

Submitted to *Econometrica*

Small sample methods for  
cluster-robust variance estimation and  
hypothesis testing in fixed effects  
models

James E. Pustejovsky and Elizabeth Tipton

February 2, 2016

# SMALL SAMPLE METHODS FOR CLUSTER-ROBUST VARIANCE ESTIMATION AND HYPOTHESIS TESTING IN FIXED EFFECTS MODELS<sup>1</sup>

JAMES E. PUSTEJOVSKY<sup>a</sup> AND ELIZABETH TIPTON<sup>b</sup>

Cluster-robust variance estimation (CRVE) is commonly used to account for heteroskedasticity and un-modeled dependence among the errors in regression models with unobserved effects. Although asymptotically consistent, CRVE can be biased downward when the number of clusters is small, leading to hypothesis tests with overly liberal rejection rates. One solution is to use bias-reduced linearization (BRL), which corrects the CRVE based on a working model, in conjunction with t-tests with Satterthwaite degrees of freedom. We propose a generalization of BRL that can be applied in models with arbitrary sets of fixed effects, where the original BRL method is undefined. We also propose a small-sample test for multiple-parameter hypotheses, which generalizes the Satterthwaite approximation for t-tests. In simulations covering a variety of scenarios, we find that conventional cluster-robust Wald tests can severely under-reject while the proposed small-sample test maintains Type I error close to nominal levels.

KEYWORDS: Robust Standard Errors, Clustering, Fixed-Effects, Small Samples.

## 1. INTRODUCTION

In a wide array of economic analyses, interest centers on the parameters of linear regression models, estimated by ordinary or weighted least squares (OLS/WLS) from a sample of units that are correlated. Such correlation among units can arise from sampling aggregate units (e.g., countries, districts, villages), each of which contains multiple observations; from repeated

---

<sup>1</sup>The authors thank Dan Knopf for helpful discussions about the linear algebra behind the cluster-robust variance estimator.

<sup>a</sup>University of Texas at Austin, Educational Psychology Department, 1912 Speedway, Stop D5800, Austin, TX 78712. Email: [pusto@austin.utexas.edu](mailto:pusto@austin.utexas.edu)

<sup>b</sup>Teachers College, Columbia University, Department of Human Development, 525 W. 120th St., Box 118, New York, NY 10027. Email: [tipton@tc.columbia.edu](mailto:tipton@tc.columbia.edu)

measurement of an outcome on a common set of units, as in panel data; or from model misspecification, as in analysis of regression discontinuity designs (e.g., Lee and Card, 2008). A common approach to inference in these settings is to use a cluster-robust variance estimator (CRVE; Arelano, 1987; Liang and Zeger, 1986; White, 1984). The advantage of CRVEs is that they produce consistent standard errors and test statistics without imposing strong parametric assumptions about the dependence structure of the errors in the model. Instead, the method relies on the weaker assumption that units can be grouped into clusters that are mutually independent. CRVEs are an extension to another economic mainstay, heteroskedasticity-robust variance estimators (Huber, 1967; White, 1980), which are used to account for non-constant variance in regression models with independent errors. In the past decade, use of CRVE has become standard practice for applied micro-economic analysis, as evidenced by coverage in major textbooks and review articles (e.g., Angrist and Pischke, 2009; Cameron and Miller, 2015; Wooldridge, 2010).

As a leading example of the application of CRVEs, consider a study of the effects on employment outcomes of several state-level policy shifts, where the policies were implemented at different time-points in each state. In a difference-in-differences analysis of state-by-year panel data, the policy effects would be parameterized in a regression model that includes indicator variables for each policy shift and perhaps additional demographic controls. It is also common to include fixed effects for states and time-periods in order to control for unobserved confounding in each dimension. The model could be estimated by OLS, with the fixed effects included as indicator variables; more commonly, the effects of the policy indicators would be estimated after absorbing the fixed effects, a computational technique that is also known as the fixed effects or within transformation (Wooldridge, 2010). Standard errors would then be clustered by state to account for residual dependence

in the errors from a given state, and these clustered standard errors would be used to test hypotheses regarding each policy or the set of policies. The need to cluster the standard errors by state, even when including state fixed effects, was highlighted by Bertrand, Duflo and Mullainathan (2004), who showed that to do otherwise can lead to inappropriately small standard errors and hypothesis tests with incorrect rejection rates.

The consistency property of CRVEs is asymptotic in the number of independent clusters (Wooldridge, 2003). Recent methodological work has demonstrated that CRVEs can be biased downward and associated hypothesis tests can have Type-I error rates considerably in excess of nominal levels when based on samples with only a small or moderate number of clusters (e.g., Webb and MacKinnon, 2013). Cameron and Miller (2015) provide a thorough review of this literature, including a discussion of current practice, possible solutions, and open problems. In particular, they demonstrate that small-sample corrections for t-tests implemented in common software packages such as Stata and SAS do not provide adequate control of Type-I error.

Bell and McCaffrey (2002, see also McCaffrey, Bell and Botts, 2001) proposed a method that improves the small-sample properties of CRVEs. Their method, called bias-reduced linearization (BRL), entails adjusting the CRVE so that it is exactly unbiased under a working model specified by the analyst, while also remaining asymptotically consistent under arbitrary true variance structures. Simulations reported by Bell and McCaffrey (2002) demonstrate that the BRL correction serves to reduce the bias of the CRVE even when the working model is misspecified. The same authors also proposed and studied small-sample corrections to single-parameter hypothesis tests using the BRL variance estimator, based on Satterthwaite (Bell and McCaffrey, 2002) or saddlepoint approximations (McCaffrey and Bell, 2006). Angrist and Lavy (2009) applied the BRL correction in an analysis

of a longitudinal cluster-randomized trial with 35 clusters, observing that the bias correction makes a difference for inferences.

Despite a growing body of simulation evidence that BRL performs well (e.g., Imbens and Kolesar, 2015), several problems with the method hinder its wider application. First, Angrist and Pischke (2009) noted that the BRL correction breaks down (i.e., cannot be calculated) in some highly parameterized models, such as state-by-year panels that include fixed effects for states and for years. Second, in models with fixed effects, the magnitude of the BRL adjustment depends on whether it is computed based on the full design matrix used in OLS estimation (i.e., including fixed effect dummies) or after absorbing the fixed effects. Cameron and Miller (2015) noted that other methods of small-sample correction suffer from the same subtle problem of depending on arbitrary computational details. Third, extant methods for hypothesis testing based on BRL are limited to single-parameter constraints (Bell and McCaffrey, 2002; McCaffrey and Bell, 2006) and small-sample methods for multiple-parameter hypothesis tests remain lacking. Multiple-parameter tests are used in a range of applications, including in panel data settings (e.g., Hausman tests for consistency of random effects estimators), seemingly unrelated regression models, and analysis of field experiments with multiple treatment groups.

This paper addresses each of these concerns in turn, with the aim of extending the BRL method so that is suitable for everyday econometric practice, including models with fixed effects. First, we describe an extension to the BRL adjustment that is well-defined in models with arbitrary sets of fixed effects, where existing BRL adjustments break down. Second, we demonstrate how to calculate the BRL adjustments so that they are invariant to whether the regression model is estimated including dummy fixed effects or after absorbing the fixed effects (i.e., using the within estimator) and identify conditions under which first-stage absorption of the

fixed effects can be safely ignored. Finally, we propose a procedure for testing multiple-parameter hypotheses by approximating the sampling distribution of the Wald statistic by Hotelling's  $T^2$  distribution with estimated degrees of freedom. The method is a generalization of the Satterthwaite correction proposed by Bell and McCaffrey (2002) for single parameter constraints.

Our work is related to a stream of recent literature that has examined methods for cluster-robust inference with a small number of clusters. Conley and Taber (2011) proposed methods for hypothesis testing in a difference-in-differences setting where the number of treated units is small and fixed, while the number of untreated units increases asymptotically. Ibragimov and Müller (2010) proposed a method for constructing robust tests of scalar parameters that maintains the nominal Type-I error rate; however, their method requires that the target parameter be identified within each independent cluster and so it is not always applicable. Cameron, Gelbach and Miller (2008) investigated a range of bootstrapping procedures that provide improved Type-I error control in small samples, finding that a cluster wild-bootstrap technique was particularly accurate in small samples. Nearly all of this work has focused on single-parameter hypothesis tests only. For multiple-parameter constraints, Cameron and Miller (2015) suggest an ad-hoc degrees of freedom adjustment and note, as an alternative, that bootstrapping techniques can in principle be applied to multiple-parameter tests. However, little methodological work has examined the accuracy of multiple-parameter tests.

The paper is organized as follows. The remainder of this section introduces our econometric framework and reviews the standard CRVE methods, as implemented in most software applications. Section 2 reviews the original BRL correction and describes modifications that make it possible to implement BRL in a broad class of models with fixed effects. Section 3 discusses methods for hypothesis testing based on the BRL-adjusted CRVE.

Section 4 reports a simulation study examining the null rejection rates of our proposed test for multiple-parameter constraints, where we find that the small-sample test offers drastic improvements over commonly implemented alternatives. Section 5 illustrates the use of the proposed hypothesis tests in three examples that cover a variety of contexts in which CRVE is commonly used. Section 6 concludes and discusses avenues for future work.

### 1.1. *Econometric framework*

We consider a linear regression model of the form,

$$(1) \quad y_{ij} = \mathbf{r}_{ij}'\boldsymbol{\beta} + \mathbf{s}_{ij}'\boldsymbol{\gamma} + \mathbf{t}_{ij}'\boldsymbol{\mu} + \epsilon_{ij}$$

where for observation  $j$  in cluster  $i$ ,  $\mathbf{r}_{ij}$  is a vector of  $r$  predictors of primary interest (e.g., policy variables) and any additional controls,  $\mathbf{s}_{ij}$  is a vector of  $s$  fixed effects that vary across clusters, and  $\mathbf{t}_{ij}$  is a vector of  $t$  fixed effects that are identified within clusters. In the state-policy example described in the introduction, the  $\mathbf{r}_{ij}$  would include indicator variables for each policy change, as well as additional demographic controls;  $\mathbf{s}_{ij}$  would include year fixed effects; and  $\mathbf{t}_{ij}$  would indicate state fixed effects. Interest would center on testing hypotheses regarding the coefficients in  $\boldsymbol{\beta}$  that correspond to the policy indicators, while  $\boldsymbol{\gamma}$  and  $\boldsymbol{\mu}$  would be treated as incidental.

For developing theory, it is often easier to work with the matrix version of this model, in which

$$(2) \quad \mathbf{y}_i = \mathbf{R}_i\boldsymbol{\beta} + \mathbf{S}_i\boldsymbol{\gamma} + \mathbf{T}_i\boldsymbol{\mu} + \boldsymbol{\epsilon}_i,$$

where for cluster  $i$ ,  $\mathbf{R}_i$  is an  $n_i \times r$  matrix of focal predictors and controls;  $\mathbf{S}_i$  is an  $n_i \times s$  matrix describing fixed effects that vary across clusters, and  $\mathbf{T}_i$  is an  $n_i \times t$  matrix describing fixed effects that are identified only within clusters. The distinction between the covariates  $\mathbf{R}_i$  versus the fixed effects  $\mathbf{S}_i$  is arbitrary and depends on the analyst's inferential goals. However, the

distinction between the two fixed effect matrices  $\mathbf{S}_i$  and  $\mathbf{T}_i$  is unambiguous, in that the within-cluster fixed effects satisfy  $\mathbf{T}_h \mathbf{T}'_i = \mathbf{0}$  for  $h \neq i$ .

We shall assume that  $E(\boldsymbol{\epsilon}_i | \mathbf{R}_i, \mathbf{S}_i, \mathbf{T}_i) = \mathbf{0}$  and  $\text{Var}(\boldsymbol{\epsilon}_i | \mathbf{R}_i, \mathbf{S}_i, \mathbf{T}_i) = \boldsymbol{\Sigma}_i$ , for  $i = 1, \dots, m$ , where the form of  $\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_m$  may be unknown but the errors are independent across clusters. For notational convenience, let  $\mathbf{U}_i = [\mathbf{R}_i \ \mathbf{S}_i]$  denote the set of predictors that vary across clusters,  $\mathbf{X}_i = [\mathbf{R}_i \ \mathbf{S}_i \ \mathbf{T}_i]$  denote the full set of predictors,  $\boldsymbol{\alpha} = (\boldsymbol{\beta}', \boldsymbol{\gamma}', \boldsymbol{\mu}')'$ , and  $p = r + s + t$ . Denote the total number of individual observations by  $N = \sum_{i=1}^m n_i$ . Let  $\mathbf{y}$ ,  $\mathbf{R}$ ,  $\mathbf{S}$ ,  $\mathbf{T}$ ,  $\mathbf{U}$ ,  $\mathbf{X}$ , and  $\boldsymbol{\epsilon}$  denote the matrices obtained by stacking their corresponding components, as in  $\mathbf{R} = (\mathbf{R}'_1 \ \mathbf{R}'_2 \ \dots \ \mathbf{R}'_m)'$ .

We assume that  $\boldsymbol{\beta}$  is estimated by weighted least squares (WLS) using symmetric, full rank weighting matrices  $\mathbf{W}_1, \dots, \mathbf{W}_m$ . Clearly, the WLS estimator includes OLS as a special case (where  $\mathbf{W}_i = \mathbf{I}_i$ , an identity matrix), as well as feasible GLS.<sup>1</sup> In the latter case, it is assumed that  $\text{Var}(\mathbf{e}_i | \mathbf{X}_i) = \boldsymbol{\Phi}_i$ , where  $\boldsymbol{\Phi}_i$  is a known function of a low-dimensional parameter. For example, an auto-regressive error structure might be posited to describe repeated measures on an individual over time. The weighting matrices are then taken to be  $\mathbf{W}_i = \hat{\boldsymbol{\Phi}}_i^{-1}$ , where the  $\hat{\boldsymbol{\Phi}}_i$  are constructed from estimates of the variance parameter. Finally, for analysis of data from complex survey designs, WLS may be used with sampling weights in order to account for unequal selection probabilities.

---

<sup>1</sup>The WLS estimator also encompasses the estimator proposed by Ibragimov and Müller (2010) for clustered data. Assuming that  $\mathbf{X}_i$  has rank  $p$  for  $i = 1, \dots, m$ , their proposed approach involves estimating  $\boldsymbol{\beta}$  separately within each cluster and taking the simple average of these estimates. The resulting average is equivalent to the WLS estimator with weights  $\mathbf{W}_i = \mathbf{X}_i (\mathbf{X}'_i \mathbf{X}_i)^{-2} \mathbf{X}_i$ .



## 1.2. Absorption

The goal of most analyses is to estimate and test hypotheses regarding the parameters in  $\beta$ , while the fixed effects  $\gamma$  and  $\mu$  are not of inferential interest. Moreover, estimating all of the parameters by WLS becomes computationally intensive and numerically inaccurate if the model includes a large number of fixed effects (i.e.,  $s+t$  large). A commonly implemented solution is to first absorb the fixed effects, which leaves only the  $r$  parameters in  $\beta$  to be estimated. Section 2 examines the implications of absorption for application of the BRL adjustment. In order to do, we now formalize the absorption method.

To begin, denote the full block-diagonal weighting matrix as  $\mathbf{W} = \text{diag}(\mathbf{W}_1, \dots, \mathbf{W}_m)$ . Let  $\mathbf{K}$  be the  $x \times r$  matrix that selects the covariates of interest, so that  $\mathbf{XK} = \mathbf{R}$  and  $\mathbf{K}'\alpha = \beta$ . For a generic matrix  $\mathbf{Z}$  of full column rank, let  $\mathbf{M}_Z = (\mathbf{Z}'\mathbf{WZ})^{-1}$  and  $\mathbf{H}_Z = \mathbf{ZM}_Z\mathbf{Z}'\mathbf{W}$ .

The absorption technique involves obtaining the residuals from the regression of  $\mathbf{y}$  on  $\mathbf{T}$  and from the multivariate regressions of  $\mathbf{U} = [\mathbf{R} \ \mathbf{S}]$  on  $\mathbf{T}$ . The  $\mathbf{y}$  residuals and  $\mathbf{R}$  residuals are then regressed on the  $\mathbf{S}$  residuals. Finally, these twice-regressed  $\mathbf{y}$  residuals are regressed on the twice-regressed  $\mathbf{R}$  residuals to obtain the WLS estimates of  $\beta$ . Let  $\ddot{\mathbf{S}} = (\mathbf{I} - \mathbf{H}_T) \mathbf{S}$ ,  $\ddot{\mathbf{R}} = (\mathbf{I} - \mathbf{H}_{\ddot{\mathbf{S}}}) (\mathbf{I} - \mathbf{H}_T) \mathbf{R}$ , and  $\ddot{\mathbf{y}} = (\mathbf{I} - \mathbf{H}_{\ddot{\mathbf{S}}}) (\mathbf{I} - \mathbf{H}_T) \mathbf{y}$ . In what follows, subscripts on  $\ddot{\mathbf{R}}$ ,  $\ddot{\mathbf{S}}$ ,  $\ddot{\mathbf{U}}$ , and  $\ddot{\mathbf{y}}$  refer to the rows of these matrices corresponding to a specific cluster. The WLS estimator of  $\beta$  can then be written as

$$(3) \quad \hat{\beta} = \mathbf{M}_{\ddot{\mathbf{R}}} \sum_{i=1}^m \ddot{\mathbf{R}}'_i \mathbf{W}_i \ddot{\mathbf{y}}_i.$$

This estimator is algebraically identical to the direct WLS estimator based on the full set of predictors,

$$\hat{\beta} = \mathbf{K}'\mathbf{M}_X \sum_{i=1}^m \mathbf{X}'_i \mathbf{W}_i \mathbf{y}_i,$$

but avoids the need to solve a system of  $p$  linear equations.

In the remainder, we focus on the more general case in which fixed effects are absorbed before estimation of  $\beta$ . For models that do not include within-cluster fixed effects, so that the full covariate matrix is  $\mathbf{U} = [\mathbf{R} \ \mathbf{S}]$ , all of the results hold after substituting  $\mathbf{U}$  for  $\ddot{\mathbf{R}}$ .

### 1.3. Standard CRVE

The WLS estimator  $\hat{\beta}$ , has true variance

$$(4) \quad \text{Var}(\hat{\beta}) = \mathbf{M}_{\ddot{\mathbf{R}}} \left( \sum_{i=1}^m \ddot{\mathbf{R}}_i' \mathbf{W}_i \boldsymbol{\Sigma}_i \mathbf{W}_i \ddot{\mathbf{R}}_i \right) \mathbf{M}_{\ddot{\mathbf{R}}},$$

which depends upon the unknown variance matrices  $\boldsymbol{\Sigma}_i$ . A model-based approach to estimating this variance would involve assuming that  $\boldsymbol{\Sigma}_i$  follows a structure defined by some low-dimensional parameter; for example, it may be assumed that the structure was hierarchical or auto-regressive. The model-based variance estimator would substitute estimates of  $\boldsymbol{\Sigma}_i$  into (4). However, if the model is misspecified, this estimator will be inconsistent and inferences based upon it will be invalid.

The CRVE involves estimating  $\text{Var}(\hat{\beta})$  empirically, without imposing structural assumptions on  $\boldsymbol{\Sigma}_i$ . While there are several versions of this approach, all can be written in the form

$$(5) \quad \mathbf{V}^{CR} = \mathbf{M}_{\ddot{\mathbf{R}}} \left( \sum_{i=1}^m \ddot{\mathbf{R}}_i' \mathbf{W}_i \mathbf{A}_i \mathbf{e}_i \mathbf{e}_i' \mathbf{A}_i' \mathbf{W}_i \ddot{\mathbf{R}}_i \right) \mathbf{M}_{\ddot{\mathbf{R}}},$$

where  $\mathbf{e}_i = \mathbf{Y}_i - \mathbf{X}_i \hat{\beta}$  is the vector of residuals from cluster  $i$  and  $\mathbf{A}_i$  is some  $n_i$  by  $n_i$  adjustment matrix.

The form of these adjustment matrices parallels those of the heteroskedasticity-consistent (HC) variance estimators proposed by MacKinnon and White (1985). The most basic CRVE, described by Liang and Zeger (1986), uses  $\mathbf{A}_i = \mathbf{I}_i$ , an  $n_i \times n_i$  identity matrix. Following Cameron and Miller (2015),

we refer to this estimator as CR0. This estimator is biased towards zero because the cross-product of the residuals  $\mathbf{e}_i \mathbf{e}_i'$  tends to under-estimate the true variance  $\Sigma_i$  in cluster  $i$ . A rough bias adjustment is to take  $\mathbf{A}_i = c \mathbf{I}_i$ , where  $c = \sqrt{(m/(m-1))}$ ; we denote this adjusted estimator as CR1. Some functions in Stata use a slightly different correction factor  $c_S = \sqrt{(mN)/[(m-1)(N-p)]}$ ; we will refer to the adjusted estimator using  $c_S$  as CR1S. When  $N \gg p$ ,  $c_S \approx \sqrt{m/(m-1)}$  and so CR1 and CR1S will be very similar. The CR1 or CR1S estimator is now commonly used in empirical applications.

Use of these adjustments still tends to under-estimate the true variance of  $\hat{\beta}$  (Cameron and Miller, 2015). Analytic work and simulation studies indicate that the degree of bias depends not only on the number of clusters  $m$ , but also on features of the covariates in  $\mathbf{X}$ . Specifically, the bias tends to be larger when the covariates are skewed or unbalanced across clusters, or when clusters vary in size (Carter, Schnepel and Steigerwald, 2013; MacKinnon, 2013). A more principled approach to bias correction would therefore take into account these features of the covariates. One such estimator uses adjustment matrices given by  $\mathbf{A}_i = \left( \mathbf{I} - \ddot{\mathbf{R}}_i \mathbf{M}_{\ddot{\mathbf{R}}} \ddot{\mathbf{R}}_i' \mathbf{W}_i \right)^{-1}$ . This estimator, denoted CR3, closely approximates the jackknife re-sampling variance estimator (Bell and McCaffrey, 2002; Mancl and DeRouen, 2001). However, CR3 tends to over-correct the bias of CR0, while the CR1 estimator tends to under-correct. The next section describes in detail the BRL approach, which makes adjustments that are intermediate in magnitude between CR1 and CR3.

## 2. BIAS REDUCED LINEARIZATION

In contrast to the CR1 or CR3 estimators, the BRL correction for CRVE is premised on a “working” model for the structure of the errors, which must be specified by the analyst. Under a given working model, adjustment

matrices  $\mathbf{A}_i$  are defined so that the variance estimator is exactly unbiased. We refer to this correction as CR2 because it is an extension of the HC2 variance estimator for regressions with uncorrelated errors, which is exactly unbiased when the errors are homoskedastic (MacKinnon and White, 1985). The idea of specifying a model may seem antithetical to the purpose of using CRVE, yet extensive simulation studies have demonstrated that the method performs better in small samples than any of the other adjustments, even when the working model is incorrect. (See Section 4 for a review of this literature.) Although the CR2 estimator may no longer be exactly unbiased when the working model is misspecified, its bias still tends to be greatly reduced compared to CR1 or CR0 (thus the name “bias reduced linearization”). Furthermore, as the number of clusters increases, reliance on the working model diminishes. In a sense, CR2 provides necessary scaffolding in the small-sample case, which falls away when there is sufficient data.

Let  $\Phi_i$  denote a working model for the covariance of the errors in cluster  $i$  (up to a scalar constant), with  $\Phi = \text{diag}(\Phi_1, \dots, \Phi_m)$ . For example, following Bell and McCaffrey (2002) we might assume  $\Phi_i = \mathbf{I}_i$ , i.e., that the errors are uncorrelated and homoskedastic. Alternatively, Imbens and Kolesar (2015) suggested using a basic random effects (i.e., compound symmetric) structure, in which  $\Phi_i$  has unit diagonal entries and off-diagonal entries of  $\rho$ , with  $\rho$  estimated using the OLS residuals (see Imbens and Kolesar, 2015, p. 16).

Based on a given working model, in the original formulation of Bell and McCaffrey (2002), the BRL adjustment matrices are chosen to satisfy the criterion

$$(6) \quad \mathbf{A}_i (\mathbf{I} - \mathbf{H}_\mathbf{X})_i \Phi (\mathbf{I} - \mathbf{H}_\mathbf{X})_i' \mathbf{A}_i' = \Phi_i$$

where  $(\mathbf{I} - \mathbf{H}_\mathbf{X})_i$  denotes the rows of  $\mathbf{I} - \mathbf{H}_\mathbf{X}$  corresponding to cluster  $i$ . If the working model and weight matrices are both taken to be identity matrices,

then the adjustment matrices simplify to  $\mathbf{A}_i = (\mathbf{I}_i - \mathbf{X}_i \mathbf{M}_{\mathbf{X}} \mathbf{X}_i')^{-1/2}$ , where  $\mathbf{Z}^{-1/2}$  denotes the symmetric square-root of the matrix  $\mathbf{Z}$ . This formulation of  $\mathbf{A}_i$  is problematic because, for some fixed effects models that are common in economic applications, Equation 6 does not have a solution. In the next two subsections, we address two problems that arise in models with fixed effects, thereby articulating a BRL methodology that is suitable for a wide range of applications.

### 2.1. Generalized Inverse

The equality defining the  $\mathbf{A}_i$  matrices cannot always be solved because it is possible that some of the matrices involved are not of full rank, and thus cannot be inverted. Angrist and Pischke (2009) note that this problem arises in balanced state-by-year panel models that include fixed effects for states and for years. In order to address this concern, we provide an alternative criterion for the adjustment matrices that can always be satisfied. Instead of criterion (6), we seek adjustment matrices  $\mathbf{A}_i$  that satisfy:

$$(7) \quad \ddot{\mathbf{R}}_i' \mathbf{W}_i \mathbf{A}_i (\mathbf{I} - \mathbf{H}_{\mathbf{X}})_i \Phi (\mathbf{I} - \mathbf{H}_{\mathbf{X}})_i' \mathbf{A}_i' \mathbf{W}_i \ddot{\mathbf{R}}_i = \ddot{\mathbf{R}}_i' \mathbf{W}_i \Phi \mathbf{W}_i \ddot{\mathbf{R}}_i.$$

A variance estimator that uses such adjustment matrices will be exactly unbiased when the working model is correctly specified.

The above criterion (7) does not uniquely define  $\mathbf{A}_i$ . One solution, which produces symmetric adjustment matrices, uses

$$(8) \quad \mathbf{A}_i = \mathbf{D}_i' \mathbf{B}_i^{+1/2} \mathbf{D}_i,$$

where  $\mathbf{D}_i$  is the upper-right triangular Cholesky factorization of  $\Phi_i$ ,

$$(9) \quad \mathbf{B}_i = \mathbf{D}_i (\mathbf{I} - \mathbf{H}_{\mathbf{X}})_i \Phi (\mathbf{I} - \mathbf{H}_{\mathbf{X}})_i' \mathbf{D}_i',$$

and  $\mathbf{B}_i^{+1/2}$  is the symmetric square root of the Moore-Penrose inverse of  $\mathbf{B}_i$ . The Moore-Penrose inverse is well-defined and unique even when  $\mathbf{B}_i$  is

not of full rank (Banerjee and Roy, 2014, Thm. 9.18). These adjustment matrices satisfy criterion (7), as stated in the following theorem.

**THEOREM 1** *Let  $\mathbf{L}_i = (\ddot{\mathbf{U}}'\ddot{\mathbf{U}} - \ddot{\mathbf{U}}_i'\ddot{\mathbf{U}}_i)$ , where  $\ddot{\mathbf{U}} = (\mathbf{I} - \mathbf{H}_{\mathbf{T}})\mathbf{U}$ , and assume that  $\mathbf{L}_1, \dots, \mathbf{L}_m$  have full rank  $r+s$ . Further assume that  $\text{Var}(\boldsymbol{\epsilon}_i | \mathbf{R}_i, \mathbf{S}_i, \mathbf{T}_i) = \sigma^2 \boldsymbol{\Phi}_i$ , for  $i = 1, \dots, m$ . Then the adjustment matrix  $\mathbf{A}_i$  defined in (8) and (9) satisfies criterion (7) and the CR2 variance estimator is exactly unbiased.*

Proof is given in Appendix A. If  $\mathbf{B}_i$  is of full rank, then the adjustment matrices also satisfy the original criterion (6). The main implication of Theorem 1 is that the CR2 variance estimator remains well-defined, even in models with large sets of fixed effects.

## 2.2. Absorption and Dummy Equivalence

A problem with existing small-sample adjustments to CRVEs is that they can result in a different estimator depending upon if the fixed effects are estimated by OLS or are first absorbed. For example, this problem arises with the CR1S estimator because it uses a multiplicative correction to the residuals that depends on the total number of covariates estimated in the model. When fixed effects are included as indicators, the constant is calculated as  $c_S = \sqrt{(mN)/[(m-1)(N-p)]}$ , where  $p$  is the total number of covariates, including fixed effects. In contrast, if the fixed effects are absorbed, the constant is calculated as  $c_S = \sqrt{(mN)/[(m-1)(N-r)]}$ , where  $r$  is the number of covariates that are not absorbed. Cameron and Miller (2015) highlight that this discrepancy can be substantial if the clusters are small. For instance, if each cluster includes  $n_i = 2$  units, then the CR1S estimator based on estimating the fixed effects by OLS is over twice as large as the estimator based on the absorbed model. Such differences between the correction based on OLS estimation of the fixed effects and the correction based on the absorbed model are problematic because the magnitude of

the variance estimator should not depend on how the model estimates are calculated.

Similar inconsistencies can arise when applying the BRL method in models with fixed effects. Consider the scenario in which absorption is used to estimate  $\beta$ ; here, the analyst might choose to calculate the CR2 correction based on the absorbed covariate matrix  $\ddot{\mathbf{R}}$ —that is, by substituting  $\mathbf{H}_{\ddot{\mathbf{R}}}$  for  $\mathbf{H}_{\mathbf{X}}$  in (9)—in order to avoid calculating the full projection matrix  $\mathbf{H}_{\mathbf{X}}$ . However, this approach can lead to differences in the adjustment matrices compared to when the full model is estimated by OLS because it is based on a subtly different working model. Essentially, calculating CR2 based on the absorbed model amounts to assuming that the working model  $\Phi$  applies not to the model errors  $\epsilon$ , but rather to the errors from the regression of  $\ddot{\mathbf{y}}$  on  $\ddot{\mathbf{R}}$ . We find this method of specifying the working model to be incoherent, and therefore recommend against taking it. Rather, the CR2 adjustment matrices should be calculated based on a working model for the errors in the full regression model, following Equations (8) and (9) as stated.

A drawback of using the CR2 adjustment matrices based on the full regression model is that it entails calculating the projection matrix  $\mathbf{H}_{\mathbf{X}}$  for the full set of  $p$  covariates (i.e., including fixed effect indicators). Given that the entire advantage of using absorption to calculate  $\hat{\beta}$  is to avoid computations involving large, sparse matrices, it is of interest to find methods for more efficiently calculating the CR2 adjustment matrices. Some efficiency can be gained by using the fact that the residual projection matrix  $\mathbf{I} - \mathbf{H}_{\mathbf{X}}$  can be factored into components as  $(\mathbf{I} - \mathbf{H}_{\mathbf{X}})_i = (\mathbf{I} - \mathbf{H}_{\ddot{\mathbf{R}}})_i (\mathbf{I} - \mathbf{H}_{\ddot{\mathbf{S}}}) (\mathbf{I} - \mathbf{H}_{\mathbf{T}})$ .

In certain cases, further computational efficiency can be achieved by computing the adjustment matrices after absorbing the within-cluster fixed effects  $\mathbf{T}$  (but not the between-cluster fixed effects  $\mathbf{S}$ ). Specifically, if the weights used for WLS estimation are the inverses of the working covariance model, so that  $\mathbf{W}_i = \Phi_i^{-1}$  for  $i = 1, \dots, m$ , then the adjustment matrices can

be calculated without accounting for the within-cluster fixed effects. This result is formalized in the following theorem.

**THEOREM 2** *Let  $\tilde{\mathbf{A}}_i = \mathbf{D}_i' \tilde{\mathbf{B}}_i^{+1/2} \mathbf{D}_i$ , where*

$$(10) \quad \tilde{\mathbf{B}}_i = \mathbf{D}_i (\mathbf{I} - \mathbf{H}_{\tilde{\mathbf{R}}})_i (\mathbf{I} - \mathbf{H}_{\tilde{\mathbf{S}}}) \boldsymbol{\Phi} (\mathbf{I} - \mathbf{H}_{\tilde{\mathbf{S}}})' (\mathbf{I} - \mathbf{H}_{\tilde{\mathbf{R}}})_i' \mathbf{D}_i'.$$

*If  $\mathbf{T}_i \mathbf{T}_k' = \mathbf{0}$  for  $j \neq k$  and  $\mathbf{W} = \boldsymbol{\Phi}^{-1}$ , then  $\mathbf{A}_i = \tilde{\mathbf{A}}_i$ .*

Proof is given in Appendix A. The main implication of Theorem 2 is that the more computationally convenient formula  $\tilde{\mathbf{B}}_i$  can be applied in the common case that the weighting matrices are the inverse of the working covariance model.

Following the working model suggested by Bell and McCaffrey (2002), in which  $\boldsymbol{\Phi} = \mathbf{I}$ , the above theorem shows that the adjustment method is invariant to the choice of method for dealing with fixed effects so long as the model is estimated by OLS (i.e.,  $\mathbf{W} = \mathbf{I}$ ). In this case, the CR2 adjustment matrices then simplify further to

$$\mathbf{A}_i = \left( \mathbf{I}_i - \ddot{\mathbf{U}}_i \left( \ddot{\mathbf{U}}' \ddot{\mathbf{U}} \right)^{-1} \ddot{\mathbf{U}}_i' \right)^{+1/2}.$$

In contrast, if the working model proposed by Imbens and Kolesar (2015) is instead used, then the above theorem implies that the CR2 adjustments will differ if the model is estimated by OLS with dummies for fixed effects versus by using absorption.

The two theorems of this section extend the BRL methodology described by Bell and McCaffrey (2002), demonstrating how the CR2 adjustment can be computed efficiently—and from a coherent working model—for a broad range of commonly used regression models, including those with within- and between-cluster fixed effects. The next section addresses a final set of concerns: how to conduct single- and multiple-parameter hypothesis tests using the CR2 estimator.



## 3. HYPOTHESIS TESTING

The CR2 correction produces a CRVE that has reduced bias (compared to other CRVEs) when the number of clusters is small, leading to more accurate standard errors. However, standard errors are of limited inherent interest—rather, their main use is for the construction of hypothesis tests and confidence intervals. Cluster-robust Wald-type test statistics are a function of the parameter estimates  $\hat{\beta}$  and the corresponding CRVE. Conventional Wald tests are justified based on the asymptotic behavior of robust Wald statistics as the number of clusters grows large (i.e., as  $m \rightarrow \infty$ ).

Like the research on the bias of the CRVE estimator, evidence from a wide variety of contexts indicates that the asymptotic limiting distribution of these statistics may be a poor approximation when the number of clusters is small, even if corrections such as CR2 or CR3 are employed (Bell and McCaffrey, 2002; Bertrand et al., 2004; Cameron et al., 2008). Like the bias of the CRVE estimator itself, the accuracy of the asymptotic approximations depends on design features such as the degree of imbalance across clusters, skewness or leverage of the covariates, and the similarity of cluster sizes (Carter et al., 2013; McCaffrey et al., 2001; Tipton and Pustejovsky, 2015; Webb and MacKinnon, 2013). This provides motivation for development of general-purpose hypothesis testing procedures that have accurate rejection rates in small samples.

In this section, we develop a general method for conducting hypothesis tests based on CRVEs. We consider linear constraints on  $\beta$ , where the null hypothesis has the form  $H_0 : \mathbf{C}\beta = \mathbf{d}$  for fixed  $q \times r$  matrix  $\mathbf{C}$  and  $q \times 1$  vector  $\mathbf{d}$ . The cluster-robust Wald statistic is then

$$(11) \quad Q = (\mathbf{C}\hat{\beta} - \mathbf{d})' (\mathbf{C}\mathbf{V}^{CR}\mathbf{C}')^{-1} (\mathbf{C}\hat{\beta} - \mathbf{d}),$$

where  $\mathbf{V}^{CR}$  is one of the cluster-robust estimators described in previous sections. The asymptotic Wald test rejects  $H_0$  if  $Q$  exceeds the  $\alpha$  critical

value from a chi-squared distribution with  $q$  degrees of freedom. In large samples, it can be shown that this test has level  $\alpha$ . However, in practice it is rarely clear how large a sample is needed for the asymptotic approximation to be accurate.

### 3.1. *Small-sample corrections for t-tests*

Consider testing the hypothesis  $H_0 : \mathbf{c}'\boldsymbol{\beta} = 0$  for a fixed  $r \times 1$  contrast vector  $\mathbf{c}$ . For this one-dimensional constraint, an equivalent to the Wald statistic given in (11) is to use the test statistic  $Z = \mathbf{c}'\hat{\boldsymbol{\beta}}/\sqrt{\mathbf{c}'\mathbf{V}^{CR}\mathbf{c}}$ , which follows a standard normal distribution in large samples. In small samples, it is common to use the CR1 or CR1S estimator and to approximate the distribution of  $Z$  by a  $t(m-1)$  distribution. Hansen (2007) provided one justification for the use of this reference distribution by identifying conditions under which  $Z$  converges in distribution to  $t(m-1)$  as the within-cluster sample sizes grow large, with  $m$  fixed (see also Donald and Lang, 2007). Ibragimov and Müller (2010) proposed a weighting technique derived so that that  $t(m-1)$  critical values are conservative (leading to rejection rates less than or equal to  $\alpha$ ). However, both of these arguments require that  $\mathbf{c}'\boldsymbol{\beta}$  be separately identified within each cluster. Outside of these circumstances, using  $t(m-1)$  critical values can still lead to over-rejection (Cameron and Miller, 2015). Furthermore, using these critical values does not take into account that the distribution of  $\mathbf{V}^{CR}$  is affected by the structure of the covariate matrix.

An alternative t-test developed by Bell and McCaffrey (2002) involves using a  $t(\nu)$  references distribution, with degrees of freedom  $\nu$  estimated by a Satterthwaite approximation. The Satterthwaite approximation (Satterthwaite, 1946) entails using degrees of freedom that are a function of the the first two moments of the sampling distribution of  $\mathbf{c}'\mathbf{V}^{CR}\mathbf{c}$ . Theoretically,

these degrees of freedom should be

$$(12) \quad \nu = \frac{2 [E (\mathbf{c}' \mathbf{V}^{CR2} \mathbf{c})]^2}{\text{Var} (\mathbf{c}' \mathbf{V}^{CR2} \mathbf{c})}.$$

Expressions for the first two moments of  $\mathbf{c}' \mathbf{V}^{CR2} \mathbf{c}$  can be derived under the assumption that the errors  $\epsilon_1, \dots, \epsilon_m$  are normally distributed.

In practice, both moments involve the variance structure  $\Sigma$ , which is unknown. Bell and McCaffrey (2002) proposed to estimate the moments based on the same working model that is used to derive the adjustment matrices. This “model-assisted” estimate of the degrees of freedom is then calculated as

$$(13) \quad \nu_M = \frac{(\sum_{i=1}^m \mathbf{p}_i' \Phi \mathbf{p}_i)^2}{\sum_{i=1}^m \sum_{j=1}^m (\mathbf{p}_i' \Phi \mathbf{p}_j)^2},$$

where  $\mathbf{p}_i = (\mathbf{I} - \mathbf{H}_X)'_i \mathbf{A}_i \mathbf{W}_i \ddot{\mathbf{R}}_i \mathbf{M}_{\ddot{\mathbf{R}}} \mathbf{c}$ . Alternately, for any of the CRVEs one could instead use an “empirical” estimate of the degrees of freedom, constructed by substituting  $\mathbf{e}_i \mathbf{e}_i'$  in place of  $\Sigma_i$ . However, Bell and McCaffrey (2002) found using simulation that this plug-in degrees of freedom estimate led to very conservative rejection rates.

The Bell and McCaffrey (2002) Satterthwaite approximation has been shown to perform well in a variety of conditions (see Section 4). These studies encompass a variety of data generation processes, covariate types, and weighting procedures. A key finding is that the degrees of freedom depend not only on the number of clusters  $m$ , but also on features of the covariates. When the covariate is balanced across clusters—as occurs in balanced panels with a dichotomous covariate with the same proportion of ones in each cluster—the degrees of freedom are  $m - 1$  even in small samples. However, when the covariate is highly unbalanced—as occurs when the panel is not balanced or if the proportion of ones varies from cluster to cluster—the degrees of freedom can be considerably smaller. Similarly,

Can we use  $\mathbf{H}_{\ddot{\mathbf{U}}}$  here instead?

covariates with large leverage points will tend to exhibit lower of degrees of freedom.

By adjusting the degrees of freedom to account for these features, the Type I error rate of the test is nearly always less than or equal to the nominal  $\alpha$ , so long as the degrees of freedom are larger than 4 or 5 (Bell and McCaffrey, 2002; Tipton, 2015). This is because when the degrees of freedom are smaller, the t-distribution approximation to the sampling distribution does not hold, and the Type I error can be higher than the stated  $\alpha$  level.<sup>2</sup> In comparison, the CR1 degrees of freedom (i.e.,  $m - 1$ ) are constant, and the test only performs well when in the cases in which the covariates are balanced. Because the degrees of freedom are covariate-dependent, it is not possible to assess whether a small-sample correction is needed based solely on the total number of clusters in the data. Consequently, Tipton (2015) argued that t-tests based on CRVE should routinely use the CR2 variance estimator and the Satterthwaite degrees of freedom, even when  $m$  appears to be large.

### 3.2. *Small-sample corrections for F-tests*

Little research has considered small-sample corrections for multiple-constraint hypothesis tests based on cluster-robust Wald statistics. Cameron and Miller highlight this problem, noting that some form of adjustment is clearly needed in light of the extensive work on single-parameter tests. We now describe an approach to multi-parameter testing that closely parallels the

---

<sup>2</sup>When the degrees of freedom are smaller than 4 or 5, Tipton (2015) suggested using a smaller  $\alpha$  level for hypothesis testing in order to partially compensate. Although degrees of freedom this small may seem unlikely, they can easily arise in practice even when the number of clusters is moderate. For example, in a state-by-year panel with  $m = 48$  states, the degrees of freedom can be quite small if a policy is only implemented in 10 percent of states.

Satterthwaite correction for t-tests.

Our approach is to approximate the sampling distribution of  $Q$  by Hotelling's  $T^2$  distribution (a multiple of an F distribution) with estimated degrees of freedom. To motivate the approximation, let  $\mathbf{G} = \mathbf{C}\mathbf{M}_{\ddot{\mathbf{R}}}\ddot{\mathbf{R}}'\mathbf{W}\Phi\mathbf{W}\ddot{\mathbf{R}}\mathbf{M}_{\ddot{\mathbf{R}}}'\mathbf{C}'$  denote the variance of  $\mathbf{C}\hat{\boldsymbol{\beta}}$  under the working model and observe that  $Q$  can be written as

$$Q = \mathbf{z}'\boldsymbol{\Omega}^{-1}\mathbf{z},$$

where  $\mathbf{z} = \mathbf{G}^{-1/2}(\mathbf{C}\hat{\boldsymbol{\beta}} - \mathbf{d})$  and  $\boldsymbol{\Omega} = \mathbf{G}^{-1/2}\mathbf{C}\mathbf{V}^{CR2}\mathbf{C}'\mathbf{G}^{-1/2}$ . Now suppose that  $\eta \times \boldsymbol{\Omega}$  follows a Wishart distribution with  $\eta$  degrees of freedom and a  $q$ -dimensional identity scale matrix. It then follows that

$$(14) \quad \left( \frac{\eta - q + 1}{\eta q} \right) Q \sim F(q, \eta - q + 1).$$

We will refer to this as the approximate Hotelling's  $T^2$  (AHT) test. We consider how to estimate  $\eta$  below. This approximation is conceptually similar to the Satterthwaite approximation for one-dimensional constraints, and in fact reduces to the Satterthwaite approximation when  $q = 1$ . For  $q > 1$ , the test depends on the multivariate distribution of  $\mathbf{V}^{CR2}$ , including both variance and covariance terms.

Tipton and Pustejovsky (2015) recently introduced this test for application in the special case of CRVE for meta-regression models. Wishart approximations have been considered as approximations in several simpler models where special cases of CRVE are used. Nel and van der Merwe (1986) proposed an AHT-type test for testing equality of multivariate means across two samples with unequal variance-covariance matrices (i.e., the multivariate Behrens-Fisher problem; see also Krishnamoorthy and Yu, 2004). Zhang (2012a) followed a similar approach in developing a test for contrasts in analysis of variance models with unequal within-cell variance, which are particularly simple cases of linear models with heteroskedastic error terms.

Zhang (2012b) extended the method to multivariate analysis of variance models where the covariance of the errors differs across groups, a special case of model (2) where the CR2 variance estimator has a particularly simple form. In each of these special cases, the robust variance estimator is a mixture of Wishart distributions that is well-approximated by a Wishart distribution with estimated degrees of freedom. Additionally, Pan and Wall (2002) described an F-test for use in GEE models, which uses the Wishart approximation to the distribution of  $\mathbf{V}^{CR0}$  but estimates the degrees of freedom using a different method than the one we describe below.

The contribution of the present paper is to extend the AHT test to the general setting of linear models with fixed effects and clustered errors. The remaining question is how to estimate the parameter  $\eta$ , which determines scalar multiplier and denominator degrees of freedom of the AHT test. To do so, we match the mean and variance of  $\mathbf{\Omega}$  to that of the approximating Wishart distribution under the working variance model  $\mathbf{\Phi}$ , just as in the degrees of freedom for the t-test. The problem that arises in doing so is that it is not possible to exactly match both moments if  $q > 1$ . Following Tipton and Pustejovsky (2015), we instead match the mean and total variance of  $\mathbf{\Omega}$ —i.e., the sum of the variances of its entries.

Let  $\mathbf{g}_1, \dots, \mathbf{g}_q$  denote the  $q \times 1$  column vectors of  $\mathbf{G}^{-1/2}$ . Let

$$\mathbf{p}_{si} = (\mathbf{I} - \mathbf{H}_X)_i' \mathbf{A}_i \mathbf{W}_i \ddot{\mathbf{R}}_i \mathbf{M}_{\ddot{\mathbf{R}}} \mathbf{C} \mathbf{g}_s$$

for  $s = 1, \dots, q$  and  $i = 1, \dots, m$ . The degrees of freedom are then estimated under the working model as

$$(15) \quad \eta_M = \frac{q(q+1)}{\sum_{s,t=1}^q \sum_{i,j=1}^m \mathbf{p}_{si}' \mathbf{\Phi} \mathbf{p}_{tj} \mathbf{p}_{ti}' \mathbf{\Phi} \mathbf{p}_{sj} + \mathbf{p}_{si}' \mathbf{\Phi} \mathbf{p}_{sj} \mathbf{p}_{ti}' \mathbf{\Phi} \mathbf{p}_{tj}}.$$

If  $q = 1$ , then  $\eta_M$  reduces to  $\nu_M$  from Equation (13).

This AHT F-test shares several features with the t-test developed by Bell and McCaffrey. As with the t-test, the degrees of freedom of this F-test

Possible to simplify  $p_{si}$  by ignoring within-cluster fixed effects?

depend not only on the number of clusters, but also on features of the covariates being tested. The degrees of freedom can be much lower than  $m - 1$ , particularly when the covariates being tested exhibit high leverage or are unbalanced across clusters. For example, if the goal is to test if there are differences across a three-arm, block-randomized experiment with clustering by block, the degrees of freedom will be largest (approaching  $m - 1$ ) when the treatment is allocated equally across the three groups within each block. When the proportion varies across clusters, the degrees of freedom are reduced, potentially into “small sample” territory even when the number of clusters is large.

A primary difference between the AHT test and the standard test is in the degrees of freedom. We expect that using the AHT degrees of freedom, which take into account features of the covariate distribution, will improve the accuracy of the rejection rates in small samples. We have also claimed that the choice of working model used in the CR2 correction does not have a strong influence on performance. In the next section, we provide evidence for these claims through a careful review of prior simulation study results and through the results of a new simulation study based upon the conditions commonly found in economic applications.

#### 4. SIMULATION EVIDENCE

Evidence from several large simulation studies indicates that hypothesis tests based on the CR2 adjustment and estimated degrees of freedom substantially out-perform the procedures that are most commonly used in empirical applications. However, existing simulations have focused almost entirely on single-parameter tests. In this section, we first review findings from previous simulations, with particular emphasis on the role of covariate features and sample size on the Type I error rates of these tests. We then describe the design and results of a new simulation study, which focused

on the rejection rates of multiple-parameter tests. Throughout, we refer to tests employing the CR2-corrected CRVE and estimated degrees of freedom as the “AHT” test; for t-tests, the estimated degrees of freedom are equivalent to the Satterthwaite approximation given in Equation (13). We refer to tests employing the CR1 correction and  $m - 1$  degrees of freedom as the “standard” test.

#### 4.1. *Review of previous simulation studies*

To date, four simulation studies have examined the performance of the CR2 t-test, across a total of nearly 100 parameter combinations and a range of application contexts. Cameron and Miller (2015) and Imbens and Kolesar (2015) focused on conditions common in economics, while Bell and McCaffrey (2002) focused on those common in complex surveys and Tipton (2015) on those in meta-analysis. Table I summarizes the results of these studies. Some of the studies focused on policy dummies in the balanced case, while others varied the degree of balance; still others examined continuous covariates that are symmetrically distributed, as well as those with high skew and leverage. These studies also examined the role of the number of clusters, with values ranging from 6 to 50, as well as the number of observations per clusters (from 1 to roughly 260). Finally, the studies used a range of both true error structures (including various combinations of heteroskedasticity and clustering) and estimation strategies (including different ‘working’ models), including scenarios in which the working model differed from the true error structure. Finally, while most previous studies focused on OLS estimation, one study (Tipton, 2015) examined the performance of t-tests based on WLS estimation.

Table I also indicates the range of Type I error rates observed across the conditions studied in each of the simulation studies, with values given for both the standard and AHT tests. Across studies, the Type I error for the



Table 1: Type I error rates of  $\chi^2$ -tests based on CPE/N

Article	Table	m	n	Data generation		Type I error ( $\alpha = .05$ )					
				Error structure	Covariate type(s)	Working model	Standard		AHT		
							Min	Max	Min	Max	
BM 2002	2	20	10	RE(3)	c[50%, 15%, M], o[M]	I	.05	.16	.03	.05	
	2	6-50	260	IPUMS-CPS	c[50%]	RE	.07	.11	.05	.05	
IK 2015	1	30	1	H(5)	c[10%]	I	.12	.22	.01	.05	
	2	30	1	H(5); LN	c[10%]	I	.07	.32	.00	.13	
	3	30	1	H(5)	c[50%]	I	.05	.05	.05	.05	
	4	5-10	10-60	RE; RE+H	c[M]	I; RE	.08	.13	.03	.06	
	5	50	6	RE	c[50%]	I; RE	.06	.06	.05	.05	
T 2015	5	50	6	RE; RE+H	c[6%]	I; RE	.15	.23	.01	.05	
	5	50	6	RE	c[K]	I; RE	.13	.13	.03	.03	
	1	10-40	1-10	RE + C(9)	c[50%]	H	.05	.05	.04	.05	
	1	10-40	1-10	RE + C(9)	c[15%]	H	.10	.15	.04	.05	
	1	10-40	1-10	RE + C(9)	o[10%]	H	.04	.19	.04	.06	
T 2015	1	10-40	1-10	RE + C(9)	c[M]	H	.07	.13	.04	.05	
	1	10-40	1-10	RE + C(9)	o[M]	H	.06	.12	.04	.06	
	1	10-40	1-10	RE + C(9)	o[K]	H	.03	.28	.01	.04	
	2	20	1-10	RE + C(9)	c[10%, M], o[10%, M]	H	.01	.12	.01	.06	
	2	20	1-10	RE + C(9)	c[30%, M], o[30%, M]	H	.02	.06	.01	.06	

Table refers to the table within the relevant article.  $m$  is the number of clusters;  $n$  is the number of observations

within each cluster;  $c$  indicates cluster-level covariate, while  $o$  indicates observation-level covariate; % = percent

taking value of one; M = symmetric continuous; K = skewed continuous; H = heteroskedastic; RE = random effects

(Moulton factor); C = correlated errors; LN = log-normal errors; (#) indicates number of different models tested.

standard t-test ranges from .01 to .34 for a stated  $\alpha$  level of .05. These values are particularly far above nominal when the covariate tested is unbalanced, skewed (i.e., high leverage), or when the number of observations per cluster varies. Although not reported in the table, high Type I error rates occur not only when the number of clusters is very small, but also at moderate sizes when the covariate is unbalanced or skewed.

In comparison, the AHT t-test performs considerably better across the range of conditions studied, with Type I error rates ranging between 0.01 and 0.13. Notably, the largest value observed here is from Imbens and Kolesar (2015), who do not break results out by degrees of freedom. Given the condition studied (30 clusters, with only 3 having a policy dummy), it is quite possible that the degrees of freedom are below the cut-off of 4 or 5 at which others have shown the t-test approximation can fail (Tipton, 2015). Putting this value aside, the maximum Type I error observed in these conditions is 0.06, only slightly higher than nominal. Crucially, these nearly nominal Type I error rates hold even when the working model is far from the true error structure, and for various types of covariates. This is because the AHT test takes into account covariate features in the degrees of freedom, which can be far less than  $m - 1$ .

In comparison to the t-test, the AHT F-test has only been studied in a single simulation focused on the meta-analytic case (Tipton and Pustejovsky, 2015). Although this study focused only on the use of CRVE with WLS estimation, it was comprehensive in other regards. In particular, it examined the effects of the number of covariates in the model (up to  $p = 5$ ) and the number of constraints tested ( $q = 2, 3, 4, 5$ ), including cases in which  $p = q$  and in which  $q < p$ . The simulations also examined models with various combinations of covariate types, including both balanced and unbalanced indicator variables, as well as symmetric or skewed continuous covariates. Like Tipton (2015), these simulations focused on true correlation

structures that included heteroskedasticity, clustering (i.e., a cluster specific random effect), and correlated errors. The working models were then chosen to be far from the true error structure (i.e., an independent-errors working model). Finally, the number of clusters was varied from 10 to 100, each with between 1 and 10 observations. Type I error rates of the standard test and the AHT F-test were compared for nominal  $\alpha$  levels of .01, .05, and .10.

The results of the simulations by Tipton and Pustejovsky (2015) indicate that the AHT F-test always has Type I error less than or equal to the stated  $\alpha$  level, except in cases with extreme model misspecification. However, even under such conditions, the Type I error was in line with rates observed for t-tests; for example, for  $\alpha = 0.05$  the error was not above 0.06. In comparison, the Type I error of the standard test was often very high, with maximum rejection rates ranging from .17 to .22, depending on the dimension of the constraint being tested. Like the t-test, the degrees of freedom of the AHT F-test were driven by covariate features, with particularly low degrees of freedom resulting from covariates that are unbalanced or skewed.

While the simulation study by Tipton and Pustejovsky (2015) included a variety of conditions, its design was focused on the types of data found in meta-analytic applications. These differ from the economic context in two ways. First, in meta-analysis, it is common to have heteroskedasticity of a known form and for analysts to incorporate weights in the analysis (typically inverse-variance weights). In comparison, unweighted, OLS estimation is more common in economic applications. Second, meta-analytic regressions often involve testing a variety of types of covariates, including continuous regressors. In comparison, many economic applications are focused on testing binary indicator variables that represent differences between policy regimes. Tests for policy effects can involve cluster-level comparisons (e.g., comparisons across states) or observation-level comparisons (e.g., pre/post comparisons within each state), or a combination of both observation-level

and cluster-level comparisons (as in difference-in-differences analysis). In light of these differences, we conducted a new study to evaluate the performance of the standard and AHT tests under conditions that more closely resemble economic applications.

## 4.2. Simulation Design

The simulation study focused on testing hypotheses about the relative effects of three policy conditions, while varying the manner in which the policy indicators are assigned following one of three distinct designs. First, we considered a randomized block (RB) design in which every policy condition is observed in every cluster. Second, we considered a cluster-randomized (CR) design in which each cluster is observed under a single policy condition. Third, we considered a difference-in-differences (DD) design in which some clusters are observed under all three policy conditions while other clusters are observed under a single condition. For each design, we simulated both balanced and unbalanced configurations, for a total of six distinct study designs, across which the performance of CRVEs is expected to vary. Appendix B describes the exact specification of each design. For each design, we simulated studies with  $m = 15, 30$ , or  $50$  clusters, each with  $n = 18$  or  $30$  units.

For a given study design, we simulated multivariate outcome data so that we could examine the performance of the proposed testing procedures for constraints of varying dimension. Specifically, we simulated a tri-variate, equi-correlated outcome from a data-generating process in which all three policy conditions produce identical average outcomes, so that all tested null hypotheses hold. Let  $y_{hijk}$  denote the measurement of outcome  $k$  at time point  $j$  for unit  $i$  under condition  $h$ , for  $h = 1, \dots, 3$ ,  $i = 1, \dots, m$ ,  $j = 1, \dots, n$ ,

and  $k = 1, \dots, 3$ . The data-generating model is then

$$(16) \quad y_{hijk} = \nu_{hi} + \epsilon_{ijk},$$

where  $\nu_{hi}$  is a random effect for unit  $i$  under condition  $h$  and  $\epsilon_{ijk}$  is the idiosyncratic error for unit  $i$  at time point  $j$  on outcome  $k$ . The random effects for unit  $i$  are taken to have variance  $\text{Var}(\nu_{hi}) = \tau^2$ . We further assumed that the random effects are correlated, which has the effect of inducing variability in the cluster-specific treatment effects and thus a degree of misspecification into the analytic models described below. Letting  $\sigma_\delta^2$  denote the degree of treatment effect variability relative to the total variability in a given outcome measurement, we simulated the random effects  $\nu_{1i}, \nu_{2i}, \nu_{3i}$  to satisfy  $\text{Var}(\nu_{gi} - \nu_{hi}) = \sigma_\delta^2$  for  $g \neq h, g, h = 1, 2, 3$ . The errors at a given time point are assumed to be correlated, with  $\text{Var}(\epsilon_{ijk}) = 1 - \tau^2$  and  $\text{corr}(\epsilon_{ijk}, \epsilon_{ijl}) = \rho$  for  $k \neq l, k, l = 1, 2, 3$ .

Under this data-generating process, we simulated data based on parameter values of  $\tau^2 = .05, .15$ , or  $.25$  for the intra-class correlation; outcomes that were either weakly ( $\rho = .2$ ) or strongly correlated ( $\rho = .8$ ); and values of  $\sigma_\delta^2 = .00, .01$ , or  $.04$  for treatment effect variability. Each combination of sample sizes and parameter levels was simulated under each of the six study designs, yielding a total of 648 simulation conditions.

Given a set of simulated data, we estimated the effects of the second and third policy conditions (relative to the first) on each outcome, using a seemingly unrelated regression framework. The general analytic model for the difference-in-differences design was

$$(17) \quad y_{hijk} = \mu_{hk} + \alpha_i + \gamma_j + \epsilon_{ijk},$$

where  $\mu_{hk}$  is the mean of outcome  $k$  under condition  $h$ ,  $\alpha_i$  is a fixed effect for each cluster,  $\gamma_j$  is a fixed effect for each unit within the cluster (i.e., per time-point), and  $\epsilon_{ijk}$  is residual error. For the cluster-randomized designs,

fixed effects for clusters were omitted because the clusters are nested within treatment conditions. For the randomized block designs, the fixed effects for time-points were omitted for simplicity. The model is estimated by OLS after absorbing any fixed effects, and so the “working” model amounts to assuming that the residuals are all independent and identically distributed. Note that the working model departs from the true data generating model both because of correlation among the outcomes ( $\rho > 0$ ) and because of treatment effect variability ( $\sigma_\delta^2 > 0$ ). The range of parameter combinations used in the true data generating model thus allow us to examine the performance of the AHT test under both small and large degrees of working model misspecification.

Analytic model (17) provided opportunities to test a range of single- and multi-parameter constraints. We first tested the single-dimensional null hypotheses that a given policy condition had no average effect on the first outcome ( $H_0 : \mu_{11} = \mu_{12}$  or  $H_0 : \mu_{11} = \mu_{13}$ ). We also tested the null hypothesis of no differences among policy conditions on the first outcome ( $H_0 : \mu_{11} = \mu_{12} = \mu_{13}$ ), which has dimension  $q = 2$ . We then tested the multi-variate versions of the above tests, which involve all three outcome measures jointly. Namely, we tested the null hypotheses that a given policy condition had no average effects on any outcome (i.e.,  $H_0 : \mu_{11} = \mu_{1h}, \mu_{21} = \mu_{2h}, \mu_{31} = \mu_{3h}$ , for  $h = 2$  or  $h = 3$ ) which has dimension  $q = 3$ , and the null hypothesis of no differences among policy conditions on any outcome ( $H_0 : \mu_{11} = \mu_{12} = \mu_{13}, \mu_{21} = \mu_{22} = \mu_{23}, \mu_{31} = \mu_{32} = \mu_{33}$ ), which has dimension  $q = 6$ . For a given combination of sample sizes, parameter levels, and study design, we simulated 10,000 datasets from model (16), estimated model (17) on each dataset, and tested all of the hypotheses described above. Simulated Type I error rates therefore have standard errors of approximately 0.001 for  $\alpha = .01$ , 0.0022 for  $\alpha = .05$ , and 0.003 for  $\alpha = .10$ .

### 4.3. *Simulation Results*

Our discussion of the simulation results is focused on four trends, each of which is depicted visually in a figure and described in the text. All of the trends are similar to the findings from Tipton and Pustejovsky (2015), which provides further support that the AHT F-test performs well across a wide range of data generating mechanisms and parameter combinations.

The first finding is that the AHT test has Type I error close to the stated  $\alpha$  level for all parameter combinations studied, whereas the standard test (based on the CR1 variance estimator and  $m - 1$  degrees of freedom) does not. Figure 1 illustrates this pattern, for constraints of varying dimension (from  $q = 1$ , in the first column, to  $q = 6$ , in the final column) and nominal  $\alpha$  level (from .01, in the first row, to .10, in the last row). In each of these figures, the number of clusters varies from 15 to 50 (on the horizontal axis), the solid horizontal line indicates the stated  $\alpha$  level and the dashed line indicates an upper confidence bound on simulation error. It can be seen that the AHT test has Type I error near the stated  $\alpha$  level, even with a small number of clusters. When the number of clusters is very small, the Type I error can be smaller than the stated  $\alpha$  level. Although there exist situations in which the error is above the simulation bound, the departures are typically small. For example, when  $m = 15$  the rejection rates do not exceed 0.021 for  $\alpha = .01$ , 0.073 for  $\alpha = .05$ , and 0.134 for  $\alpha = .10$ . The rejection rates are even closer to nominal for lower-dimensional constraints. In comparison, the Type I error for the standard test can be markedly higher than the stated  $\alpha$  level, particularly when the number of clusters is small or the dimension of the hypothesis is large. For example, for nominal  $\alpha = .05$ , the maximum Type I error ranges from 0.124 ( $q = 1$ ) to 0.686 ( $q = 6$ ) for data sets with 15 clusters. Perhaps even more important for practice, even when there are 50 clusters, the rejection rate of the standard test can be far above the stated  $\alpha$  level. Again, focusing on the  $\alpha = 0.05$  case, the

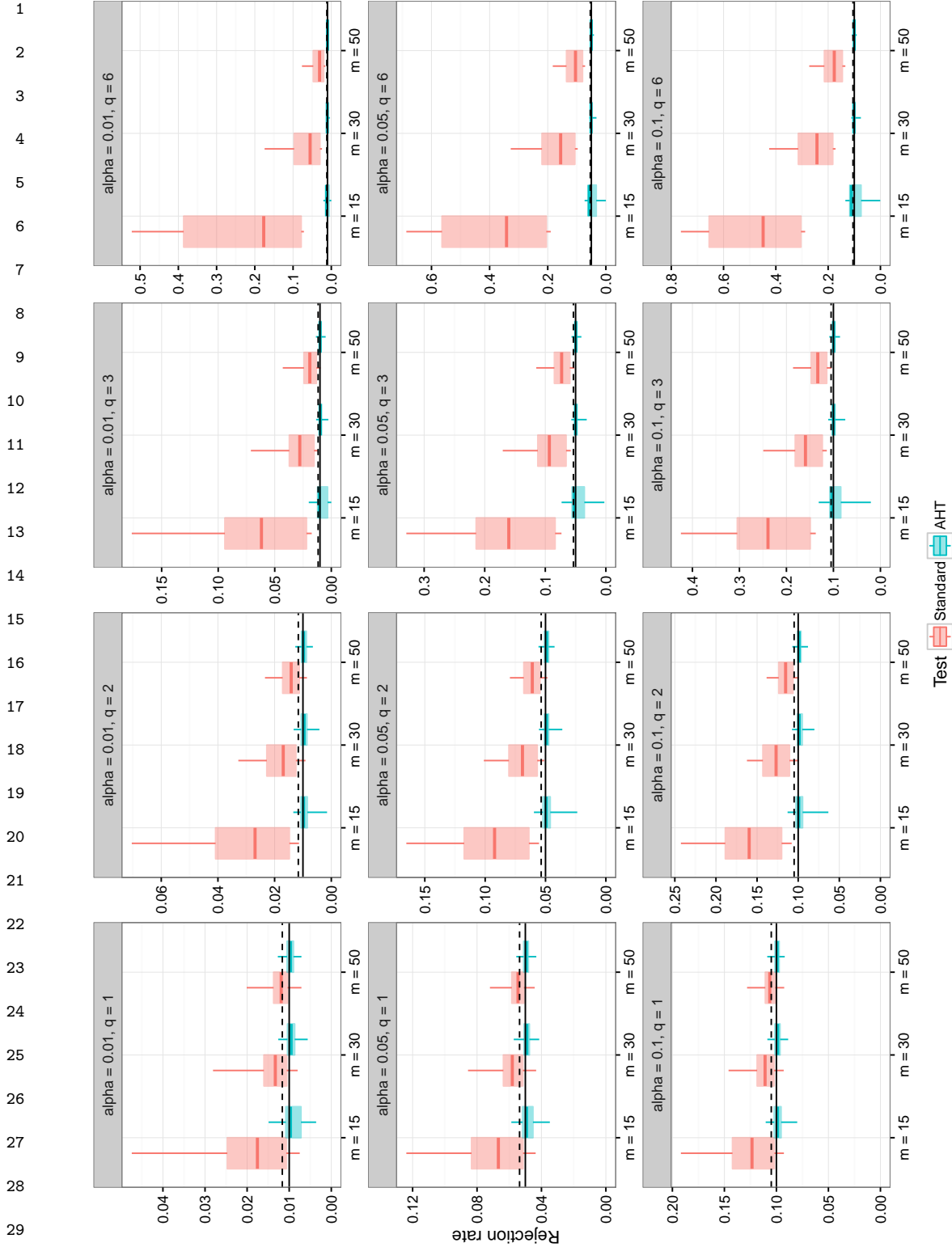


Figure 1: Rejection rates of Standard and AHT tests, by dimension of hypothesis ( $q$ ) and nominal type I error ( $\alpha$ ).



maximum error ranges from 0.072 ( $q = 1$ ) to 0.183 ( $q = 6$ ).

In order to better understand the effects of different parameter combinations on the performance of both tests, Figure 2 focuses on the  $\alpha = 0.05$  case and breaks out the results by study design. The top row of the figure depicts the standard test, while the bottom row depicts the AHT test; columns correspond to the number of clusters. Within each graph, results are given by study design (on the horizontal axis), with colors corresponding to the dimension of the test. In the top panel, it can be seen that the rejection rate of the standard test increases with the dimension of the test ( $q$ ) and the degree of unbalance in the study design. Differences between the balanced and unbalanced designs are largest for the CR and DD designs, with smaller discrepancies in RB designs. The bottom row of Figure 2 displays results for the AHT test; here we focus on three trends. First, the rejection rate of the AHT test usually increases as the dimension of the test increases, though never above 0.073. Second, unbalanced designs led to rejection rates that were usually below the nominal  $\alpha$ —just the opposite of how the standard test is affected by unbalance. This trend is the strongest for CR and DD designs, where Type I error can be close to 0 at its minimum. Third, for studies with at least 30 clusters, rejection rates are very close to nominal (between 0.032 and 0.057) across all conditions studied.

Next, by simulating the errors across a variety of parameter combinations, we were also able to test the impact of misspecification of the working model on Type I error. Because the CR2 correction and AHT degrees of freedom are both based on a working model with independent, homoskedastic errors, model misspecification increases with the true level of treatment effect variance ( $\sigma_\delta^2$ ) and intra-class correlation ( $\tau^2$ ). Figure 3 depicts Type I error rates for  $\alpha = 0.05$  for the AHT test, with separate graphs according to the dimension  $q$  of the test (columns) and the number of clusters (rows). Within each panel, results are separated by the 9 combinations of  $\sigma_\delta^2$  and

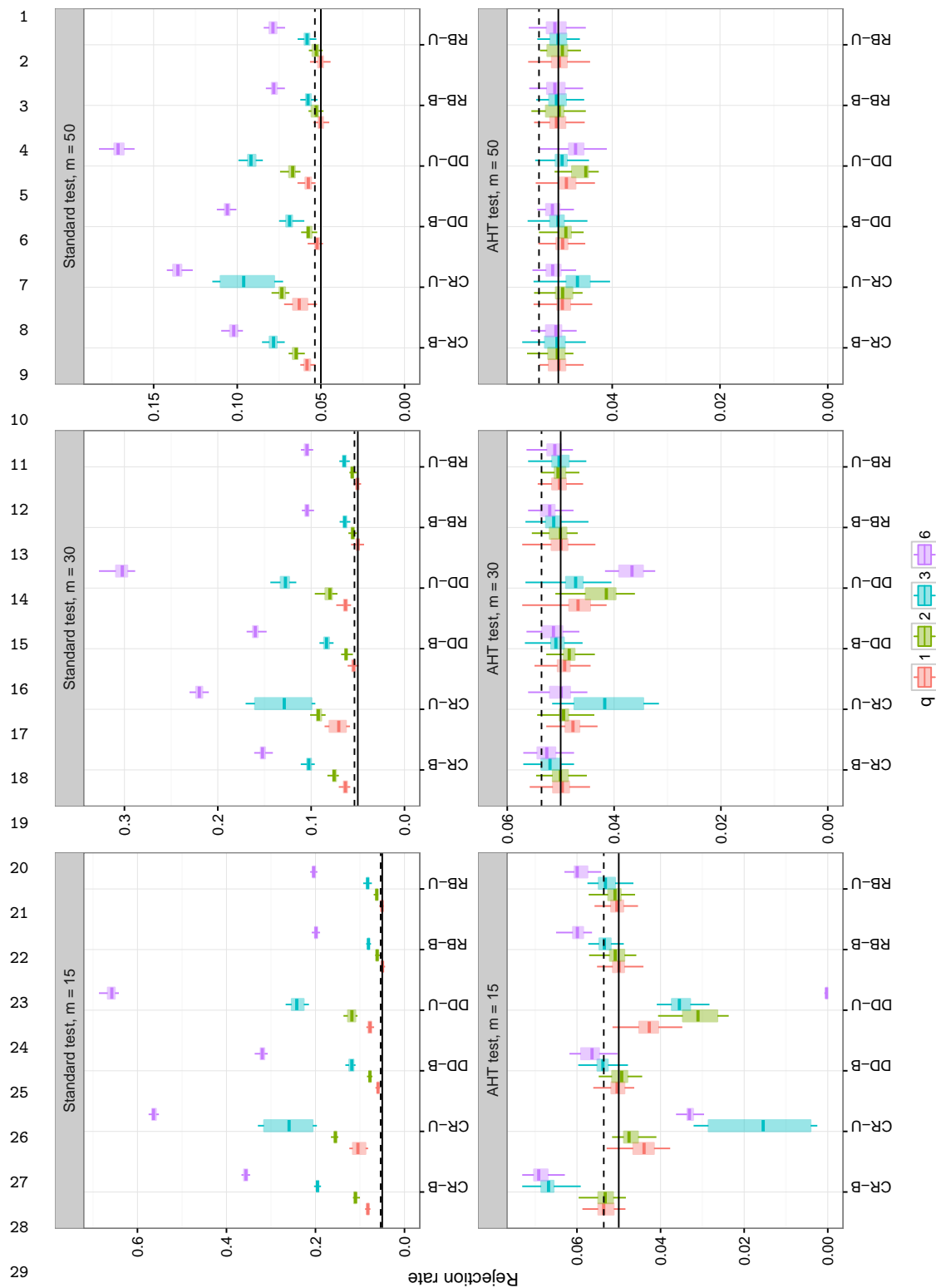
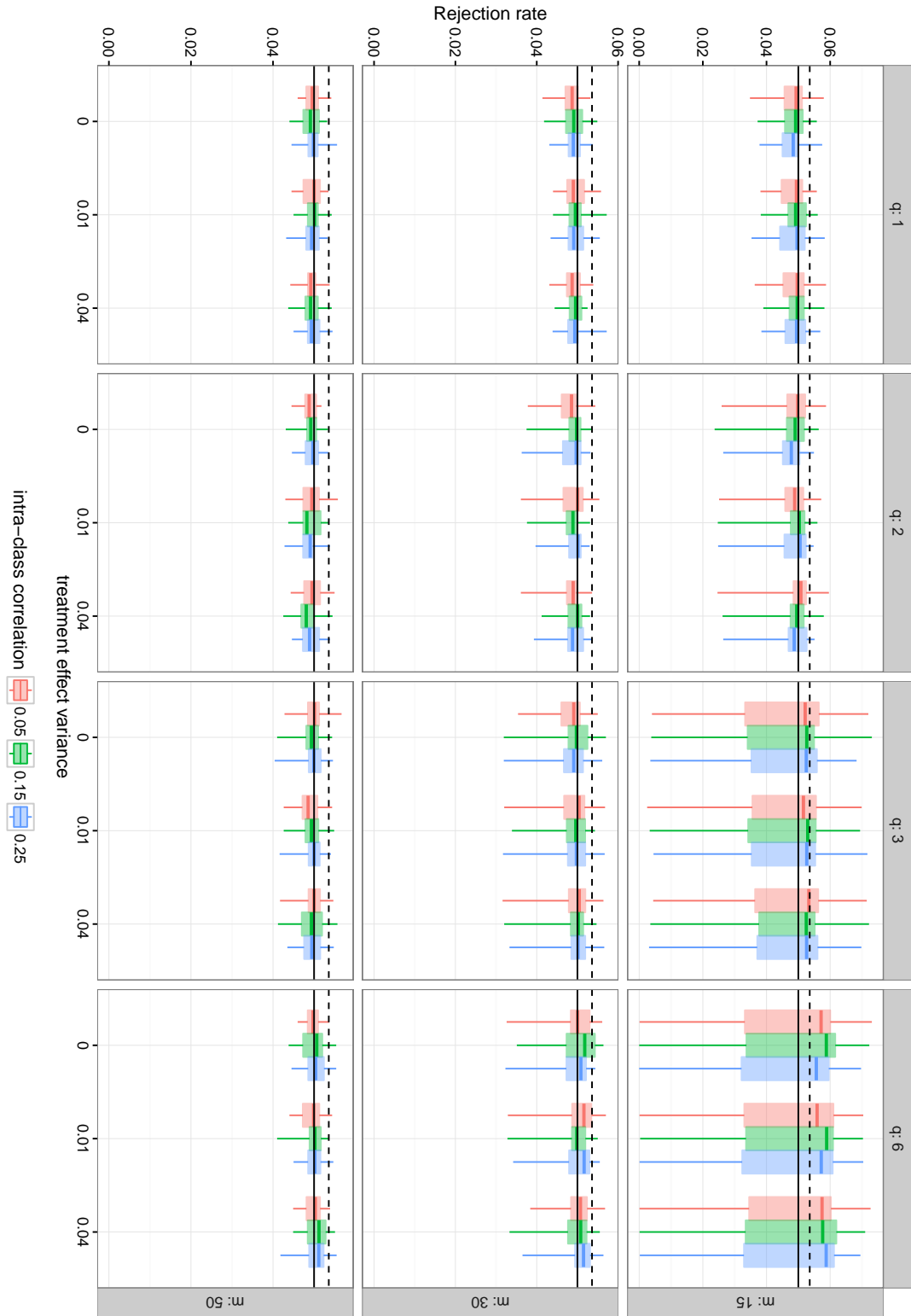


Figure 2: Rejection rates of Standard and AHT tests, by study design and dimension of hypothesis ( $q$ ). CR = cluster-randomized design; DD = difference-in-differences design; RB = randomized block design; B = balanced; U = unbalanced.

Figure 3: Rejection rates of AHT test, by treatment effect variance and intra-class correlation.



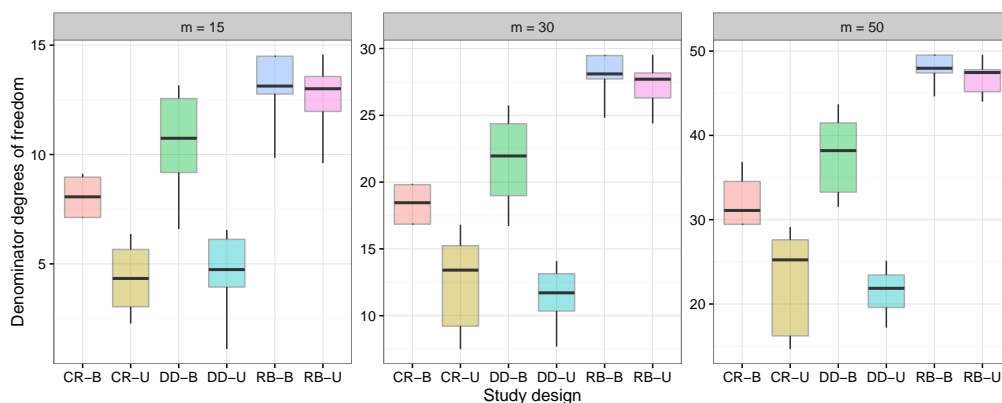


Figure 4: Range of denominator degrees of freedom for AHT test, by number of clusters and study design. CR = cluster-randomized design; DD = difference-in-differences design; RB = randomized block design; B = balanced; U = unbalanced.

$\tau^2$ . It can be seen that the range of rejection rates remains very similar across the 9 error structures, with no clear pattern to the small differences that emerge. These results follow closely those from Tipton and Pustejovsky (2015), which also found that even with extreme model misspecification the Type I error of the CR2S test was close to nominal.

Finally, Figure 4 depicts the range of estimated degrees of freedom for the AHT test as a function of the simulated study design and number of clusters ( $m$ ). Within each box plot, the degrees of freedom vary depending on the hypothesis tested, with constraints of larger dimension having lower degrees of freedom. It can be seen that the AHT degrees of freedom are often far less than  $m - 1$  and that they are strongly affected by the pattern of treatment assignment and the degree of balance. The balanced and unbalanced RB designs generally had AHT degrees of freedom closest to  $m - 1$  because the treatment effects being tested are all identified within each cluster. The balanced DD design usually had the next largest degrees of freedom because it involved contrasts between two patterns of treatment configuration, fol-

lowed by the balanced CR design, which involved contrasts between three patterns of treatment configurations. For both of these designs, unbalance led to sharply reduced degrees of freedom.

These new simulation results have demonstrated that the standard robust Wald test, using the CR1 correction and  $m - 1$  degrees of freedom, produces a wide range of rejection rates, often far in excess of the nominal Type I error. In contrast, the rejection rates of the AHT test are below or at most slightly above nominal, across the conditions that we have examined. This is because the AHT test incorporates information about the covariate features into its estimated degrees of freedom, whereas the standard test does not. An important question that remains then is how much the AHT and standard tests diverge in actual application. In the next section, we compare the two tests in several examples, drawn from a range of recent empirical research.

## 5. EXAMPLES

This section presents three short examples that illustrate the performance of CRVE across a variety of applied contexts. In the first example, the effects of substantive interest involve between-cluster contrasts. The second example involves a cluster-robust Hausman test for differences between within- and across-cluster information. In the final example, the effects are identified within each cluster. In each example, we demonstrate the proposed AHT test for single- and multiple-parameter hypotheses and compare the results to the standard test based on the CR1 variance estimator and  $m - 1$  degrees of freedom. The focus here is on providing insight into the conditions under which the AHT and standard estimators diverge in terms of three quantities of interest: the standard error estimates, the degrees of freedom estimates, and the stated p-values. Data files and replication code (in R) are available for each analysis as an online supplement.

### 5.1. *Achievement Awards demonstration*

Angrist and Lavy (2009) reported results from a randomized trial in Israel that aimed to increase completion rates of the Bagrut, the national matriculation certificate for post-secondary education, among low-achieving high school students. In the Achievement Awards demonstration, 40 non-vocational high schools with low rates of Bagrut completion were selected from across Israel, including 10 Arab and 10 Jewish religious schools and 20 Jewish secular schools. The schools were then pair-matched based on 1999 rates of Bagrut completion, and within each pair one school was randomized to receive a cash-transfer program. In these treatment schools, seniors who completed certification were eligible for payments of approximately \$1,500. Student-level covariate and outcome data were drawn from administrative records for the school years ending in June of 2000, 2001, and 2002. The incentive program was in effect for the group of seniors in treatment schools taking the Bagrut exams in Spring of 2001, but the program was discontinued for the following year. We therefore treat completion rates for 2000 and 2002 as being unaffected by treatment assignment. The primary outcome of interest is Bagrut completion.

This study provides an opportunity to examine the AHT test in a situation in which the treatment was assigned at the cluster level, with a smaller number of clusters. For simplicity, we restrict our analysis to the sample of female students, which reduces the total sample to 35 schools. Following the original analysis of Angrist and Lavy (2009), we allow the program's effects to vary depending on whether a students was in the upper or lower half of the distribution of prior-year academic performance. Letting  $h = 1, 2, 3$  index the sector of each school (Arab religious, Jewish religious, or Jewish secular), we consider the following analytic model:

$$(18) \quad y_{hitj} = z_{hit} \mathbf{r}'_{hitj} \boldsymbol{\beta}_h + \mathbf{s}'_{hitj} \boldsymbol{\gamma} + \gamma_{ht} + \mu_{hi} + \epsilon_{hitj}$$

In this model for student  $j$  in year  $t$  in school  $i$  in sector  $h$ ,  $z_{hit}$  is an indicator equal to one in the treatment schools for the 2001 school year and otherwise equal to zero;  $\mathbf{r}_{hitj}$  is a vector of indicators for whether the student is in the lower or upper half of the distribution of prior academic performance; and  $\beta_h = (\beta_{1h}, \beta_{2h})$  is a vector of average treatment effects for schools in sector  $h$ . The vector  $\mathbf{s}_{hitj}$  includes the following individual student demographic measures: mother's and father's education, immigration status, number of siblings, and indicators for each quartile in the distribution of prior-year academic performance. The model also includes fixed effects  $\gamma_{ht}$  for each sector in each year and  $\mu_{hi}$  for each school.

Based on Model (18), we test four hypotheses, again with the goal of exploring the use of the AHT tests under a range of conditions. First, we assume that the program effects are constant across sector (i.e.,  $\beta_1 = \beta_2 = \beta_3 = \beta$ ) and test for whether the program affected completion rates for students in the upper half of the prior achievement distribution ( $H_0 : \beta_2 = 0$ , with  $q = 1$ ). Second, we test for whether the program was effective in either half of the prior academic performance ( $H_0 : \beta = 0$ , with  $q = 2$ ), still assuming that program effects are constant across sector. Third, we test for whether program effects in the upper half of the prior achievement distribution are moderated by school sector ( $H_0 : \beta_{21} = \beta_{22} = \beta_{23}$ , with  $q = 3$ ). Finally, we conduct a joint test for whether program effects in either half of the prior achievement distribution are moderated by school sector ( $H_0 : \beta_1 = \beta_2 = \beta_3$ , with  $q = 4$ ).

Table II reports the results of all four hypothesis tests. These results indicate three important trends. First, in the case of the first two hypotheses, the AHT test statistics are only slightly smaller than their standard counterparts, but the degrees of freedom are considerably smaller. These differences in degrees of freedom arise because the treatment was assigned at the cluster level, while the subgroups varied within each cluster. Second, the third

TABLE II  
TESTS OF TREATMENT EFFECTS IN THE ACHIEVEMENT AWARDS DEMONSTRATION

Hypothesis	Test	F	df	p
ATE - upper half ( $q = 1$ )	Standard	5.746	34.00	0.02217
	AHT	5.169	15.86	0.03726
ATE - joint ( $q = 2$ )	Standard	3.848	34.00	0.03116
	AHT	3.371	15.46	0.06096
Moderation - upper half ( $q = 2$ )	Standard	3.186	34.00	0.05393
	AHT	0.091	3.19	0.91520
Moderation - joint ( $q = 4$ )	Standard	8.213	34.00	0.00010
	AHT	2.895	3.21	0.19446

and fourth hypotheses tests, which compared treatment effects across sectors and subgroups, are cases in which the AHT and standard tests diverge markedly. For these cases, the AHT test statistic and degrees of freedom are both considerably smaller than those from the standard test. This reflects the degree of unbalance in allocations across sectors (19 Jewish secular, 7 Jewish religious, and 9 Arab religious schools), combined with cluster-level randomization. In combination, these smaller test statistics and degrees of freedom result in larger p-values for the AHT test when compared to the standard test.

## 5.2. *Effects of minimum legal drinking age on mortality*

Our second example focuses on panel data, using an example described in Angrist and Pischke (2014, see also Carpenter and Dobkin, 2011). Based on data from the Fatal Accident Reporting System maintained by the National Highway Traffic Safety Administration, we estimated the effects of changes in the minimum legal drinking age over the time period of 1970-1983 on state-level death rates resulting from motor vehicle crashes. A standard



difference-in-differences specification for such a state-by-year panel is

$$(19) \quad y_{it} = \mathbf{r}_{it}'\boldsymbol{\beta} + \gamma_t + \mu_i + \epsilon_{it}.$$

In this model, time-point  $t$  is nested within state  $i$ ; the outcome  $y_{it}$  is the number of deaths in motor vehicle crashes (per 100,000 residents) in state  $i$  at time  $t$ ;  $\mathbf{r}_{it}$  is a vector of covariates;  $\gamma_t$  is a fixed effect for time point  $t$ ; and  $\mu_i$  is an effect for state  $i$ . The vector  $\mathbf{r}_{it}$  consists of a measure of the proportion of the population between the ages of 18 and 20 years who can legally drink alcohol and a measure of the beer taxation rate, both of which vary across states and across time.

We apply both random effects (RE) and fixed effects (FE) approaches to estimate the effect of lowering the legal drinking age. For the RE estimates, we use WLS with weights derived under the assumption that  $\mu_1, \dots, \mu_m$  are mutually independent, normally distributed, and independent of  $\epsilon_{it}$  and  $\mathbf{r}_{it}$ . We also report an artificial Hausman test (Arellano, 1993; Wooldridge, 2002) for correlation between the covariates  $\mathbf{r}_{it}$  and the state effects  $\mu_i$ . Such correlation creates bias in the RE estimator of the policy effect, thus necessitating the use of the FE estimator. The artificial Hausman test amends model (19) to include within-cluster deviations for the variables of interest, so that the estimating equation is

$$(20) \quad y_{it} = \mathbf{r}_{it}\boldsymbol{\beta} + \ddot{\mathbf{r}}_{it}\boldsymbol{\delta} + \gamma_t + \mu_i + \epsilon_{it},$$

where  $\ddot{\mathbf{r}}_{it}$  denotes the within-cluster deviations of the covariate. The parameter  $\boldsymbol{\delta}$  captures the difference between the between-cluster and within-cluster estimates of  $\boldsymbol{\beta}$ . With this setup, the artificial Hausman test amounts to testing the null hypothesis that  $\boldsymbol{\delta} = \mathbf{0}$ , where  $\boldsymbol{\delta}$  is estimated using RE.

Table III displays the results of the tests for the policy variable and the Hausman tests for each model specification. The results of the policy effect tests are quite similar across specifications and versions of the test. Of note

Possible to get replication data from Carpenter and Dobkin (2011) instead?

TABLE III  
TESTS OF EFFECTS OF MINIMUM LEGAL DRINK AGE AND HAUSMAN SPECIFICATION

TEST				
Hypothesis	Test	F	df	p
Random effects	Standard	8.261	49.00	0.00598
	AHT	7.785	24.74	0.00999
Fixed effects	Standard	9.660	49.00	0.00313
	AHT	9.116	22.72	0.00616
Hausman test	Standard	2.930	49.00	0.06283
	AHT	2.489	8.69	0.13980

is that, for both the RE and FE estimates, the AHT tests have only half the degrees of freedom of the corresponding standard tests. For the artificial Hausman test, the AHT test has fewer than 9 degrees of freedom, which leads to a much larger p-value compared to using the standard test based on CR1.

### 5.3. *Tennessee STAR class-size experiment.*

The final example demonstrates an application in which the AHT and standard tests lead to similar results. The Tennessee STAR class size experiment is one of the most intensively studied interventions in education (for a detailed review, see Schanzenbach, 2006). The experiment involved students in kindergarten through third grade across 79 schools. Within each school, students and their teachers were randomized equally to one of three conditions: small class-size (targeted to have 13-17 students), regular class-size, or regular class-size with an aide. Subsequent research has focused on the effects of these conditions on kindergarten reading, math, and word recognition (Achilles, Bain, Bellott, Boyd-Zaharias, Finn, Folger, Johnston and Word, 2008); high school test scores (Schanzenbach, 2006); college entrance exam participation (Krueger and Whitmore, 2001); and home ownership

and earnings (Chetty, Friedman, Hilger, Saez, Schanzenbach and Yagan, 2011), among other outcomes.

The STAR experiment involved three treatment conditions and multiple outcomes, providing a scenario where both t-tests (with  $q = 1$ ) and F-tests with varying constraint dimensions can be applied. For simplicity, we focus only on the subgroup of students who were in kindergarten during the first year of the study, and on three outcomes measured at the end of the kindergarten year: reading, word recognition, and math (Achilles et al., 2008). Outcome scores are standardized to percentile ranks, following Krueger and Whitmore (2001). The analytic model is:

$$(21) \quad y_{ijk} = \mathbf{r}_{ij}'\boldsymbol{\beta}_k + \mathbf{s}_{ij}'\boldsymbol{\gamma}_0 + \gamma_k + \mu_i + \epsilon_{ijk},$$

where  $y_{ijk}$  is the percentile rank on outcome  $k$  for student  $j$  in school  $i$ ;  $\mathbf{r}_{ij}$  includes indicators for the small-class and regular-plus-aide conditions;  $\mathbf{s}_{ij}$  includes student demographic covariates (i.e., free or reduced-price lunch status; race; gender; age);  $\gamma_k$  is a fixed effect for outcome  $k$ ; and  $\mu_i$  is a fixed effect for school  $i$ . In this model,  $\beta_{1k}$  represents the average effect of being in a small class and  $\beta_{2k}$  represents the average of effect of being in a regular class with an aid, in each case compared to a regular-size class without an aid.

Using this model, we test four distinct hypotheses that vary in dimension from  $q = 1$  to  $q = 6$ . First, using only the math achievement scores, we test the effects of small class size ( $H_0 : \beta_{11} = 0$ ) while maintaining the assumption that the additional classroom aide has no effect on student achievement (i.e., constraining  $\beta_{21} = 0$ ). Second, again only using the data for outcome  $k$ , we test the hypothesis that there are no differences across the three class-size conditions (i.e.,  $H_0 : \boldsymbol{\beta}_1 = \mathbf{0}$ ). Third, combining the data across all three outcomes, we test the hypothesis that small class size (vs regular and regular plus aide) had no effects on any outcome (i.e.,  $\beta_{11} =$

$\beta_{12} = \beta_{13} = 0$ ). Finally, we test the hypothesis that there are no differences across the three class-size conditions on any outcome (i.e.,  $H_0 : \beta_1 = \beta_2 = \beta_3 = \mathbf{0}$ ). The third and fourth tests use the seemingly unrelated regression (SUR) framework, in which separate treatment effects are estimated for each outcome, but the student demographic effects and school fixed effects are pooled across outcomes. In all models, we estimated  $\beta_k$  and  $\gamma$  after absorbing the school fixed effects and clustering the standard errors by school.

TABLE IV  
TESTS OF TREATMENT EFFECTS IN THE TENNESSEE STAR CLASS SIZE EXPERIMENT

Outcome	Effect	Test	F	df	p
Math	Small class (q=1)	Standard	13.624	78.0	0.00041
		AHT	13.590	69.0	0.00045
	Small class and classroom aide (q=2)	Standard	6.838	78.0	0.00183
		AHT	6.725	68.6	0.00215
Combined	Small class (q=3)	Standard	6.408	78.0	0.00062
		AHT	6.206	67.0	0.00088
	Small class and classroom aide (q=6)	Standard	3.284	78.0	0.00622
		AHT	3.042	64.9	0.01103

Table IV displays the results for a representative subset of these hypothesis tests, using either the standard test (with CR1) or the AHT test (with CR2). These results illustrate two important points regarding the use of the AHT test in practice. First, across all three analyses, the AHT t- and F-tests are only typically slightly smaller than the corresponding standard test. Second, if treatment is randomly allocated in approximately equal proportions within each cluster—as occurred in the TN STAR experiment—the degrees of freedom for the AHT tests are only slightly smaller than those for the standard tests. In combination with the rather large sample size of 79 schools, these differences have only a minimal effect on the p-values for

these tests. As the previous two examples illustrate, however, the similarity between these tests is not common, and is a result only of the design of the study, indicating that the standard test is best used only in experiments randomized within clusters.

## 6. CONCLUSION

Across the field of economics, empirical studies often involve modeling data with a correlated error structure. Correlated errors arise in the analysis of multi-stage samples, cluster-randomized trials, panel data, and regression discontinuities with discrete forcing variables, among other study designs. It is now routine to handle dependent error structures by using cluster-robust variance estimation, which provides asymptotically valid standard errors and hypothesis tests without making strong parametric assumptions about the error structure. However, a growing body of recent work has drawn attention to the shortcomings of CRVE methods when the data include only a small or moderate number of independent clusters (Cameron et al., 2008; Cameron and Miller, 2015; Imbens and Kolesar, 2015; Webb and MacKinnon, 2013). In particular, Wald tests based on CRVE can have rejection rates far in excess of the nominal Type I error. This problem is compounded by the fact that the performance of standard Wald tests depends on features of the study design beyond just the total number of clusters, which can make it difficult to determine whether standard, asymptotic valid CRVE methods are accurate.

One promising solution to this problem is to use the bias-reduced linearization variance estimator (i.e., CR2) proposed by Bell and McCaffrey (2002), which corrects the CRVE so that it is exactly unbiased under an analyst-specified working model for the error structure, together with degrees of freedom estimated based on the same working model. In this paper, we have demonstrated that the CR2 variance estimator is a fully general

solution, which can be applied even in models with fixed effects in multiple dimensions. Our re-formulation of the bias-reduced linearization criteria also makes clear how to calculate the CR2 correction when the model includes fixed effects, whether those fixed effects are estimated by OLS or are instead absorbed before estimating the target regression parameters. Finally, we have proposed a method for testing hypotheses that involve multiple constraints on regression parameters, based on an approximation that generalizes the existing Satterthwaite approximation for t-tests. With the modifications and extensions proposed in this paper, the CR2 variance estimator and small-sample testing procedures can be applied in a wide range of analytic models—essentially, any model estimated by ordinary or weighted least squares.

We join Imbens and Kolesar (2015) in arguing that the CR2 estimator and corresponding estimated degrees of freedom for hypothesis tests should be applied routinely, whenever analysts use CRVE and hypothesis tests based thereon. Because the performance of standard CRVE methods depends on features of the study design, the total number of clusters in the data is an insufficient guide to whether small-sample corrections are needed. Instead, the clearest way to determine whether small-sample corrections are needed is simply to calculate them. The proposed AHT test involves two adjustments: use of the CR2 adjustment for the variance estimator and use of estimated degrees of freedom. Our simulation study illustrates that the combined result of these adjustments results in an AHT test with Type I error close to the stated  $\alpha$  level. Furthermore, our empirical examples illustrate that the degrees of freedom adjustment has a relatively larger influence on small-sample performance. These degrees of freedom can be much smaller than the number of clusters, particularly when the covariates involved in the test involve high leverage or are unbalanced across clusters. The estimated degrees of freedom are indicative of the precision of the standard errors, and

thus provide diagnostic information that is similar to the effective sample size measure proposed by Carter et al. (2013). We therefore recommend that the degrees of freedom be reported along with standard errors and  $p$ -values whenever the method is applied.

The idea of developing small-sample adjustments based on a working model may seem strange to analysts accustomed to using CRVE—after all, the whole point of clustering standard errors is to avoid making assumptions about the error structure. However, simulation studies reported here and elsewhere (Tipton, 2015; Tipton and Pustejovsky, 2015) have demonstrated that the approach is actually robust to a high degree of misspecification in the working model. Furthermore, while the working model provides necessary “scaffolding” when the number of clusters is small, its influence tends to fall away as the number of clusters increases, so that the CR2 estimator and AHT maintain the same asymptotic robustness as standard CRVE methods.

One outstanding problem with the CR2 variance estimator is that it can become computationally costly (or even infeasible) when the within-cluster sample sizes are large (Mackinnon, 2014). For example, Bertrand et al. (2004) analyzed micro-level data from a 21-year panel of current population survey data, with clustering by state. Their data included some state-level clusters with over  $n_i = 10,000$  individual observations. The CR2 adjustment matrices have dimension  $n_i \times n_i$ , and would be very expensive to compute in this application. Methods for improving the computational efficiency of the CR2 variance estimator (or alternative estimators that have similar performance to CR2), should be investigated further.

This paper has developed the CR2 estimator and AHT testing procedure for weighted least squares estimation of linear regression models. Extensions to linear regression models with clustering in multiple, non-nested dimensions (cf. Cameron, Gelbach and Miller, 2011) appear to be possi-

ble, and their utility should be further investigated. McCaffrey and Bell (2006) have proposed extensions to bias-reduced linearization for use with generalized estimating equations, and future work should consider further extensions to other classes of estimators, including two-stage least squares and generalized method of moments. McCaffrey and Bell (2006) also found that for single-parameter hypotheses, a saddlepoint approximation to the Wald test statistic provides even more accurate rejection rates than the Satterthwaite approximation given in Equation (13). It would be interesting to investigate whether the saddlepoint approximation could be extended to handle multiple-parameter constraints, although this appears to be far from straight-forward.

## APPENDIX A: BRL ADJUSTMENT MATRICES

This appendix provides proof of the two theorems from Section 2.

### A.1. Proof of Theorem 1

The Moore-Penrose inverse of  $\mathbf{B}_i$  can be computed from its eigen-decomposition. Let  $b \leq n_i$  denote the rank of  $\mathbf{B}_i$ . Let  $\mathbf{\Lambda}$  be the  $b \times b$  diagonal matrix of the positive eigenvalues of  $\mathbf{B}_i$  and  $\mathbf{V}$  be the  $n_i \times b$  matrix of corresponding eigenvectors, so that  $\mathbf{B}_i = \mathbf{V}\mathbf{\Lambda}\mathbf{V}'$ . Then  $\mathbf{B}_i^+ = \mathbf{V}\mathbf{\Lambda}^{-1}\mathbf{V}'$  and  $\mathbf{B}_i^{+1/2} = \mathbf{V}\mathbf{\Lambda}^{-1/2}\mathbf{V}'$ . Now, observe that

$$\begin{aligned} \ddot{\mathbf{R}}_i' \mathbf{W}_i \mathbf{A}_i (\mathbf{I} - \mathbf{H}_{\mathbf{X}})_i \mathbf{\Phi} (\mathbf{I} - \mathbf{H}_{\mathbf{X}})_i' \mathbf{A}_i' \mathbf{W}_i \ddot{\mathbf{R}}_i &= \ddot{\mathbf{R}}_i' \mathbf{W}_i \mathbf{D}_i \mathbf{B}_i^{+1/2} \mathbf{B}_i \mathbf{B}_i^{+1/2} \mathbf{D}_i' \mathbf{W}_i \ddot{\mathbf{R}}_i \\ (22) \qquad \qquad \qquad &= \ddot{\mathbf{R}}_i' \mathbf{W}_i \mathbf{D}_i \mathbf{V} \mathbf{V}' \mathbf{D}_i' \mathbf{W}_i \ddot{\mathbf{R}}_i. \end{aligned}$$

Because  $\mathbf{D}_i$ , and  $\mathbf{\Phi}$  are positive definite and  $\mathbf{B}_i$  is symmetric, the eigenvectors  $\mathbf{V}$  define an orthonormal basis for the column span of  $(\mathbf{I} - \mathbf{H}_{\mathbf{X}})_i$ . We now show that  $\ddot{\mathbf{U}}_i$  is in the column space of  $(\mathbf{I} - \mathbf{H}_{\mathbf{X}})_i$ . Let  $\mathbf{Z}_i$  be an  $n_i \times (r + s)$  matrix of zeros. Let  $\mathbf{Z}_k = -\ddot{\mathbf{U}}_k \mathbf{L}^{-1} \mathbf{M}_{\ddot{\mathbf{U}}}^{-1}$ , for  $k \neq j$  and take



$\mathbf{Z} = (\mathbf{Z}'_1, \dots, \mathbf{Z}'_m)'$ . Now observe that  $(\mathbf{I} - \mathbf{H}_\mathbf{T}) \mathbf{Z} = \mathbf{Z}$ . It follows that

$$\begin{aligned} (\mathbf{I} - \mathbf{H}_\mathbf{X})_i \mathbf{Z} &= (\mathbf{I} - \mathbf{H}_{\ddot{\mathbf{U}}})_i (\mathbf{I} - \mathbf{H}_\mathbf{T}) \mathbf{Z} = (\mathbf{I} - \mathbf{H}_{\ddot{\mathbf{U}}})_i \mathbf{Z} \\ &= \mathbf{Z}_i - \ddot{\mathbf{U}}_i \mathbf{M}_{\ddot{\mathbf{U}}} \sum_{k=1}^m \ddot{\mathbf{U}}'_k \mathbf{W}_k \mathbf{Z}_k = \ddot{\mathbf{U}}_i \mathbf{M}_{\ddot{\mathbf{U}}} \left( \sum_{k \neq j} \ddot{\mathbf{U}}'_k \mathbf{W}_k \ddot{\mathbf{U}} \right) \mathbf{L}^{-1} \mathbf{M}_{\ddot{\mathbf{U}}}^{-1} \\ &= \ddot{\mathbf{U}}_i. \end{aligned}$$

Thus, there exists an  $N \times (r + s)$  matrix  $\mathbf{Z}$  such that  $(\mathbf{I} - \mathbf{H}_{\ddot{\mathbf{X}}})_i \mathbf{Z} = \ddot{\mathbf{U}}_i$ , i.e.,  $\ddot{\mathbf{U}}_i$  is in the column span of  $(\mathbf{I} - \mathbf{H}_\mathbf{X})_i$ . Because  $\mathbf{D}_i \mathbf{W}_i$  is positive definite and  $\ddot{\mathbf{R}}_i$  is a sub-matrix of  $\ddot{\mathbf{U}}_i$ ,  $\mathbf{D}_i \mathbf{W}_i \ddot{\mathbf{R}}_i$  is also in the column span of  $(\mathbf{I} - \mathbf{H}_\mathbf{X})_i$ .

It follows that

$$(23) \quad \ddot{\mathbf{R}}'_i \mathbf{W}_i \mathbf{D}_i \mathbf{V} \mathbf{V}' \mathbf{D}'_i \mathbf{W}_i \ddot{\mathbf{R}}_i = \ddot{\mathbf{R}}'_i \mathbf{W}_i \Phi_i \mathbf{W}_i \ddot{\mathbf{R}}_i.$$

Substituting (23) into (22) demonstrates that  $\mathbf{A}_i$  satisfies criterion (7).

Under the working model, the residuals from cluster  $i$  have mean  $\mathbf{0}$  and variance

$$\text{Var}(\ddot{\mathbf{e}}_i) = (\mathbf{I} - \mathbf{H}_\mathbf{X})_i \Phi (\mathbf{I} - \mathbf{H}_\mathbf{X})'_i,$$

It follows that

$$\begin{aligned} \text{E}(\mathbf{V}^{CR2}) &= \mathbf{M}_{\ddot{\mathbf{R}}} \left[ \sum_{i=1}^m \ddot{\mathbf{R}}'_i \mathbf{W}_i \mathbf{A}_i (\mathbf{I} - \mathbf{H}_\mathbf{X})_i \Phi (\mathbf{I} - \mathbf{H}_\mathbf{X})'_i \mathbf{A}_i \mathbf{W}_i \ddot{\mathbf{R}}_i \right] \mathbf{M}_{\ddot{\mathbf{R}}} \\ &= \mathbf{M}_{\ddot{\mathbf{R}}} \left[ \sum_{i=1}^m \ddot{\mathbf{R}}'_i \mathbf{W}_i \Phi_i \mathbf{W}_i \ddot{\mathbf{R}}_i \right] \mathbf{M}_{\ddot{\mathbf{R}}} \\ &= \text{Var}(\hat{\beta}) \end{aligned}$$

## A.2. Proof of Theorem 2

From the fact that  $\ddot{\mathbf{U}}'_i \mathbf{W}_i \mathbf{T}_i = \mathbf{0}$  for  $i = 1, \dots, m$ , it follows that

$$\begin{aligned} \mathbf{B}_i &= \mathbf{D}_i (\mathbf{I} - \mathbf{H}_{\ddot{\mathbf{U}}})_i (\mathbf{I} - \mathbf{H}_\mathbf{T}) \Phi (\mathbf{I} - \mathbf{H}_\mathbf{T})' (\mathbf{I} - \mathbf{H}_{\ddot{\mathbf{U}}})'_i \mathbf{D}'_i \\ &= \mathbf{D}_i (\mathbf{I} - \mathbf{H}_{\ddot{\mathbf{U}}} - \mathbf{H}_\mathbf{T})_i \Phi (\mathbf{I} - \mathbf{H}_{\ddot{\mathbf{U}}} - \mathbf{H}_\mathbf{T})'_i \mathbf{D}'_i \\ &= \mathbf{D}_i \left( \Phi_i - \ddot{\mathbf{U}}_i \mathbf{M}_{\ddot{\mathbf{U}}} \ddot{\mathbf{U}}'_i - \mathbf{T}_i \mathbf{M}_\mathbf{T} \mathbf{T}'_i \right) \mathbf{D}'_i \end{aligned}$$

and

$$(24) \quad \mathbf{B}_i^+ = (\mathbf{D}'_i)^{-1} \left( \Phi_i - \ddot{\mathbf{U}}_i \mathbf{M}_{\ddot{\mathbf{U}}} \ddot{\mathbf{U}}'_i - \mathbf{T}_i \mathbf{M}_{\mathbf{T}} \mathbf{T}'_i \right)^+ \mathbf{D}_i^{-1}.$$

Let  $\Psi_i = \left( \Phi_i - \ddot{\mathbf{U}}_i \mathbf{M}_{\ddot{\mathbf{U}}} \ddot{\mathbf{U}}'_i \right)^+$ . Using a generalized Woodbury identity (Henderson and Searle, 1981),

$$\Psi_i = \mathbf{W}_i + \mathbf{W}_i \ddot{\mathbf{U}}_i \mathbf{M}_{\ddot{\mathbf{U}}} \left( \mathbf{M}_{\ddot{\mathbf{U}}} - \mathbf{M}_{\ddot{\mathbf{U}}} \ddot{\mathbf{U}}'_i \mathbf{W}_i \ddot{\mathbf{U}}_i \mathbf{M}_{\ddot{\mathbf{U}}} \right)^+ \mathbf{M}_{\ddot{\mathbf{U}}} \ddot{\mathbf{U}}'_i \mathbf{W}_i.$$

It follows that  $\Psi_i \mathbf{T}_i = \mathbf{W}_i \mathbf{T}_i$ . Another application of the generalized Woodbury identity gives

$$\begin{aligned} \left( \Phi_i - \ddot{\mathbf{U}}_i \mathbf{M}_{\ddot{\mathbf{U}}} \ddot{\mathbf{U}}'_i - \mathbf{T}_i \mathbf{M}_{\mathbf{T}} \mathbf{T}'_i \right)^+ &= \Psi_i + \Psi_i \mathbf{T}_i \mathbf{M}_{\mathbf{T}} (\mathbf{M}_{\mathbf{T}} - \mathbf{M}_{\mathbf{T}} \mathbf{T}'_i \Psi_i \mathbf{T}_i \mathbf{M}_{\mathbf{T}})^+ \mathbf{M}_{\mathbf{T}} \mathbf{T}'_i \Psi_i \\ &= \Psi_i + \mathbf{W}_i \mathbf{T}_i \mathbf{M}_{\mathbf{T}} (\mathbf{M}_{\mathbf{T}} - \mathbf{M}_{\mathbf{T}} \mathbf{T}'_i \mathbf{W}_i \mathbf{T}_i \mathbf{M}_{\mathbf{T}})^+ \mathbf{M}_{\mathbf{T}} \mathbf{T}'_i \mathbf{W}_i \\ &= \Psi_i. \end{aligned}$$

The last equality follows from the fact that  $\mathbf{T}_i \mathbf{M}_{\mathbf{T}} (\mathbf{M}_{\mathbf{T}} - \mathbf{M}_{\mathbf{T}} \mathbf{T}'_i \mathbf{W}_i \mathbf{T}_i \mathbf{M}_{\mathbf{T}})^- \mathbf{M}_{\mathbf{T}} \mathbf{T}'_i = \mathbf{0}$  because the fixed effects are nested within clusters. Substituting into (24), we then have that  $\mathbf{B}_i^+ = (\mathbf{D}'_i)^{-1} \Psi_i \mathbf{D}_i^{-1}$ . But

$$\tilde{\mathbf{B}}_i = \mathbf{D}_i (\mathbf{I} - \mathbf{H}_{\ddot{\mathbf{U}}})_i \Phi (\mathbf{I} - \mathbf{H}_{\ddot{\mathbf{U}}})'_i \mathbf{D}'_i = \mathbf{D}_i \left( \Phi_i - \ddot{\mathbf{U}}_i \mathbf{M}_{\ddot{\mathbf{U}}} \ddot{\mathbf{U}}'_i \right) \mathbf{D}'_i = \mathbf{D}_i \Psi_i^+ \mathbf{D}'_i,$$

and so  $\mathbf{B}_i^+ = \tilde{\mathbf{B}}_i^+$ . It follows that  $\mathbf{A}_i = \tilde{\mathbf{A}}_i$  for  $i = 1, \dots, m$ .

## APPENDIX B: DETAILS OF SIMULATION STUDY

This appendix provides further details regarding the design of the simulations reported in Section 4. The simulations examined six distinct study designs. Outcomes are measured for  $n$  units (which may be individuals, as in a cluster-randomized or block-randomized design, or time-points, as in a difference-in-differences panel) in each of  $m$  clusters under one of three treatment conditions. Suppose that there are  $G$  groups of units that share an identical pattern of treatment assignments, each of size  $m_g$ . Let  $n_{ghi}$  denote the number of units at which cluster  $i$  in group  $g$  is observed under

condition  $h$ , for  $i = 1, \dots, m$ ,  $g = 1, \dots, G$ , and  $h = 1, 2, 3$ . The following six designs were simulated:

1. A balanced, block-randomized design, with an un-equal allocation within each block. In the balanced design, the treatment allocation is identical for each block, with  $G = 1$ ,  $m_1 = m$ ,  $n_{11i} = n/2$ ,  $n_{12i} = n/3$ , and  $n_{13i} = n/6$ .
2. An unbalanced, block-randomized design, with two different patterns of treatment allocation. Here,  $G = 2$ ,  $m_1 = m_2 = m/2$ ,  $n_{11i} = n/2$ ,  $n_{12i} = n/3$ ,  $n_{13i} = n/6$ ,  $n_{21i} = n/3$ ,  $n_{22i} = 5n/9$ , and  $n_{23i} = n/9$ .
3. A balanced, cluster-randomized design, in which units are nested within clusters and an equal number of clusters are assigned to each treatment condition. Here,  $G = 3$ ,  $m_g = m/3$ , and  $n_{ghi} = n$  for  $g = h$  and zero otherwise.
4. An unbalanced, cluster-randomized design, in which units are nested within clusters but the number of clusters assigned to each condition is not equal. Here,  $G = 3$ ;  $m_1 = 0.5m$ ,  $m_2 = 0.3m$ ,  $m_3 = 0.2m$ ; and  $n_{ghi} = n$  for  $g = h$  and zero otherwise.
5. A balanced difference-in-differences design, with two patterns of treatment allocation ( $G = 2$ ) and clusters allocated equally to each pattern ( $m_1 = m_2 = m/2$ ). Here, half of the clusters are observed under the first treatment condition only ( $n_{11i} = n$ ) and the remaining half are observed under all three conditions, with  $n_{21i} = n/2$ ,  $n_{22i} = n/3$ , and  $n_{23i} = n/6$ .
6. An unbalanced difference-in-differences design, again with two patterns of treatment allocation ( $G = 2$ ), but where  $m_1 = 2m/3$  clusters are observed under the first treatment condition only ( $n_{11i} = n$ ) and the remaining  $m_2 = m/3$  clusters are observed under all three conditions, with  $n_{21i} = n/2$ ,  $n_{22i} = n/3$ , and  $n_{23i} = n/6$ .

## REFERENCES

- Achilles, C. M., Bain, H. P., Bellott, F., Boyd-Zaharias, J., Finn, J., Folger, J., Johnston, J. and Word, E. (2008), ‘Tennessee’s Student Teacher Achievement Ratio (STAR) project’.  
**URL:** <http://hdl.handle.net/1902.1/10766>
- Angrist, J. D. and Lavy, V. (2009), ‘The effects of high stakes high school achievement awards : Evidence from a randomized trial’, *American Economic Review* **99**(4), 1384–1414.
- Angrist, J. D. and Pischke, J. (2009), *Mostly harmless econometrics: An empiricist’s companion*, Princeton University Press, Princeton, NJ.
- Angrist, J. D. and Pischke, J.-S. (2014), *Mastering’metrics: The Path from Cause to Effect*, Princeton University Press.
- Arellano, M. (1987), ‘Computing robust standard errors for within-groups estimators’, *Oxford Bulletin of Economics and Statistics* **49**(4), 431–434.
- Arellano, M. (1993), ‘On the testing of correlated effects with panel data’, *Journal of Econometrics* **59**(1-2), 87–97.
- Banerjee, S. and Roy, A. (2014), *Linear Algebra and Matrix Analysis for Statistics*, Taylor & Francis, Boca Raton, FL.
- Bell, R. M. and McCaffrey, D. F. (2002), ‘Bias reduction in standard errors for linear regression with multi-stage samples’, *Survey Methodology* **28**(2), 169–181.
- Bertrand, M., Duflo, E. and Mullainathan, S. (2004), ‘How much should we trust differences-in-differences estimates?’, *Quarterly Journal of Economics* **119**(1), 249–275.
- Cameron, A. C., Gelbach, J. B. and Miller, D. (2008), ‘Bootstrap-based improvements for inference with clustered errors’, *The Review of Economics and Statistics* **90**(3), 414–427.
- Cameron, A. C., Gelbach, J. B. and Miller, D. L. (2011), ‘Robust inference with multiway clustering’, *Journal of Business & Economic Statistics* **29**(2), 238–249.
- Cameron, A. C. and Miller, D. L. (2015), A practitioner’s guide to cluster-robust inference.
- Carpenter, C. and Dobkin, C. (2011), ‘The minimum legal drinking age and public health’, *Journal of Economic Perspectives* **25**(2), 133–156.
- Carter, A. V., Schnepel, K. T. and Steigerwald, D. G. (2013), Asymptotic Behavior of a t Test Robust to Cluster Heterogeneity.

- Chetty, R., Friedman, J. N., Hilger, N., Saez, E., Schanzenbach, D. W. and Yagan, D. (2011), ‘How does your kindergarten classroom affect your earnings? Evidence from Project STAR’, *The Quarterly Journal of Economics* **126**(4), 1593–1660.
- Conley, T. G. and Taber, C. R. (2011), ‘Inference with Difference in Differences with a Small Number of Policy Changes’, *Review of Economics and Statistics* **93**(1), 113–125.
- Donald, S. G. and Lang, K. (2007), ‘Inference with difference-in-differences and other panel data’, *Review of Economics and Statistics* **89**(2), 221–233.
- Hansen, C. B. (2007), ‘Asymptotic properties of a robust variance matrix estimator for panel data when T is large’, *Journal of Econometrics* **141**, 597–620.
- Henderson, H. V. and Searle, S. R. (1981), ‘On deriving the inverse of a sum of matrices’, *Siam Review* **23**(1), 53–60.
- Huber, P. J. (1967), The behavior of maximum likelihood estimates under nonstandard conditions, in ‘Proceedings of the fifth Berkeley symposium on Mathematical Statistics and Probability’, University of California Press, Berkeley, CA, pp. 221–233.
- Ibragimov, R. and Müller, U. K. (2010), ‘t-Statistic based correlation and heterogeneity robust inference’, *Journal of Business & Economic Statistics* **28**(4), 453–468.
- Imbens, G. W. and Kolesar, M. (2015), Robust Standard Errors in Small Samples: Some Practical Advice.  
**URL:** <https://www.princeton.edu/~mkolesar/papers/small-robust.pdf>
- Krishnamoorthy, K. and Yu, J. (2004), ‘Modified Nel and Van der Merwe test for the multivariate BehrensFisher problem’, *Statistics & Probability Letters* **66**(2), 161–169.
- Krueger, A. and Whitmore, D. (2001), ‘The effect of attending a small class in the early grades on college-test taking and middle school test results: Evidence from Project STAR’, *The Economic Journal* **111**(468), 1–28.
- Lee, D. S. and Card, D. (2008), ‘Regression discontinuity inference with specification error’, *Journal of Econometrics* **142**(2), 655–674.
- Liang, K.-Y. and Zeger, S. L. (1986), ‘Longitudinal data analysis using generalized linear models’, *Biometrika* **73**(1), 13–22.
- MacKinnon, J. G. (2013), Thirty years of heteroskedasticity-robust inference, in X. Chen and N. R. Swanson, eds, ‘Recent Advances and Future Directions in Causality, Prediction, and Specification Analysis’, Springer New York, New York, NY.
- Mackinnon, J. G. (2014), Wild cluster bootstrap confidence intervals.

- MacKinnon, J. G. and White, H. (1985), ‘Some heteroskedasticity-consistent covariance matrix estimators with improved finite sample properties’, *Journal of Econometrics* **29**, 305–325.
- Mancl, L. A. and DeRouen, T. A. (2001), ‘A covariance estimator for GEE with improved small-sample properties’, *Biometrics* **57**(1), 126–134.
- McCaffrey, D. F. and Bell, R. M. (2006), ‘Improved hypothesis testing for coefficients in generalized estimating equations with small samples of clusters.’, *Statistics in medicine* **25**(23), 4081–98.
- McCaffrey, D. F., Bell, R. M. and Botts, C. H. (2001), Generalizations of biased reduced linearization, in ‘Proceedings of the Annual Meeting of the American Statistical Association’, number 1994.
- Nel, D. and van der Merwe, C. (1986), ‘A solution to the multivariate Behrens-Fisher problem’, *Communications in Statistics - Theory and Methods* **15**(12), 3719–3735.
- Pan, W. and Wall, M. M. (2002), ‘Small-sample adjustments in using the sandwich variance estimator in generalized estimating equations.’, *Statistics in medicine* **21**(10), 1429–41.
- Satterthwaite, F. E. (1946), ‘An approximate distribution of estimates of variance components’, *Biometrics bulletin* **2**(6), 110–114.
- Schanzenbach, D. W. (2006), ‘What have researchers learned from Project STAR?’, *Brookings Papers on Education Policy* **2006**(1), 205–228.
- Tipton, E. (2015), ‘Small sample adjustments for robust variance estimation with meta-regression.’, *Psychological Methods* **20**(3), 375–393.
- Tipton, E. and Pustejovsky, J. E. (2015), ‘Small-sample adjustments for tests of moderators and model fit using robust variance estimation in meta-regression’, *Journal of Educational and Behavioral Statistics* **40**(6), 604–634.
- Webb, M. and MacKinnon, J. G. (2013), Wild Bootstrap Inference for Wildly Different Cluster Sizes.
- White, H. (1980), ‘A heteroskedasticity-consistent covariance matrix estimator and a direct test for heteroskedasticity’, *Econometrica* **48**(4), 817–838.
- White, H. (1984), *Asymptotic theory for econometricians*, Academic Press, Inc., Orlando, FL.
- Wooldridge, J. M. (2002), *Econometric Analysis of Cross Section and Panel Data*, MIT Press, Cambridge, MA.
- Wooldridge, J. M. (2003), ‘Cluster-sample methods in applied econometrics’, *The American Economic Review* **93**(2), 133–138.

- Wooldridge, J. M. (2010), *Econometric Analysis of Cross Section and Panel Data*, 2nd edn, MIT Press, Cambridge, MA.
- Zhang, J.-T. (2012*a*), ‘An approximate degrees of freedom test for heteroscedastic two-way ANOVA’, *Journal of Statistical Planning and Inference* **142**(1), 336–346.
- Zhang, J.-T. (2012*b*), ‘An approximate Hotelling T<sup>2</sup> -test for heteroscedastic one-way MANOVA’, *Open Journal of Statistics* **2**, 1–11.