Z}}' \tilde{\mathbf{M}} \mathbf{q}$ ,  $\mathbf{t}_2 = \mathbf{H}_2^{*'} \tilde{\mathbf{Z}}' \tilde{\mathbf{M}} \mathbf{q}$ , and  $\xi_1 = \Phi^{-1/2} \mathbf{H}_1^{*'} \tilde{\mathbf{Z}}' \tilde{\mathbf{M}} \boldsymbol{\varepsilon} \sim \text{N}(0, \mathbf{I}_r)$ . Thus

$$T_5^2 \geq \frac{v^2 \xi_0^2}{\sum_{j=1}^r \phi_j (t_j + \xi_j)^2}$$

where  $\phi_j$ ,  $t_j$ , and  $\xi_j$  are the components of  $\Phi$ ,  $\mathbf{t}_1$ , and  $\xi_1$ .

Notice that  $\xi_1$  is independent of  $(\xi_0, \mathbf{t}_1)$ . Thus, conditioning on  $(\xi_0, \mathbf{t}_1)$ , by Theorem 3 of Mathew and Nordström (1997) (see equation (72)),  $\sum_{j=1}^r \phi_j (t_j + \xi_j)^2$  stochastically dominates  $\sum_{j=1}^r \phi_j \xi_j^2$ , so

$$\mathbb{P} [T_5^2 \leq x] \geq \mathbb{P} \left[ \frac{v^2 \xi_0^2}{\sum_{j=1}^r \phi_j (t_j + \xi_j)^2} \leq x \right] \geq \mathbb{P} \left[ \frac{v^2 \xi_0^2}{\sum_{j=1}^r \phi_j \xi_j^2} \leq x \right]. \quad (95)$$

The random variables  $(\xi_0, \xi_1, \dots, \xi_r)$  are independent  $\text{N}(0, 1)$ . This is the inequality in (30).

The Satterthwaite approximation states that since  $Q = \sum_{j=1}^r \phi_j \xi_j^2$  is a weighted sum of independent chi-squares with non-negative weights, then  $Q \approx bQ_K/K$ , with

$$b = \sum_{j=1}^G \phi_j = \text{tr}(\mathbf{L}),$$

$$K = \frac{\left( \sum_{j=1}^r \phi_j \right)^2}{\sum_{j=1}^r \phi_j^2} = \frac{(\text{tr}(\mathbf{L}))^2}{\text{tr}(\mathbf{L}\mathbf{L})},$$

and  $Q_K \sim \chi_K^2$ . It follows that  $\xi_0^2 / (Q_K/K)$  has distribution function  $F(x; 1, K)$ . Then the right side of (95) equals

$$\mathbb{P} \left[ \frac{v^2 \xi_0^2}{Q} \leq x \right] \approx \mathbb{P} \left[ \frac{v^2 \xi_0^2}{bQ_K/K} \leq x \right] = \mathbb{P} \left[ \frac{\xi_0^2}{Q_K/K} \leq \frac{bx}{v^2} \right] = F(a^2 x; 1, K). \quad (96)$$

with  $a = \sqrt{b/v^2}$ . This is the approximation in (30) as stated. ■

**Proof of Theorem 6:** Recalling definition (82), with some algebra we can write

$$\tilde{\mathbf{M}}' \tilde{\mathbf{Z}} = \begin{bmatrix} \mathbf{Z}_1 & -\mathbf{X}_1 \mathbf{U}_2 & \cdots & -\mathbf{X}_1 \mathbf{U}_G \\ -\mathbf{X}_2 \mathbf{U}_1 & \mathbf{Z}_2 & \cdots & -\mathbf{X}_2 \mathbf{U}_G \\ \vdots & \vdots & \ddots & \vdots \\ -\mathbf{X}_G \mathbf{U}_1 & -\mathbf{X}_G \mathbf{U}_2 & \cdots & \mathbf{Z}_G \end{bmatrix} = \mathbf{S} - \mathbf{X} \mathbf{U}'.$$

Under  $\boldsymbol{\Sigma} = \mathbf{I}_n$ , (93) equals  $\boldsymbol{\Omega} = \mathbf{I}_n - \mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'$ . Then (94) equals

$$\mathbf{L} = \tilde{\mathbf{Z}}' \tilde{\mathbf{M}} \boldsymbol{\Omega} \tilde{\mathbf{M}}' \tilde{\mathbf{Z}} = (\mathbf{S}' - \mathbf{U} \mathbf{X}') \boldsymbol{\Omega} (\mathbf{S} - \mathbf{X} \mathbf{U}') = \mathbf{S}' \boldsymbol{\Omega} \mathbf{S} + \mathbf{U} \hat{\mathbf{X}}' \hat{\mathbf{X}} \mathbf{U}' - \mathbf{V} \mathbf{U}' - \mathbf{U} \mathbf{V}'$$

where we use  $\hat{\mathbf{X}} = \boldsymbol{\Omega} \mathbf{X}$  and  $\mathbf{V} = \mathbf{S}' \hat{\mathbf{X}}$ .

Using (93) and the relationship  $T = S'Z$  we find

$$\begin{aligned}\text{tr}[L] &= \text{tr}[S'S] - \text{tr}\left[S'Z(Z'Z)^{-1}Z'S\right] + \text{tr}\left[U\hat{X}'\hat{X}U'\right] - \text{tr}[VU'] - \text{tr}[UV'] \\ &= \sum_{g=1}^G \mathbf{s}'_g \mathbf{s}_g - \frac{T'T}{Z'Z} + \text{tr}\left[U'U\hat{X}'\hat{X}\right] - 2\text{tr}[U'V]\end{aligned}$$

which is (34).

Similarly,

$$\begin{aligned}\text{tr}[LL] &= \text{tr}\left[\left(S'\Omega S + U\hat{X}'\hat{X}U' - VU' - UV'\right)\left(S'\Omega S + U\hat{X}'\hat{X}U' - VU' - UV'\right)\right] \\ &= \text{tr}\left[(S'\Omega S)(S'\Omega S)\right] + \text{tr}\left[U\hat{X}'\hat{X}U'U\hat{X}'\hat{X}U'\right] + 2\text{tr}[VU'VU'] \\ &\quad + 2\text{tr}\left[S'\Omega S U\hat{X}'\hat{X}U'\right] - 4\text{tr}\left[S'\Omega S VU'\right] - 4\text{tr}\left[U\hat{X}'\hat{X}U'VU'\right] + 2\text{tr}[VU'UV'].\end{aligned}$$

Using (93) and the relationship  $W = S'SU$  we find that this equals

$$\begin{aligned}&= \text{tr}\left[(S'S)(S'S)\right] - 2\frac{T'S'ST}{Z'Z} + \left(\frac{T'T}{Z'Z}\right)^2 \\ &\quad + \text{tr}\left[\hat{X}'\hat{X}U'U\hat{X}'\hat{X}U'U\right] + 2\text{tr}[V'UV'U] + 2\text{tr}\left[U'W\hat{X}'\hat{X}\right] - 2\frac{T'U\hat{X}'\hat{X}U'T}{Z'Z} \\ &\quad - 4\text{tr}[W'V] + 4\frac{T'VU'T}{Z'Z} - 4\text{tr}\left[U'U\hat{X}'\hat{X}U'V\right] + 2\text{tr}[U'UV'V]\end{aligned}$$

which equals (35). ■

## References

- [1] Andrews, Donald W. K. (1991): "Asymptotic normality of series estimators for nonparametric and semiparametric regression models," *Econometrica*, 59, 307-345.
- [2] Arellano, Manuel (1987): "Computing robust standard errors for within groups estimators," *Oxford Bulletin of Economics and Statistics* 49, 431-434.
- [3] Bell, Robert M., and Daniel F. McCaffrey (2002): "Bias reduction in standard errors for linear regression with multi-stage samples," *Survey Methodology*, 28, 169-181.
- [4] Bera, Anil K., Totok Suprayitno, and Gamini Premaratne (2002): "On some heteroskedasticity-robust estimators of variance-covariance matrix of the least-squares estimators," *Journal of Statistical Planning and Inference*, 108, 121-136.
- [5] Cameron, A. Colin, Jonah B. Gelbach, and Douglas L. Miller (2008): "Bootstrap-based improvements for inference with clustered errors," *Review of Economics and Statistics*, 90, 414-427.
- [6] Canay, Ivan A., Andres Santos, and Azeem M. Shaikh (2021): "The wild bootstrap with a small number of large clusters," *Review of Economics and Statistics*, 103, 346-363.
- [7] Cattaneo, Matias D., Michael Jansson, and Whitney K. Newey (2018): "Inference in linear regression models with many covariates and heteroskedasticity," *Journal of the American Statistical Association*, 113, 1350-1361.
- [8] Chesher, Andrew D. (1989): "Hájek inequalities, measures of leverage, and the size of heteroskedasticity robust Wald tests," *Econometrica*, 57, 971-977.
- [9] Chesher, Andrew D. and Gerard Austin (1991): "The finite-sample distributions of heteroskedasticity robust Wald statistics," *Journal of Econometrics*, 47, 153-173.
- [10] Chesher, Andrew D. and Ian D. Jewitt (1987): "The bias of the heteroskedasticity consistent covariance matrix estimator," *Econometrica*, 55, 1217-1272.
- [11] Cochran, William G. (1977): *Sampling Techniques*, 3rd Edition, Wiley.
- [12] Conley, Timothy G. and Christopher R. Taber (2011): "Inference with 'difference in differences' with a small number of policy changes," *Review of Economics and Statistics*, 93, 113-125.
- [13] Davidson, Russell, and James G. MacKinnon (1993): *Estimation and Inference in Econometrics*, Oxford University Press.
- [14] Djogbenou, Antoine. A., James G. MacKinnon, and Morten Ørregaard Nielsen (2019): "Asymptotic theory and wild bootstrap inference with clustered errors," *Journal of Econometrics*, 212, 393-412.

- [15] Duchesne, Pierre and Pierre Lafaye de Micheaux (2010): “Computing the distribution of quadratic forms: Further comparisons between the Liu-Tang-Zhang approximation and exact methods,” *Computational Statistics and Data Analysis*, 54, 858-862.
- [16] Efron, Bradley (1982): *The Jackknife, the Bootstrap, and Other Resampling Plans*, Society for Industrial and Applied Mathematics.
- [17] Efron, Bradley, and Charles Stein (1981): “The jackknife estimate of variance,” *The Annals of Statistics*, 9, 586-596.
- [18] Eicker, Friedhelm (1963): “Asymptotic normality and consistency of the least squares estimators for families of linear regressions,” *Annals of Mathematical Statistics*, 34, 447-456.
- [19] Ferman, Bruno and Cristine Pinto (2019): “Inference in differences-in-differences with few treated groups and heteroskedasticity,” *Review of Economics and Statistics*, 101, 452-467.
- [20] Hagemann, Andreas (2019): “Placebo inference on treatment effects when the number of clusters is small,” *Journal of Econometrics*, 213, 190-209.
- [21] Hagemann, Andreas (2023): “Inference with a single treated cluster,” working paper.
- [22] Hansen, Bruce E. (2022): *Probability and Statistics for Economists*, Princeton University Press.
- [23] Hinkley, David V. (1977): “Jackknifing in unbalanced situations,” *Technometrics*, 19, 285-292.
- [24] Horn, Roger A. and Charles R. Johnson (2013): *Matrix Analysis*, Second Edition, Cambridge University Press.
- [25] Huber, Peter J. (1967): “The behavior of maximum likelihood estimates under nonstandard conditions,” *Proceedings of the Berkeley Symposium on Mathematical Statistics and Probability*, Lucien M. Le Cam and Jerzy Neyman, editors, 1, 221-223.
- [26] Ibragimov, Rustam and Ulrich K. Müller (2016): “Inference with a few heterogeneous clusters,” *Review of Economics and Statistics*, 98, 83-96.
- [27] Imbens, Guido W. and Michal Kolesár (2016): “Robust standard errors in small samples: Some practical advice,” *Review of Economics and Statistics*, 98, 701-712.
- [28] Imhof, J. P. (1961): “Computing the distribution of quadratic forms in normal variables,” *Biometrika*, 48, 419-426.
- [29] Kline, Patrick, Raffaele Saggio, and Mikkel Solvsten (2020): “Leave-out estimation of variance components,” *Econometrica*, 88, 1859-1898.
- [30] Kolesár, Michal (2023): “Robust standard errors in small samples,” unpublished R vignette.

- [31] Liang, Kung-Yee, and Scott L. Zeger (1986): “Longitudinal Data Analysis Using Generalized Linear Models,” *Biometrika*, 73, 13-22.
- [32] Long, J. Scott, and Laurie H. Ervin (2000): “Using heteroscedasticity consistent standard errors in the linear regression model,” *The American Statistician*, 54, 217-224.
- [33] MacKinnon, James G., and Matthew D. Webb (2020): “Randomization inference for difference-in-differences with few treated clusters,” *Journal of Econometrics*, 218, 435-450.
- [34] MacKinnon, James G. and Halbert White (1985): “Some heteroskedasticity-consistent covariance matrix estimators with improved finite sample properties,” *Journal of Econometrics*, 29, 305-325.
- [35] MacKinnon, James G., Morten Ørregaard Nielsen, and Matthew D. Webb (2023a): “Cluster-robust inference: A guide to empirical practice,” *Journal of Econometrics*, 232, 272-299.
- [36] MacKinnon, James G., Morten Ørregaard Nielsen, and Matthew D. Webb (2023b): “Fast and reliable jackknife and bootstrap methods for cluster-robust inference,” *Journal of Applied Econometrics*.
- [37] MacKinnon, James G., Morten Ørregaard Nielsen, and Matthew D. Webb (2023c): “Leverage, influence, and the jackknife in clustered regression models: Reliable inference using `summclust`,” *Stata Journal*, forthcoming.
- [38] Mathew, Thomas, and Kenneth Nordström (1997): “Inequalities for the probability content of a rotated ellipse and related stochastic domination results,” *The Annals of Applied Probability*, 7, 1106-1117.
- [39] Meng, Xin, Nancy Qian, and Pierre Yared (2015): “The institutional causes of China’s Great Famine, 1959-1961,” *Review of Economic Studies*, 82, 1568-1611.
- [40] Niccodemi, Gianmaria and Tom Wansbeek (2022): “A new estimator for standard errors with a few unbalanced clusters,” *Econometrics*, 10, 6.
- [41] Pötscher, Benedikt M. and David Preinerstorfer (2023): “Valid heteroskedasticity robust testing,” *Econometric Theory*.
- [42] Pustejovsky, James E. and Elizabeth Tipton (2018): “Small sample methods for cluster-robust variance estimation and hypothesis testing in fixed effects models,” *Journal of Business and Economic Statistics*, 36, 672-683.
- [43] Rokicki, Slawa, Jessica Cohen, Günther Fink, Joshua A. Salomon, and Mary Beth Landrum (2018): “Inference with difference-in-differences with a small number of groups: A review, simulation study, and empirical application using SHARE data,” *Medical Care*, 56, 97-105.
- [44] Rust, Keith F. and J. N. K. Rao (1996): “Variance estimation for complex surveys using replication techniques,” *Statistical Methods in Medical Research*, 5, 283-310.

- [45] Satterthwaite, F. E. (1946): "An approximate distribution of estimates of variance components," *Biometrics Bulletin*, 2, 110-114.
- [46] Shao, Jun and Dongsheng Tu (1995): *The Jackknife and Bootstrap*, Springer.
- [47] Tukey, John (1958): "Bias and confidence in not quite large samples," *Annals of Mathematical Statistics*, 29, 614.
- [48] White, Halbert (1980): "A heteroskedasticity-consistent covariance matrix estimator and a direct test for heteroskedasticity," *Econometrica*, 48, 817-838.
- [49] Young, Alwyn (2019): "Channeling Fisher: Randomization tests and the statistical insignificance of seemingly significant experimental results," *Quarterly Journal of Economics*, 134, 557-598.
- [50] Zhang, Fuzhen and Qingling Zhang (2006): "Eigenvalue inequalities for matrix product," *IEEE Transactions on Automatic Control*, 51, 1506-1509.