

MacKinnon, James G.; Webb, Matthew

Working Paper

When and how to deal with clustered errors in regression models

Queen's Economics Department Working Paper, No. 1421

Provided in Cooperation with:

Queen's University, Department of Economics (QED)

Suggested Citation: MacKinnon, James G.; Webb, Matthew (2019) : When and how to deal with clustered errors in regression models, Queen's Economics Department Working Paper, No. 1421, Queen's University, Department of Economics, Kingston (Ontario)

This Version is available at:

<http://hdl.handle.net/10419/230574>

Standard-Nutzungsbedingungen:

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

Terms of use:

Documents in EconStor may be saved and copied for your personal and scholarly purposes.

You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.

If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.



Queen's Economics Department Working Paper No. 1421

When and How to Deal with Clustered Errors in Regression Models

James G. MacKinnon
Queen's University

Matthew D. Webb
Carleton University

Department of Economics
Queen's University
94 University Avenue
Kingston, Ontario, Canada
K7L 3N6

8-2019

When and How to Deal With Clustered Errors in Regression Models*

James G. MacKinnon[†] Matthew D. Webb
Queen's University Carleton University
jgm@econ.queensu.ca matt.webb@carleton.ca

August 26, 2019

Abstract

We discuss when and how to deal with possibly clustered errors in linear regression models. Specifically, we discuss situations in which a regression model may plausibly be treated as having error terms that are arbitrarily correlated within known clusters but uncorrelated across them. The methods we discuss include various covariance matrix estimators, possibly combined with various methods of obtaining critical values, several bootstrap procedures, and randomization inference. Special attention is given to models with few treated clusters and clusters that vary in size, where inference may be problematic. Two empirical examples and a simulation experiment illustrate the methods we discuss and the concerns we raise.

*We thank the Social Sciences and Humanities Research Council of Canada (SSHRC) for financial support. We are grateful to Mehtab Hanzroh for his excellent research assistance. We benefited from the comments of Alfonso Flores-Lagunes and of participants at the CIREQ 2019 Bootstrap Conference and the 2019 Canadian Economics Association Annual Meeting.

[†]Corresponding author. Address: Department of Economics, 94 University Avenue, Queen's University, Kingston, Ontario K7L 3N6, Canada. Email: jgm@econ.queensu.ca. Tel. 613-533-2293. Fax 613-533-6668.

1 Introduction

When estimating regression models for cross-section data, it used to be common for investigators to assume that the error terms (or disturbances) for any pair of observations are uncorrelated. Although this assumption may seem natural, it is actually very strong, and it has two important implications. First, it means that inference can safely be based on “robust” (that is, heteroskedasticity-robust) covariance matrix estimators. For large samples, these estimators are usually quite reliable, although there can be exceptions when a few observations have high leverage (MacKinnon 2013).

A more profound implication of the assumption that error terms are uncorrelated is that information about the parameters accumulates at a rate proportional to the square root of the sample size. If we estimate the same model on two datasets, one with M observations and one with $N = \phi M$ observations, where $\phi \gg 1$, the (true) standard errors for the second set of estimates should be approximately $\phi^{-1/2}$ times those for the first set. Of course, this statement assumes that the investigator does not take advantage of the larger sample size by estimating a more complicated model for the sample of size N than for the one of size M . In practice, models tend to become more complicated as sample sizes increase, so that standard errors are not actually proportional to one over the square root of the sample size.

As we discuss in Section 3, it is much less common than it once was to assume that the error terms of regression models are uncorrelated. Instead, investigators commonly assume that they are “clustered.” The sample is divided into clusters (which might be associated, for example, with schools, firms, villages, counties, or states), and the disturbances for observations within each cluster are allowed to be correlated. This requires the use of a covariance matrix estimator that is robust to arbitrary patterns of both heteroskedasticity and intra-cluster correlation; see Section 2.

The use of cluster-robust variance estimators in empirical microeconomics began after such an estimator became available in Stata (Rogers 1993). It became much more widespread after a very influential paper (Bertrand, Duflo, and Mullainathan 2004) showed that inferences for difference-in-differences (DiD) estimators based on standard errors that ignore autocorrelation within geographical clusters can be extremely unreliable; see Section 6. Cameron and Miller (2015) is an influential survey. More recent surveys include Esarey and Menger (2019) and MacKinnon (2019).

Failing to allow for intra-cluster correlation has particularly serious consequences when the sample size is large. Thus one important reason for the increased use of cluster-robust standard errors in recent years is that sample sizes have become larger. When cluster sizes are growing with the sample size N , information about the parameters accumulates at a rate slower than \sqrt{N} . In extreme cases (MacKinnon 2016), it can accumulate very much more slowly. Whether or not there is intra-cluster correlation, heteroskedasticity-robust standard errors are always roughly proportional to $1/\sqrt{N}$. Therefore, when intra-cluster correlation is actually present, the ratio of a true (cluster-robust) standard error to one that is only heteroskedasticity-robust increases without limit as $N \rightarrow \infty$. This implies that errors of inference become more severe as the sample size increases.

In Section 2, we discuss methods of cluster-robust inference based on t -statistics and Wald statistics. In Section 3, we discuss why it often makes sense to divide the sample into clusters and allow for intra-cluster correlation. In Section 4, we discuss how to cluster. The

investigator has to choose the appropriate dimension(s) and level(s) of clustering, and this is often not easy. In Section 5, we discuss several commonly-encountered cases in which using cluster-robust standard errors in the usual way can lead to very serious errors of inference. We also discuss methods that can be used to obtain more reliable inferences, including the wild cluster bootstrap (Cameron, Gelbach, and Miller 2008), the wild bootstrap (MacKinnon and Webb 2018), and randomization inference. In Section 6, we discuss two empirical examples that illustrate some of the important issues. Section 7 presents some simple Monte Carlo simulations which demonstrate the consequences of getting the level of clustering correct or incorrect. Section 8 concludes.

2 Regression models with clustered errors

Consider the linear regression model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}, \quad \text{E}(\mathbf{u}|\mathbf{X}) = \mathbf{0}, \quad \text{E}(\mathbf{u}\mathbf{u}') = \boldsymbol{\Omega}, \quad (1)$$

where \mathbf{y} and \mathbf{u} are $N \times 1$ vectors of observations and disturbances, \mathbf{X} is an $N \times K$ matrix of exogenous covariates, and $\boldsymbol{\beta}$ is a $K \times 1$ parameter vector. When the $N \times N$ covariance matrix $\boldsymbol{\Omega}$ is equal to $\sigma^2\mathbf{I}$, the ordinary least squares (OLS) estimator $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ is the best linear unbiased estimator of $\boldsymbol{\beta}$, and we can make inferences based on the estimated covariance matrix $s^2(\mathbf{X}'\mathbf{X})^{-1}$, where s^2 is $1/(N - K)$ times the sum of squared residuals. When $\boldsymbol{\Omega}$ is diagonal with diagonal elements that differ, OLS is no longer efficient, but it is still consistent, and we can make inferences by using “robust” standard errors based on a heteroskedasticity-consistent covariance matrix estimator, or HCCME (White 1980).

In many cases, however, as we discuss in Sections 3 and 4, there are very good reasons to believe that $\boldsymbol{\Omega}$ is not a diagonal matrix. Suppose instead that the data can be divided into G clusters, indexed by g , where the g^{th} cluster has N_g observations. Then $\boldsymbol{\Omega}$ is assumed to be block-diagonal, with G diagonal blocks that correspond to the G clusters:

$$\boldsymbol{\Omega} = \begin{bmatrix} \boldsymbol{\Omega}_1 & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Omega}_2 & \dots & \mathbf{0} \\ \vdots & \vdots & & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \boldsymbol{\Omega}_G \end{bmatrix}. \quad (2)$$

As the notation suggests, each of the $\boldsymbol{\Omega}_g$ is an $N_g \times N_g$ positive semidefinite matrix, and every element of the off-diagonal blocks in $\boldsymbol{\Omega}$ is assumed to be zero.

The true covariance matrix of the OLS estimator $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ in the model given by (1) and (2) is

$$\text{Var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\Omega}\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} = (\mathbf{X}'\mathbf{X})^{-1}\left(\sum_{g=1}^G \mathbf{X}'_g\boldsymbol{\Omega}_g\mathbf{X}_g\right)(\mathbf{X}'\mathbf{X})^{-1}, \quad (3)$$

where the $N_g \times K$ matrix \mathbf{X}_g contains the rows of \mathbf{X} that belong to the g^{th} cluster. Thus the middle factor is actually the sum of G matrices, each of them $K \times K$.

The matrix (3) can be estimated by using the outer product of the residual vector $\hat{\mathbf{u}}_g$ with itself to estimate $\boldsymbol{\Omega}_g$ for all g . This yields a cluster-robust variance estimator, or CRVE.

By far the most widely-used version is

$$\text{CV}_1: \quad \frac{G(N-1)}{(G-1)(N-K)} (\mathbf{X}'\mathbf{X})^{-1} \left(\sum_{g=1}^G \mathbf{X}'_g \hat{\mathbf{u}}_g \hat{\mathbf{u}}'_g \mathbf{X}_g \right) (\mathbf{X}'\mathbf{X})^{-1}. \quad (4)$$

The first factor here is analogous to the factor $N/(N-K)$ used in the conventional heteroskedasticity-robust HC_1 covariance matrix (MacKinnon and White 1985), which replaces the middle matrix in (4) by $\sum_{i=1}^N \hat{u}_i^2 \mathbf{X}'_i \mathbf{X}_i$. This factor makes CV_1 larger when either G or N becomes smaller, in order to offset the tendency for OLS residuals to be too small. CV_1 evidently reduces to HC_1 when $G = N$, so that each cluster contains just one observation.

The CRVE (4), like all robust covariance matrix estimators, is a “sandwich” estimator. The filling in the sandwich is supposed to estimate the corresponding filling in (3). However, unlike the individual matrices in the summation in (3), the ones in (4) have rank 1, even though they are $K \times K$. Therefore, the individual components of the filling in (4) cannot possibly provide consistent estimators of the corresponding components of the filling in (3). Moreover, unless $G \geq K$, the matrix (4) cannot have full rank. It will have rank at most G in some cases and rank at most $G-1$ in others.

Cluster-robust variance estimators were first proposed by Liang and Zeger (1986) and Arellano (1987). They became available in Stata about half a decade later (Rogers 1993). However, econometricians did not study their properties under general assumptions until much later. Bester, Conley, and Hansen (2011) showed that, under quite restrictive conditions with N increasing and G fixed, cluster-robust t -statistics for $\beta_j = 0$, where β_j is any element of $\boldsymbol{\beta}$, are asymptotically distributed as $t(G-1)$. This result justifies the use of the $t(G-1)$ distribution for calculating critical values and P values, something that has been the default in Stata since 1993.

More recently, Djogbenou, MacKinnon, and Nielsen (2019) proved that cluster-robust t -statistics are asymptotically normally distributed under rather weak conditions. These require G to increase with N and allow the N_g to increase as well, but not too fast. There are also limits on how much the cluster sizes can vary. Using a similar framework, Hansen and Lee (2019) proved the asymptotic validity of cluster-robust inference based on the standard normal distribution combined with covariance matrix estimators similar to (4) for a wide variety of linear and nonlinear econometric models, including ones estimated by two-stage least squares, the generalized method of moments (GMM), and maximum likelihood.

Although it is by far the most widely used CRVE, the matrix CV_1 defined in (4) is not the only one. An estimator with somewhat better finite-sample properties, which was proposed by Bell and McCaffrey (2002) and advocated by Imbens and Kolesár (2016), is

$$\text{CV}_2: \quad (\mathbf{X}'\mathbf{X})^{-1} \left(\sum_{g=1}^G \mathbf{X}'_g \mathbf{M}_{gg}^{-1/2} \hat{\mathbf{u}}_g \hat{\mathbf{u}}'_g \mathbf{M}_{gg}^{-1/2} \mathbf{X}_g \right) (\mathbf{X}'\mathbf{X})^{-1}, \quad (5)$$

where $\mathbf{M}_{gg}^{-1/2}$ is the inverse symmetric square root of the matrix $\mathbf{M}_{gg} \equiv \mathbf{I}_{N_g} - \mathbf{X}_g (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'_g$, which is the g^{th} diagonal block of the $N \times N$ matrix $\mathbf{M}_{\mathbf{X}} \equiv \mathbf{I} - \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'$. Instead of multiplying by a scalar factor, CV_2 replaces the residual subvectors $\hat{\mathbf{u}}_g$ by rescaled subvectors $\mathbf{M}_{gg}^{-1/2} \hat{\mathbf{u}}_g$. It reduces to the HC_2 HCCME discussed in MacKinnon and White (1985) when $G = N$. CV_2 can be calculated efficiently in R using the package `clubSandwich` (Pustejovsky 2017). Although CV_2 seems to yield larger and more accurate standard errors than

CV₁, it is considerably more expensive to compute when the clusters are large, because it requires finding the inverse symmetric square root of \mathbf{M}_{gg} for each cluster. For sufficiently large clusters, this can be numerically infeasible (MacKinnon and Webb 2018). Recently, Jackson (2019) proposed an alternative estimator which estimates the “filling” of the CRVE by estimating a common variance and common correlation of residuals within each cluster.

As noted above, the standard way to make inferences about an individual element of $\boldsymbol{\beta}$, say β_j , is to use the cluster-robust t -statistic

$$t_j = \frac{\hat{\beta}_j - \beta_{j0}}{s_j}, \quad (6)$$

where β_{j0} is the value under the null hypothesis and s_j is the square root of the j^{th} diagonal element of either CV₁ or CV₂. The statistic t_j is then compared with the $t(G-1)$ distribution. A $(1 - \alpha)\%$ confidence interval for β_j would be

$$\left[\hat{\beta}_j - s_j C_{t(G-1)}(1 - \alpha/2), \hat{\beta}_j + s_j C_{t(G-1)}(1 - \alpha/2) \right], \quad (7)$$

where $C_{t(G-1)}(1 - \alpha/2)$ is the $1 - \alpha/2$ quantile of the $t(G-1)$ distribution. When G is small, the latter can be considerably larger than the corresponding quantile of the standard normal distribution. In combination with the fact that cluster-robust standard errors are often much larger than heteroskedasticity-robust ones, this can make the interval (7) much wider than a corresponding “robust” interval.

When there are two or more restrictions to be tested, we can use a Wald test. In order to test the hypothesis that $\mathbf{R}\boldsymbol{\beta} = \mathbf{r}$, where \mathbf{R} is an $r \times k$ matrix and \mathbf{r} is an $r \times 1$ vector, we compute the Wald statistic

$$W(\hat{\boldsymbol{\beta}}) = (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})'(\mathbf{R}\hat{\mathbf{V}}\mathbf{R}')^{-1}(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r}), \quad (8)$$

where $\hat{\mathbf{V}}$ could be either CV₁ or CV₂. The statistic $W(\hat{\boldsymbol{\beta}})$ is then compared with critical values from either the $\chi^2(r)$ distribution or, preferably, r times the $F(r, G-1)$ distribution. $W(\hat{\boldsymbol{\beta}})$ cannot be computed if $r > G$, and perhaps not if $r = G$. This Wald test is likely to over-reject severely when r is not much smaller than G . In such cases, it is strongly advised to use the wild cluster bootstrap; see Section 5.

2.1 Multi-way clustering

The CRVEs given in (4) and (5) are designed to handle arbitrary within-cluster correlation in a single dimension. Cameron, Gelbach, and Miller (2011) and Thompson (2011) independently proposed extensions to handle clustering in two or more dimensions. For example, there might be one set of clusters based on geography and another set based on time. Every observation is assumed to belong to one cluster in each of the two dimensions.

The two 2011 papers proposed covariance matrix estimators but did not work out their asymptotic properties. Recent work has developed the theory of multi-way cluster-robust estimators. Davezies, D’Haultfoeulle, and Guyonvarch (2018) proposed an alternative multi-way CRVE and proved its asymptotic validity, which is technically challenging. Menzel (2018) proposed a multi-way bootstrap procedure for inference on sample means. MacKinnon, Nielsen, and Webb (2019) compared the properties of two forms of two-way CRVE and showed that several variants of the wild cluster bootstrap (Subsection 5.1) can be combined with a two-way CRVE to obtain more reliable inferences about regression coefficients.

3 When to cluster

The simplest way to model intra-cluster correlation is to assume that there are cluster-specific random effects, say v_g . The i^{th} observation in the g^{th} cluster is then equal to

$$y_{gi} = \mathbf{X}_{gi}\boldsymbol{\beta} + u_{gi} = \mathbf{X}_{gi}\boldsymbol{\beta} + v_g + \epsilon_{gi}, \quad (9)$$

where the v_g are independently distributed with variance σ_v^2 and the ϵ_{gi} are independently distributed with variance σ_ϵ^2 . This implies that the variance of u_{gi} is $\sigma_v^2 + \sigma_\epsilon^2$, the correlation between disturbances in different clusters is zero, and the correlation between disturbances within the same cluster is $\rho_u = \sigma_v^2 / (\sigma_v^2 + \sigma_\epsilon^2)$.

Although the random-effects model (9) is simple and appealing, there are probably not many datasets for which it is actually appropriate. By assuming that all of the correlation within each cluster comes from a single cluster-specific effect v_g , which affects all observations equally, it rules out any variation in intra-cluster correlations. More realistically, we might expect there to be several cluster-specific effects for each cluster, and for them to affect different observations differently. It might well also be the case that the v_g , the ϵ_{gi} , or both of them are heteroskedastic, with variances that depend on the regressors.

The random-effects model assumes that the v_g are uncorrelated with all the regressors. This is often a very strong assumption, and it will lead to inconsistent estimates if it is false. The classic way to solve this problem is to treat the v_g as fixed effects, that is, as constants to be estimated instead of as random effects. Then the original regression (1) becomes

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{D}\mathbf{v} + \boldsymbol{\epsilon}, \quad (10)$$

where \mathbf{D} is an $N \times G$ matrix, with $D_{gi} = 1$ for observations that belong to cluster g and $D_{gi} = 0$ otherwise. Of course, one column of \mathbf{D} must be omitted if \mathbf{X} contains a constant term or the equivalent.

The fixed-effects model (10) is very popular. However, it cannot be used whenever any of the regressors varies only at the cluster level, because that regressor would simply be a linear combination of the columns of \mathbf{D} . Such a model was estimated by Riddell (1979), which led Kloek (1981) to study models in which a dependent variable measured at the individual level is regressed on data measured at the cluster level. The paper showed that conventional standard errors for OLS estimates are biased downwards, often very seriously so, when the disturbances follow the random-effects model (9) with $\rho_u > 0$. Under the assumptions of Kloek (1981), the appropriate way to make inferences is to use feasible generalized least squares (FGLS) for the random effects model. In some cases, this is numerically equal to OLS, but with a different covariance matrix.

Even when a model includes fixed effects, there is no reason to believe that they account for all of the intra-cluster correlation. For example, the log-earnings equations of Bertrand, Duflo, and Mullainathan (2004) using U.S. data contain both state and year fixed effects, and yet failing to cluster at the state level leads to very severe over-rejections in their placebo-law experiments.¹ Note that the sample size is over 500,000 in this case. This illustrates the

¹Placebo-law experiments are an ingenious way to evaluate the performance of inferential procedures using real data. Every replication uses the same data for the regressand and all but one of the regressors. The only thing that differs across replications is the regressor of interest, a treatment dummy that affects

fact, discussed below, that even very small amounts of intra-cluster correlation can have a large effect on the accuracy of inferences when N is large.

In two influential papers, [Moulton \(1986, 1990\)](#) demonstrated via empirical examples that intra-cluster correlation is widespread and that failing to account for it can lead to standard errors that are much too small. Moreover, [Moulton \(1986\)](#) showed that, in the context of the random-effects model, the square of the ratio of the true standard error to the conventional OLS standard error is

$$1 + \rho_x \rho_u \left(\frac{\text{Var}(N_g)}{\bar{N}_g} + \bar{N}_g - 1 \right), \quad (11)$$

where ρ_x is the intra-cluster correlation of the regressor of interest (after it is projected off all other regressors), and \bar{N}_g is the mean of the N_g . The quantity (11) is sometimes called the ‘‘Moulton factor.’’ It is evidently one when either ρ_u or ρ_x is zero, increases with both ρ_u and ρ_x , and increases without limit as either \bar{N}_g or the ratio of $\text{Var}(N_g)$ to \bar{N}_g increases.

Although the Moulton factor (11) strictly applies only to the random-effects model (9), it provides useful guidance in many cases. In particular, it makes it clear that the extent of intra-cluster correlation for the regressors is just as important as the extent of intra-cluster correlation for the disturbances, and it shows that the errors of inference we make by not allowing for intra-cluster correlation become more severe as the clusters become larger and/or more variable in size.

Another way to see why cluster sizes matter is to consider the matrices $\mathbf{X}'_g \boldsymbol{\Omega}_g \mathbf{X}_g$ that appear in (3). We can write a typical element of one of these matrices as

$$\sum_{i=1}^{N_g} \Omega_{g,ii} X_{gki}^2 + 2 \sum_{i=1}^{N_g} \sum_{j=i+1}^{N_g} \Omega_{g,ij} X_{gki} X_{gkj}, \quad (12)$$

where $\Omega_{g,ij}$ is the ij^{th} element of $\boldsymbol{\Omega}_g$, and X_{gki} is the k^{th} element of the row of \mathbf{X} corresponding to the i^{th} observation within the g^{th} cluster. If all the off-diagonal elements of $\boldsymbol{\Omega}_g$ were zero, the second term here would be zero, and expression (12) would be $O(N_g)$.² This would also be the case if $E(X_{gki} X_{gkj}) = 0$. But when there is intra-cluster correlation of both the disturbances and the regressors, the second term, which is $O(N_g^2)$, does not vanish. Even if the $\Omega_{g,ij}$ are very small, the second term will ultimately dominate, and expression (12) itself must be $O(N_g^2)$.

The implications of this result are profound. When there is no intra-cluster correlation, the covariance matrix (3) is $O(N^{-1})O(N)O(N^{-1}) = O(N^{-1})$, as usual. But when there is intra-cluster correlation and the number of clusters G is fixed, the covariance matrix is instead $O(N^{-1})O(N^2)O(N^{-1}) = O(1)$, because all of the N_g must be proportional to N . This implies that $\hat{\boldsymbol{\beta}}$ is not a consistent estimator when the sample size goes to infinity with a fixed number of clusters. For this reason, any proof of consistency, such as the one in [Djogbenou et al. \(2019\)](#), requires that G tend to infinity along with N .

certain clusters in certain years. Since the treatment dummies are generated randomly, valid statistical procedures should reject the null hypothesis about as often as the level of the test. See also [MacKinnon \(2016\)](#), [MacKinnon and Webb \(2017a\)](#), and [Brewer, Crossley, and Joyce \(2018\)](#).

²Here we have used the ‘‘same-order’’ or ‘‘big O’’ notation, which is a convenient way to indicate how a quantity changes with the sample size N (or, in this case, the cluster size N_g). Informally, $x = O(n)$ means that x behaves like n as n becomes large.

The number of clusters G does not have to be proportional to N for $\hat{\beta}$ to be consistent. But when G is increasing more slowly than N , both (3) and the CRVEs (4) and (5) that estimate it consistently will tend to zero at a rate slower than $O(N^{-1})$ as $N \rightarrow \infty$. This can make reliable inference more difficult for large samples than for small ones. Because a covariance matrix that is robust only to heteroskedasticity is always $O(N^{-1})$, the errors of inference that we make if we fail to allow for intra-cluster correlation may be very severe when N is large. The same thing is true if we cluster at too fine a level, for instance, clustering by city rather than by state.

As an example, suppose that there are actually G equal-sized clusters, but we divide the data more finely so that there are instead $4G$ equal-sized clusters. Then the true covariance matrix (3) will involve G expressions like (12), but the CRVE will involve $4G$ such expressions. For the off-diagonal terms, each of these will have $1/16$ as many elements in the double summation. The net effect is that the CRVE will sum over only $1/4$ of the off-diagonal elements that it should be summing over. Unless the off-diagonal elements that it missed happen to be very small, it is likely that the CRVE will seriously underestimate (3). The Monte Carlo simulation results in Section 7 provide an example of this.

3.1 Fixed effects and clustered standard errors

Let us assume for the moment that the appropriate level at which to cluster is known. In practice, of course, it often is not known, but we postpone that issue to Section 4. Then two important issues are whether or not to include cluster fixed effects and whether or not to compute cluster-robust standard errors. In our view, these are quite different issues. Whether or not to include fixed effects is a modeling decision. Including them may or may not be feasible and may or may not be desirable. In contrast, whether to compute cluster-robust or merely heteroskedasticity-robust standard errors is an inference decision. Especially when cluster sizes are large, it is usually a bad idea not to allow for intra-cluster correlation (but see Subsection 3.3 and Section 5).

In the past, it was often considered sufficient to include cluster fixed effects and simply use heteroskedasticity-robust standard errors. If the pattern of intra-cluster correlation follows the simple random-effects model of (9), then the fixed effects will explain the cluster-specific v_g terms, and all that remains of the error terms will be the ϵ_{gi} , which by assumption are independent. However, if the pattern of intra-cluster correlation is more complicated, as common sense and the placebo-law results discussed above suggest, it is better to combine cluster fixed effects with cluster-robust standard errors. In our view, this should generally be the default approach, unless there are good reasons not to include fixed effects.

Another approach is just to use a CRVE, without any fixed effects. This makes sense in two cases. The first case arises whenever it is impossible to include cluster fixed effects because of invariant regressors (Kloek 1981; Moulton 1986). If the number of clusters is large, cluster sizes do not vary excessively, and the estimates do not depend strongly on a small subset of atypical clusters, then inference based directly on cluster-robust standard errors should be reliable. If any of these conditions is violated, however, it can be essential to use more sophisticated methods, such as bootstrap or randomization inference methods; see Section 5.

The second case in which it makes sense to use a CRVE without fixed effects is when

there is a large number of clusters, all roughly the same size. Unless the distribution of the estimated fixed effects is of interest (for example, they might be interpreted as measures of classroom or teacher quality), it is often better not to estimate them, both to save computer time and because including fixed effects can have undesirable consequences. In particular, it can lead to severe biases when estimating short dynamic panel models (Nickell 1981), such as models for earnings at the individual level, with lagged earnings as an explanatory variable, where there are many individuals and few time periods. Using a CRVE with one cluster per individual allows for the disturbances associated with each individual to be correlated over time in a more flexible way than individual fixed effects, and it avoids the possibly severe biases associated with the latter.

3.2 Spatial (auto)correlation

Another cause for concern is that the disturbances may display spatial (auto)correlation, so that they are correlated across as well as within geographic clusters. Barrios, Diamond, Imbens, and Kolesár (2012) suggested that many regressors are correlated beyond state boundaries and that researchers should investigate this possibility. More recently, Kelly (2019) argued that many “long difference” or “persistence” papers are likely to have disturbances that are spatially correlated. The paper suggested that the procedure proposed in Conley (1999) will offer improved but imperfect inference, provided a large bandwidth is specified. Similarly, Ferman (2019) argued that commonly used datasets such as the Current Population Survey (CPS) and the American Community Survey (ACS) exhibit spatial correlation. The paper suggested that multi-way clustering, with one dimension being the cross section, can correct many issues of spatial correlation. However, it is not robust to situations in which there is correlation of errors within both different time periods and different groups. The best approach to dealing with unknown spatial correlation in addition to “conventional” clustering is something that clearly warrants further study.

3.3 Design-based and sample-based uncertainty

Up to this point, we have implicitly assumed that any sample we may have is infinitesimally small relative to the population. This evidently makes sense when the population is actually very large. But suppose our sample constitutes a substantial fraction of the population in which we are interested. Should our standard errors take account of this fact? If so, it would seem that they must tend to zero as the sample size tends to the size of the population. But cluster-robust standard errors like those based on (4) do not have this property.

Abadie, Athey, Imbens, and Wooldridge (2019) explores this setting in the absence of clustering. It argues that many samples constitute substantial fractions of the populations of interest and that standard errors should take this into account. However, the paper points out that, even when the sample consists of the entire population, there may still be uncertainty that needs to be accounted for. In particular, if treatment were randomly assigned to some observations but not to others, outcomes would be stochastic because of the random assignment. This is called design-based uncertainty, because it depends on the experimental design. Under the assumption that observations are independent, the paper

provides methods for making valid inferences. These generally yield narrower confidence intervals than conventional (heteroskedasticity-robust) procedures, especially when the sample is large relative to the population.

Abadie, Athey, Imbens, and Wooldridge (2017) extends the design-based approach of Abadie et al. (2019) to the case in which disturbances may be correlated within clusters, as often happens with “natural” experiments that are analyzed using difference-in-differences models. As we discussed above, it can be important to take clustering into account when there is a high degree of correlation in an explanatory variable within clusters, as measured by ρ_x in expression (11). In the DiD setting, some states or groups are treated, at least for some observations, while others are not. Thus ρ_x can be quite high, and using a CRVE can greatly improve inference (but see Section 5, especially Subsection 5.3). These situations also occur frequently in the settings of lab and field experiments. Abadie et al. (2017) argue that cluster-robust inference is necessitated by the design-based uncertainty in these settings. Specifically, they argue that, when the assignment to treatment is correlated within clusters, then cluster-robust standard errors are required.

Cluster-robust inference is also needed whenever there is a clustered sampling design. In many cases, the population contains a large number of groups of observations, and only some of those groups are included in the sample. For example, many education samples contain observations from all students within some schools but no students from other schools. Clustered sampling also often occurs within larger surveys, such as the CPS, where sampling is typically done by census tract. For each state, all of the households in a given sample may be drawn from a relatively modest proportion of the census tracts within the state. For this sort of survey, the survey design is often quite complex, and it can create both heteroskedasticity and within-state correlation; see Kolenikov (2010), among others.

Abadie et al. (2017) suggested that clustering may be too conservative in settings where there is neither cluster-specific treatment assignment nor a clustered sampling design. In other words, it may be too conservative when there is neither design-based nor sample-based uncertainty. In the former case, ρ_x of the treatment variable will be close to zero. In the latter case, we might have a nationally representative dataset that, unlike the CPS, does not involve sampling by census tract or other geographical subclusters.

The arguments in Abadie et al. (2019, 2017) are interesting and provocative. However, they depend critically on the assumption that the sample is large relative to a finite population in which the investigator is interested. We believe that, in many cases, economists are implicitly interested not in an actual population but in a meta-population from which they imagine the former to have been drawn. For example, if we have data for 50 U.S. states, with no sampling and no experimental design involved, then we can either view those data as non-random quantities, or we can view them as 50 draws from a meta-population of states. If we take the former view, then all we can do is to report some numbers that characterize the population. But if we take the latter view, then we can perform statistical inference in the usual way.

Although the idea of a meta-population may seem odd, economists implicitly make use of it whenever they analyze time-series data (and many other types of data). For instance, if history gives us aggregate inflation data for some country from 1969 to 2018, then we cannot obtain another dataset for the same time period by drawing a new sample. However, we can imagine that the data were generated by some sort of data-generating process that

characterizes the meta-population of inflation rates and other macro variables, and we can attempt to estimate key features of that DGP.

What [Abadie et al. \(2017\)](#) calls the “model-based” approach implicitly involves the idea of a meta-population. The DGP is simply a model like (1) accompanied by a way of obtaining the matrix \mathbf{X} and the vector \mathbf{u} from random submatrices \mathbf{X}_g and subvectors \mathbf{u}_g associated with clusters chosen at random from a meta-population of clusters. In general, cluster-robust inference is valid within a meta-population framework, and the finite-population arguments of [Abadie et al. \(2017\)](#) do not apply.

4 How to cluster

In order to obtain reliable cluster-robust inferences, the most important decision to be made in many cases is how to divide the sample into clusters. In this section, we attempt to provide some guidance. We take the model-based approach and assume that some level of clustering is appropriate (but see Subsections 3.2 and 3.3). We consider several guiding principles for determining the level of clustering.

The key assumption for cluster-robust inference to be valid is that the error terms are arbitrarily correlated within clusters but uncorrelated across clusters. It is important to specify the level of clustering in such a way that this assumption is likely to be true, or at least to provide a good approximation.

4.1 Cluster at the broadest or most aggregate level

When there is more than one level at which to cluster, and the levels are nested, then one should generally cluster at the broadest feasible level ([Cameron and Miller 2015](#)). Suppose, for example, that we can cluster either by city or by state. Clustering by state captures all the within-city correlation, and it also allows for the disturbances to be correlated within states but across cities. In contrast, clustering by city assumes that all of the correlations across cities but within states are zero. If that assumption is false, then standard errors are very likely to be too small.

Of course, there is a downside to clustering at too coarse a level. The smaller is the number of clusters, for a given sample size, the larger is the number of elements of $\mathbf{\Omega}$ that implicitly have to be estimated. If we cluster too coarsely, many of these elements are actually zero, and trying to estimate a large number of zeros inevitably makes the CRVE noisier. Even in the ideal case in which the $t(G - 1)$ distribution provides a good approximation, making G smaller will cause test power to fall and confidence intervals to become wider, simply because the critical values for the t distribution increase as $G - 1$ becomes smaller. However, unless clustering at the highest feasible level means using a value of G that is very small, the loss of power from clustering at that level is likely to be modest in comparison with the severe size distortions that can occur from clustering at too low a level; see Section 7. This becomes more true as the sample size becomes larger.

Another problem with clustering at too coarse a level is that, when G is small, the $t(G - 1)$ distribution often does not provide a good approximation. But, except in the most extreme

cases, this problem can generally be overcome by using bootstrap methods; see Section 5 for a discussion of these methods and Section 7 for simulation evidence.

4.2 Cluster at least at the level of a policy change

When the null hypothesis of interest involves a treatment variable resulting from a policy change, one should always cluster at a level no lower than the one to which the policy was applied. As mentioned in Subsection 3.3, Abadie et al. (2017) suggest that clustering is necessary whenever treatment is assigned at the cluster level. From this perspective, one would want to match the level of clustering done in the analysis to the level at which treatment was assigned. For instance, in a randomized control trial, if treatment assignment were done at the village level, then one would want to cluster at the village level. However, based on the arguments of Subsection 4.1, one might choose to cluster at a still higher level, especially if the number of villages were large and they naturally fell into a reasonable number of larger groups. Ideally, both approaches would yield similar results.

Whether or not there is a policy change, it always makes sense to cluster at a level no lower than the one at which observations were included in the sample. For example, if classrooms were chosen at random for inclusion in the sample, then one would want to cluster at either the classroom level or the school level. But if schools were chosen at random, then one would want to cluster at the school level. In both cases, of course, one might choose to cluster at an ever higher level, such as school districts.

4.3 Cluster at the cross-section level for panel data

When working with panels and repeated cross sections, it is important never to cluster below the cross-section level. As shown first in Bertrand et al. (2004), clustering at the level of the cross section allows for arbitrary autocorrelation of the error terms within cross-sectional units. In many contexts, this means that clustering at the state level will result in much more reliable inference than, say, clustering at the state \times year level.

An even more general approach for this sort of data would be to use two-way clustering by cross-sectional unit and time. MacKinnon (2019) provides evidence that this seems to be appropriate in the context of an earnings equation using CPS data that includes both state and year fixed effects.

4.4 There is no golden number of clusters

Early simulation results such as those in Bertrand et al. (2004) and Cameron et al. (2008) concerned models with balanced clusters. This gave a false sense of how well cluster-robust variance estimators perform in finite samples. A rule of thumb emerged that $G \geq 50$ would allow for reliable inference, which was changed (in jest) to $G \geq 42$ in Angrist and Pischke (2008). However, any such rule of thumb can be extremely misleading, because all CRVEs tend to become less reliable as the clusters become more unbalanced; see Carter et al. (2017), MacKinnon and Webb (2017a), and Djogbenou et al. (2019). The problems associated with unbalanced clusters are discussed in Section 5.2.

It is far more important to get the level of clustering right, as we have discussed in Subsections 4.1, 4.2, and 4.3, than it is to ensure that G is large enough for t -statistics to have their namesake distribution. When we cluster at too fine a level, standard errors will typically be too small by a factor that increases with the sample size. In contrast, when we cluster at the right level, inference based on the $t(G - 1)$ distribution may be seriously unreliable, but other methods of inference (notably the bootstrap methods discussed in Subsection 5.1) often provide quite reliable inferences.

4.5 In many settings, over-clustering is mostly harmless

In a model-based context, over-clustering (within reason) tends to be relatively harmless, except in one important special case (Subsection 5.3). By over-clustering, we mean either clustering at a coarser level than is actually appropriate or clustering in two dimensions when just one is needed. Simulation results suggest that, in most cases, a moderate amount of over-clustering should have little impact on size (provided the wild cluster bootstrap is used) but some impact on power; see Section 7.

Of course, there are limits to the amount of over-clustering that can be handled safely, even when using the wild cluster bootstrap (Subsection 5.1). Although bootstrap P values are often very reliable even when G is quite small (for example, they work remarkably well in the simulations of Section 7 when $G = 10$ and cluster sizes are quite unbalanced), there are extreme cases, discussed in the next section, where they cannot be relied upon.

4.6 Tests for the level of clustering

Choosing the appropriate level of clustering can be difficult. It would be desirable if this choice could be aided by the use of formal statistical tests. Ibragimov and Müller (2016) provided a procedure for testing a null hypothesis of fine clustering against an alternative of coarse clustering. For example, it can test a null of heteroskedasticity against an alternative of city-level clustering, or a null of city-level clustering against an alternative of state-level clustering. However, this test cannot be used when the regressor of interest is invariant within a cluster, such as in a difference-in-differences model. It also cannot be used to test a null of one-way clustering versus an alternative of multi-way clustering. Work in progress by the authors of this chapter and Morten Nielsen seeks to develop a test that can handle these situations.

5 What can go wrong

Although both CV_1 and CV_2 , given in expressions (4) and (5), estimate the true covariance matrix (3) consistently under moderately weak conditions (Djogbenou et al. 2019), they do not always provide reliable estimates, even when the sample size is very large. This reflects the fact that all CRVEs differ fundamentally from most other covariance matrix estimators in one important respect. The latter usually converge to the true value as the number of observations, N , tends to infinity. But a CRVE converges to the true value as the number of clusters, G , tends to infinity. Therefore, no matter how large N may be, inference

based on cluster-robust standard errors (for the correct set of clusters) can sometimes be problematical, even seriously misleading, when G is not large.

In this section, we assume that the clusters have been chosen correctly, with no correlation of disturbances across clusters. Nevertheless, there are three situations in which cluster-robust inference may be unreliable. The first is when the number of clusters is small. The second is when cluster sizes, or other features of the clusters, are unbalanced. The third, which can be thought of as a special case of the second, is when the model focuses on the effects of a treatment dummy, and few clusters are “treated.” In the first two cases, we recommend using a particular bootstrap method, which we describe in Subsection 5.1. This method can also work well in the third case, but it can sometimes fail disastrously. When it does, randomization inference or an alternative bootstrap procedure may be able to provide reliable inferences.

5.1 Few clusters

When the number of clusters is reasonably large (say, a few hundred), and each cluster provides roughly the same amount of information, then inference based on cluster-robust standard errors and the $t(G - 1)$ distribution is likely to be very reliable. There are also cases in which this type of inference works well even when G is quite small (Bester, Conley, and Hansen 2011), but it would usually be unwise to rely on it.

Because the middle factor in any CRVE is a sum over G matrices, each with rank one, it should be obvious that a CRVE may not provide reliable inferences when G is small. Ideally, there would be more clusters than parameters, so that the CRVE could potentially have full rank. This is more important for Wald tests of several restrictions than for t -tests of just one restriction. But it makes sense that G should need to be larger for reliable inference when K (the number of regression coefficients) is large than when it is small.³

Carter et al. (2017) proposed the concept of an “effective number of clusters” G^* and provided a way to compute an approximation to it; a Stata package is discussed in Lee and Steigerwald (2018). When G^* is not much less than G , and G itself is not too small (say, 50 or more), then inference based on a cluster-robust t -statistic and the $t(G - 1)$ distribution generally works well. In contrast, when G^* is much smaller than G , that type of inference can be very unreliable.

A better approach, in our view, is to rely on bootstrap tests and bootstrap confidence intervals. The basic idea of bootstrap testing is to compare a test statistic with the empirical distribution of a large number of bootstrap test statistics computed from simulated samples. Conceptually, a bootstrap confidence interval is then a set of parameter values for which a bootstrap test does not reject. Accessible introductions to bootstrap methods include MacKinnon (2002), Davidson and MacKinnon (2006), and Horowitz (2019).

Suppose that we wish to test the restriction $\mathbf{a}'\boldsymbol{\beta} = \mathbf{0}$, where \mathbf{a} is a known vector. For example, if we are testing the hypothesis that $\beta_k = 0$, the k^{th} element of \mathbf{a} would be 1 and all the rest would be 0. We first calculate a cluster-robust t -statistic, say t_a , for the hypothesis

³To our knowledge, there have been no simulation studies that focus on the relationship between G and K for possibly large K . However, simulations in Appendix C.2 of Djogbenou et al. (2019) suggest that adding either 4 or 8 additional regressors when $G = 25$ makes tests based on the $t(24)$ distribution noticeably more prone to over-reject.

that $\mathbf{a}'\boldsymbol{\beta} = \mathbf{0}$. We then generate a large number (B) of bootstrap samples indexed by b and use each of them to compute a bootstrap t -statistic t_a^{*b} . Sensible values of B are numbers like 999 and 9999 (Davidson and MacKinnon 2000). Then a symmetric two-tailed test is based on the bootstrap P value

$$p^*(t_a) = \frac{1}{B} \sum_{b=1}^B \mathbb{I}(|t_a^{*b}| > |t_a|), \quad (13)$$

where $\mathbb{I}(\cdot)$ is the indicator function. The more extreme is $|t_a|$, the fewer of the $|t_a^{*b}|$ should exceed it by chance. If, for example, 17 out of 999 do so, then $p^*(t_a) = 0.017$, and we can confidently reject the null hypothesis at the .05 level.

In the context of the model (1), we want to calculate t_a^{*b} from the b^{th} bootstrap sample in precisely the same way as we calculated t_a from the actual sample. How well bootstrap tests perform then depends on how the bootstrap samples used to calculate the t_a^{*b} are generated. In principle, there are many ways to do so. However, both theory and simulation evidence currently favor a particular method, namely, the restricted wild cluster bootstrap that uses the Rademacher distribution.

The wild cluster bootstrap was proposed in Cameron et al. (2008) and studied extensively in MacKinnon and Webb (2017a). Its asymptotic validity was proved in Djogbenou et al. (2019). The equation used to generate the b^{th} bootstrap sample for the restricted wild cluster (WCR) bootstrap is

$$\mathbf{y}_g^{*b} = \mathbf{X}_g \tilde{\boldsymbol{\beta}} + \mathbf{u}_g^{*b} = \mathbf{X}_g \tilde{\boldsymbol{\beta}} + \tilde{\mathbf{u}}_g v_g^{*b}, \quad g = 1, \dots, G, \quad (14)$$

where quantities with g subscripts are scalars, subvectors, or submatrices associated with the g^{th} cluster. In (14), $\tilde{\boldsymbol{\beta}}$ is the vector of least squares estimates of $\boldsymbol{\beta}$ subject to the restriction or restrictions to be tested, and $\tilde{\mathbf{u}}_g$ is the restricted residual vector for cluster g . Finally, v_g^{*b} is a scalar auxiliary random variable that follows the Rademacher distribution, which takes values 1 and -1 , each with probability $1/2$.

The key idea of the WCR bootstrap, or WCRB, is that the bootstrap disturbances \mathbf{u}_g^{*b} for cluster g are generated by multiplying the residual subvector $\tilde{\mathbf{u}}_g$ by the scalar random variable v_g^{*b} . This ensures that, asymptotically, the disturbances for cluster g in the bootstrap DGP (14) have (on average) covariance matrix $\boldsymbol{\Omega}_g$. In consequence, estimates based on the bootstrap sample \mathbf{y}^{*b} have the same asymptotic distribution as ones based on the actual sample \mathbf{y} , assuming the restrictions are true. This implies that, if we reject the null hypothesis whenever the bootstrap P value in (13) is less than α , and $\alpha(B+1)$ is an integer, the asymptotic level of the test is α .

Nothing in the above arguments implies that a WCRB test will always perform better in finite samples than a test based on the $t(G-1)$ distribution. However, higher-order theory in Djogbenou et al. (2019) does strongly suggest that this is likely to be the case. It also suggests that the Rademacher distribution (Davidson and Flachaire 2008) is usually the best choice for the auxiliary distribution and that failing to impose the null hypothesis on the bootstrap DGP is a bad idea. In all cases, simulation results support these implications of the theory. Canay, Santos, and Shaikh (2018) studied cases in which the WCRB yields exact results for large samples even when the number of clusters is fixed. Their results impose conditions on the distribution of the covariates and require the use of the Rademacher distribution.

Equation (14) suggests that we need to generate B bootstrap samples of size N and compute a bootstrap t -statistic t_a^{*b} for each of them. This can be computationally challenging when N is large, especially if K is also large. Luckily, there is a way to reduce the computational burden dramatically, especially when G is small. Since the bootstrap DGP (14) satisfies the restrictions, the numerator of the bootstrap t -statistic for $\mathbf{a}'\boldsymbol{\beta} = \mathbf{0}$ is

$$\mathbf{a}'\hat{\boldsymbol{\beta}}^{*b} = \mathbf{a}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u}^{*b} = \sum_{g=1}^G \mathbf{a}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'_g\tilde{\mathbf{u}}_g v_g^{*b}. \quad (15)$$

The quantities $\mathbf{a}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'_g\tilde{\mathbf{u}}_g$ are scalars that can be calculated after the restricted model has been estimated but before bootstrapping begins. The rightmost expression in (15) is then just the sum over the G clusters of those scalars times the realized random variables v_g^{*b} . This expression can be used to calculate $\mathbf{a}'\hat{\boldsymbol{\beta}}^{*b}$ extremely rapidly.

A similar, but much trickier, procedure can be used to calculate the denominator of the bootstrap t -statistic efficiently. The computations for each bootstrap sample are now $O(G^2)$ instead of $O(G)$. However, for moderate and even quite large values of G , this is still very much less expensive than generating a bootstrap sample according to (14) and computing $\mathbf{a}'\hat{\boldsymbol{\beta}}^{*b}$ and its cluster-robust standard error in the usual way. The Stata package `boottest` (Roodman, MacKinnon, Nielsen, and Webb 2019) implements this efficient algorithm. Unless N is extremely large (so that the preliminary calculations are time-consuming) or G is greater than several hundred (in which case bootstrapping is probably unnecessary), it is usually very inexpensive to compute the bootstrap P value (13), even for $B = 9999$.⁴

Because of this, we recommend using the WCRB almost all the time. For sufficiently large values of N , this will actually be cheaper than many alternative procedures, such as computing CV_2 . When the bootstrap test yields essentially the same conclusion as a test based on the $t(G - 1)$ distribution, then we can usually be confident that both results are reasonably reliable. On the other hand, when the bootstrap test provides weaker evidence against the null hypothesis than the asymptotic test (alas, it very rarely provides stronger evidence), then we should usually disregard the latter (but see Subsection 5.3). It is difficult to say whether we should rely on the results of the bootstrap test when they differ sharply from those of the asymptotic test. This will depend on how well the WCRB is known to perform in similar circumstances.

The WCRB can also be used to form confidence intervals by “inverting” the bootstrap test, and `boottest` does this by default for tests of a single restriction. For details, see Roodman et al. (2019, Section 3.5). These bootstrap confidence intervals can be much more accurate than conventional ones based on cluster-robust standard errors and the $t(G - 1)$ distribution (MacKinnon 2015).

The WCRB often works surprisingly well even for quite small values of G . One problem, however, is that the number of distinct bootstrap samples with the Rademacher distribution (or any other two-point distribution) is just 2^G . When $G < 10$, this may be too small to estimate p^* reliably. Webb (2014) therefore proposed a six-point distribution which largely solves this problem, because $6^G \gg 2^G$. When 2^G is reasonably large but smaller than the

⁴At time of writing, the R package `clusterSEs` (Esarey 2018) implements the WCRB but does not employ the computational tricks used by `boottest`.

chosen value of B , it is better to enumerate all possible bootstrap samples than to draw them at random, and `boottest` does this by default.

Some situations in which inference is particularly difficult are discussed in the next two subsections. Even the WCRB can be unreliable, especially in the case dealt with in Subsection 5.3. When there is doubt about its reliability, it would be wise to confirm its results using other methods. In particular, [Imbens and Kolesár \(2016\)](#) suggested a procedure based on CV_2 that computes a number smaller than $G - 1$ to be used as the degrees of freedom for the t distribution. This procedure is computationally burdensome when N is large ([MacKinnon and Webb 2018](#)). A related procedure that is based on CV_1 and is computationally feasible even for very large samples was suggested by [Young \(2016\)](#). Limited simulation evidence suggests that the Imbens-Kolesár and Young procedures do not, in general, perform as well as the WCRB, but they often perform quite well, and they can yield different results in some cases. It is probably safe to accept the inferences from the WCRB when they agree with those from these alternative procedures.

All of the procedures we have discussed are based on OLS estimation of (1) using the entire sample. The procedure proposed in [Ibragimov and Müller \(2010\)](#) is based on estimating the model on a cluster-by-cluster basis, but this is only feasible if no regressors are invariant within clusters. The procedure in [Ibragimov and Müller \(2016\)](#) overcomes these problems by partitioning the sample into groups of clusters. One could instead solve the problem of invariant regressors by combining the original clusters into fewer and larger ones. For reasons of space, however, we do not discuss these methods further.

5.2 Unbalanced clusters

The asymptotic validity of inference based on both t -statistics and the wild cluster bootstrap depends on the properties of the score vectors $\mathbf{s}_g = \mathbf{X}'_g \mathbf{u}_g$ ([Djogbenou et al. 2019](#)). Ideally, all of them would follow the same multivariate distribution, with covariance matrices $\mathbf{X}'_g \boldsymbol{\Omega}_g \mathbf{X}_g$ the same for all g . Nevertheless, the theory allows for considerable heterogeneity. But when the score vectors are too heterogeneous across clusters, the conditions for the bootstrap and asymptotic distributions to coincide asymptotically are no longer satisfied. This suggests that inference will become less reliable as the data become more heterogeneous across clusters.

When cluster sizes vary too much, standard theoretical results may not apply. Some of the simulations in [Djogbenou et al. \(2019\)](#) have one cluster that contains half the observations. In this extreme case, tests based on the $t(G - 1)$ distribution over-reject more severely, not less, as G increases. WCRB tests also over-reject more severely as G increases, but to a much lesser extent. In less extreme cases, where the single large cluster becomes a smaller proportion of the sample as G increases, the performance of WCRB tests always improves with G . Once the large cluster becomes small enough (on the order of 20% of the sample in these experiments), the bootstrap tests work very well. But this can require quite a large number of clusters, perhaps on the order of several hundred.

Variation in cluster sizes is not the only sort of heterogeneity that is likely to cause cluster-robust tests to be misleading. Even if all clusters are roughly the same size, the covariance matrices $\mathbf{X}'_g \boldsymbol{\Omega}_g \mathbf{X}_g$ can vary across them for other reasons. Perhaps a few clusters contain a lot more information than the remaining ones, or perhaps the disturbances are heteroskedastic across clusters. In both cases, we would expect all methods to make more

serious inferential errors than they would if the model had the same number of homogeneous clusters. However, simulation evidence always suggests that the WCRB is less affected by heterogeneity than $t(G - 1)$ tests. This reinforces our earlier recommendation to employ the WCRB almost all the time.

Although we cannot directly observe the $\mathbf{X}'_g \boldsymbol{\Omega}_g \mathbf{X}_g$ matrices, we can observe the N_g and the $\mathbf{X}'_g \mathbf{X}_g$ matrices, which provide measures of how much information each cluster contains. The effective number of clusters G^* (Carter et al. 2017) is quite sensitive to heterogeneity across the $\mathbf{X}'_g \mathbf{X}_g$. Thus finding that $G^* \ll G$ provides a useful warning. However, G^* is not sensitive to heteroskedasticity across clusters. To see whether that is a problem, we can calculate the variance of the residuals for each cluster separately.

Because the finite-sample properties of all inferential methods, including the WCRB, depend in complicated ways on the model and dataset, it is difficult to provide a rule of thumb for when it is safe to rely on WCRB P values. These are most likely to be unreliable (often too small, but sometimes too large, as we will see in Subsection 5.3) when G is small, when cluster sizes vary a lot, when the $\mathbf{X}'_g \mathbf{X}_g$ vary a lot (which is likely to cause $G^* \ll G$), and/or when the variance of the residuals differs sharply across clusters. In such cases, as we recommended above, it is important to employ other methods as well. Of course, we are not recommending the use of multiple inferential methods as a fishing expedition, but rather as a way of verifying (or casting doubt on) the validity of the WCRB.

5.3 Few Treated Clusters

Many applications of cluster-robust inference involve treatment effects estimated at the cluster level. In what we call the pure treatment case, some schools or villages or experimental subjects are treated, and others are not. Thus every observation in every treated cluster is treated. In contrast, for difference-in-differences (DiD) models, some jurisdictions are never treated, and others are treated during some, but not all, time periods. Unfortunately, when the number of treated clusters is small, cluster-robust standard errors can be very much too small, even if the total number of clusters is large. Because this situation is commonly encountered, it is worth discussing the issues associated with few treated (or few control) clusters in some detail.

Following MacKinnon and Webb (2017a), consider the pure treatment model

$$y_{gi} = \beta_1 + \beta_2 d_{gi} + u_{gi}, \quad g = 1, \dots, G, \quad i = 1, \dots, N_g, \quad (16)$$

where d_{gi} equals 1 for the first G_1 clusters and 0 for the remaining $G_0 = G - G_1$ clusters. Every observation in the g^{th} cluster is either treated ($d_{gi} = 1$) or not treated ($d_{gi} = 0$). The analysis would be more complicated if we included additional regressors, but it would not change in any fundamental way.

The key problem is that, because the dummy variable d_{gi} must be orthogonal to the residuals, the latter must sum to zero over all the treated observations. This has some unfortunate implications. In the extreme case in which only cluster 1 is treated, $\mathbf{d}'_1 \hat{\mathbf{u}}_1 = 0$, where \mathbf{d}_1 is the subvector of the dummy variable for cluster 1. In other words, the residuals for cluster 1 must sum to zero. This implies that $\mathbf{d}'_1 \hat{\mathbf{u}}_1 \hat{\mathbf{u}}'_1 \mathbf{d}_1 = 0$. But that quantity is supposed to estimate the element $\mathbf{d}'_1 \boldsymbol{\Omega}_1 \mathbf{d}_1$ of the matrix $\mathbf{X}'_1 \boldsymbol{\Omega}_1 \mathbf{X}_1$, and that element has to be estimated accurately if we are to obtain a reliable standard error for $\hat{\beta}_2$, because

most of the information about β_2 is coming from cluster 1, the only treated cluster. Since $\mathbf{d}'_1 \hat{\mathbf{u}}_1 \hat{\mathbf{u}}'_1 \mathbf{d}_1 = 0$, the cluster-robust standard error for $\hat{\beta}_2$ is very much too small.⁵ This causes the cluster-robust t -statistic to be very much too large. When $G_1 = 1$, it would not be unusual for the t -statistic to be five times as large as it should be.

When there are two treated clusters, $\mathbf{d}'_1 \hat{\mathbf{u}}_1 + \mathbf{d}'_2 \hat{\mathbf{u}}_2 = 0$. This implies that $\mathbf{d}'_1 \hat{\mathbf{u}}_1 \hat{\mathbf{u}}'_1 \mathbf{d}_1$ and $\mathbf{d}'_2 \hat{\mathbf{u}}_2 \hat{\mathbf{u}}'_2 \mathbf{d}_2$, although both non-zero, underestimate $\mathbf{d}'_1 \boldsymbol{\Omega}_1 \mathbf{d}_1$ and $\mathbf{d}'_2 \boldsymbol{\Omega}_2 \mathbf{d}_2$ severely. The problem diminishes as the number of treated clusters increases; see [MacKinnon and Webb \(2017a, Appendix A.3\)](#). Although this argument is for the pure treatment model, essentially the same argument applies to DiD models. Whenever only cluster 1 is treated, $\mathbf{d}'_1 \hat{\mathbf{u}}_1 = 0$, because the residuals for the treated observations sum to zero, and the residuals for the untreated observations are multiplied by elements of \mathbf{d}_1 that equal 0.

Unfortunately, bootstrapping does not solve the problem. As [MacKinnon and Webb \(2017a, Section 6\)](#) explains, the WCRB always under-rejects very severely when $G_1 = 1$. In simulation experiments with 400,000 replications ([MacKinnon and Webb 2017b](#)), there are often no rejections at all for tests at the .05 level. There is also typically very severe under-rejection for $G_1 = 2$. The bootstrap fails in this case because the absolute values of the actual and bootstrap test statistics are strongly positively correlated. When $|t_a|$ is large, the $|t_a^{*b}|$ tend to be large as well, so that the bootstrap P value in [\(13\)](#) is not likely to be small. The extent of the under-rejection depends on G , G_1 , the cluster sizes, and the numbers of treated observations within the treated clusters. For given values of G and G_1 , the problem tends to be most severe when the number of treated observations is small.

In order to avoid this problem, [MacKinnon and Webb \(2018\)](#) suggested using the ordinary wild bootstrap, which uses one auxiliary random variable per observation instead of one per cluster. Surprisingly, even though the distribution of the $\hat{\beta}^{*b}$ does not coincide asymptotically with the distribution of $\hat{\beta}$, tests based on the ordinary wild bootstrap are asymptotically valid ([Djogbenou et al. 2019](#)). Moreover, in some circumstances, these tests can perform very well, even when $G_1 \leq 2$. In the pure treatment case, the key requirement is that all clusters be the same size. However, there are also many cases in which ordinary wild bootstrap tests (based on cluster-robust t -statistics) either over-reject or under-reject systematically. In our view, these tests are worth trying when G_1 is very small and the WCRB does not reject. However, one should not rely on their results unless simulation evidence suggests that they perform well for the case at hand. Empiricists may want to conduct their own simulations to assess the validity of inference procedures given the structure of their data.

A very different approach, which can work extraordinarily well in some cases even when G_1 is very small, is to employ randomization inference, or RI. The basic idea of RI is to compare an actual parameter estimate (or test statistic) with a set of hypothetical estimates obtained by re-randomizing, that is, pretending that control clusters were actually treated. In the context of a regression model like [\(16\)](#), the values of the dependent variable do not change across re-randomizations, but the values of the treatment dummy do change. There is a large literature on RI, dating back to [Fisher \(1935\)](#). For modern treatments, see [Lehmann and Romano \(2008\)](#) and [Imbens and Rubin \(2015\)](#).

There is more than one way to use RI in the context of treatment models. The simplest

⁵Here we mean the standard error based on CV_1 . In this case, the CV_2 covariance matrix cannot be computed ([MacKinnon and Webb 2018](#)).

approach is to consider all possible assignments of treatment to clusters in a model like (16). For example, if there are 15 clusters of which 2 are treated, then there are $(15 \cdot 14)/2 = 105$ ways in which treatment could have been assigned. One of them corresponds to the actual treatment, and the other 104 correspond to re-randomizations. If the estimate $\hat{\beta}_2$ is sufficiently extreme compared with the 104 estimates associated with the re-randomizations, then it seems reasonable to reject the null hypothesis of no treatment effect.

An RI procedure like the one just outlined (but not identical to it) was suggested by Conley and Taber (2011). It is based on parameter estimates. MacKinnon and Webb (2020) pointed out that a similar procedure could be based on cluster-robust t -statistics and studied the properties of both procedures. When G is sufficiently large, all clusters are identical, and treatment is assigned at random, both procedures work extremely well, but the one based on coefficient estimates has more power. However, when clusters are not identical and the investigator knows which ones were treated, both procedures can either over-reject or under-reject. The one based on t -statistics generally performs better for $G_1 > 1$, especially when the treated clusters are larger or smaller than the controls. MacKinnon and Webb (2019) studied a bootstrap variant of RI that can be used when RI yields imprecise results because the number of possible randomizations is small.

Alternative RI procedures, which are too complicated to be discussed here, have been proposed in Canay, Romano, and Shaikh (2017) and Hagemann (2019a, b). The procedures developed in the two Hagemann papers may be particularly attractive because they can both be used even when G is very small, and the latter can be used even when there is substantial heterogeneity across clusters. A Stata package for performing randomization inference is described in Hess (2017).

Other procedures have been suggested to deal with the problems of inference for few treated clusters. Alternative CRVEs (Bell and McCaffrey 2002; Imbens and Kolesár 2016; Young 2016), combined with ways of estimating degrees of freedom, work better than simply using CV_1 ; simulation results for these procedures can be found in the appendix of MacKinnon and Webb (2018). Donald and Lang (2007) proposed collapsing the data into pre-treatment and post-treatment means and comparing averages. The procedure of Ibragimov and Müller (2016) allows for inference so long as there are at least four treated groups. Ferman and Pinto (2019) suggested a DiD procedure that works for few treated and many control groups when there is heteroskedasticity of known form across clusters.

6 Empirical Examples

We now present two empirical examples that highlight some of the issues discussed above, such as having few clusters, having few treated or control clusters, and/or misspecifying the level of clustering. We do not present these examples either to confirm or overturn the original results, but rather to highlight the consequences of these common situations. All of the analysis was conducted using replication files in Stata and employing the `boottest` package discussed in Roodman et al. (2019).

6.1 Experimental Evidence on Female Voting Behavior

Our first example illustrates the effect of unbalanced clusters. [Giné and Mansuri \(2018\)](#) examines the effect on women’s voting behavior of informing women about the voting process and the importance of voting. The regression model that we examine is

$$Y_i = \beta_1 T_{1,T} + \beta_2 T_{2,T} + \beta_3 T_{1,U} + \beta_4 T_{2,U} + \mathbf{X}\gamma + v + \epsilon_i.$$

Here $T_{i,T}$ is an indicator variable for individuals in the treated region who are in treatment group i , and $T_{i,U}$ is similarly defined for the untreated individuals. The parameters of interest are β_1 through β_4 . β_1 and β_2 capture the direct treatment effect of each treatment, while β_3 and β_4 capture the spillover effects of each treatment.

The authors clustered by neighborhood, but we also cluster by village. When clustering by neighborhood, there are 57 actual clusters, but only between 14.2 and 15.7 effective clusters, depending on which coefficient G^* is calculated for. Thus, at the neighborhood level, the clusters are apparently quite unbalanced. In contrast, when clustering by village, there are 9 actual clusters and between 8.6 and 8.8 effective clusters, so that the clusters are well balanced. Only 7 of the 57 neighborhoods are controls, which must be partly responsible for the small value of G^* when clustering by neighborhood.

Table 1: Comparison of Reported and Bootstrapped P values

Variable	Coefficient	Neighborhood		Village	
		$t(56)$	WCRB	$t(8)$	WCRB
$T_{1,T}$	$\hat{\beta}_1 = -0.092$	0.024	0.0723	0.049	0.0610
$T_{2,T}$	$\hat{\beta}_2 = -0.119$	0.004	0.0203	0.025	0.0171
$T_{1,U}$	$\hat{\beta}_3 = -0.141$	0.003	0.0103	0.024	0.0238
$T_{2,U}$	$\hat{\beta}_4 = -0.112$	0.043	0.0664	0.131	0.1550

These results are for the analysis in Panel A of Table 9 of [Giné and Mansuri \(2018\)](#). Columns 3 and 5 report P values based on the $t(G - 1)$ distribution, and columns 4 and 6 report WCRB P values calculated using 99,999 replications. There are 2637 observations, 57 neighborhood clusters, and 9 village clusters. The bootstrap DGPs use Rademacher weights when clustering by neighborhood and Webb (6-point) weights when clustering by village.

Table 1 shows the results for both levels of clustering. Not surprisingly, the $t(G - 1)$ P values always increase when we move from neighborhood-level to village-level clustering. Because there are only 9 village clusters, and because the effective number of clusters at the neighborhood level is so much smaller than the actual number, these P values are probably not reliable. WCRB P values are always larger, often much larger, than the ones based on the $t(G - 1)$ distribution. Interestingly, WCRB P values based on village-level clustering are sometimes smaller and sometimes larger than the ones based on neighborhood-level clustering. Whatever the level of clustering, $\hat{\beta}_2$ and $\hat{\beta}_3$ appear to be significant at the .05 level, but the other two coefficients do not.

6.2 Private Equity Management Practices

Our second example illustrates the problem of few treated clusters. [Bloom, Sadun, and Van Reenen \(2012\)](#) studies the impact of private-equity ownership of a firm on its management practices. The dependent variable, `AllOps`, is a normalized score of the extent to which the firm’s manager is responsible for operations, as assessed in an interview. The regressor `PEOwned` is a dummy for private-equity ownership. The model we estimate is

$$\text{AllOps}_i = \beta_0 + \beta_1 \text{PEOwned}_i + \beta_2 \mathbf{X}_i + \epsilon_i, \quad (17)$$

where the vector \mathbf{X}_i contains country and industry fixed effects, as well as various firm controls, such as the number of employees and the age of the firm, along with several variables related to the interview used to determine the value of `AllOps`. The authors clustered the standard errors at the firm level, which is almost the same as not clustering at all, because there is just one observation for most firms.

Table 2: Conventional and Bootstrap P values and Confidence Intervals

Cluster	$t(G - 1)$	WCRB	WCUB	WRB	WUB
Firm P value	0.0000	0.0000	0.0000	0.0001	0.0000
Firm C.I.	[0.077, 0.221]	[0.078, 0.220]	[0.078, 0.221]	[0.077, 0.221]	[0.078, 0.221]
Ownership P value	0.0364	0.6159	0.1651	0.5085	0.4855
Ownership C.I.	[0.012, 0.287]	[-2.71, 2.77]	[-0.054, 0.353]	[-0.527, 0.763]	[-0.362, 0.661]

These results are for the regression in Table 3 Column (1) of [Bloom, Sadun, and Van Reenen \(2012\)](#). β_1 is estimated to be 0.147. Standard errors are clustered by either firm or ownership type. There are 15,038 observations, 11,370 firms, and 10 ownership types. There are 9999 bootstraps. Webb (6-point) weights were used for WCRB and WCUB at the ownership level. Rademacher weights were used for all other cases.

When we cluster by firm, there appears to be very strong evidence against the null hypothesis that $\beta_1 = 0$ no matter whether or how we bootstrap. In the table, WCUB means the unrestricted wild cluster bootstrap (which we do not recommend), and WRB and WUB mean the ordinary wild bootstrap, restricted and unrestricted, respectively. We use them here because there are few treated clusters when clustering by ownership ([MacKinnon and Webb 2018](#)). Since all methods yield P values of zero, we also report confidence intervals. It is evident that all the confidence intervals are essentially the same, and they confirm what the P values tell us about the null hypothesis. Because there are a great many firms (11,370), the fact that only 3% of firms in the sample are owned by private equity does not cause a problem with few treated clusters when clustering is at the firm level.

We can think of (17) as a treatment model, where only firms owned by private equity are “treated.” This suggests that we should cluster at the level of the treatment, that is, ownership type. Doing so increases the conventional P value substantially, from 0.0000 to 0.0364. Because there are just 10 clusters at the ownership level, and Private Equity Ownership is the only treated cluster, the few-treated problem is extreme in this case. Therefore, the number 0.0364 should not be believed. The WCRB and WCUB tests yield much larger P values, although the discrepancy between them suggests that neither is reliable. Indeed, it

would be astonishing if the WCRB and WCUB P values agreed in this case (MacKinnon and Webb 2017a). In contrast, the WRB and WUB tests yield very similar results. Both of their P values are around 0.5, nowhere near conventional significance levels. Interestingly, the WCRB, WRB, and WUB confidence intervals differ quite a lot even though the bootstrap P values are similar.

This example illustrates a problem that occurs all too often. When we cluster at a fine level (in this case, by firm), we obtain results that look sensible but depend on a strong independence assumption. When instead we cluster at a much coarser level (in this case, by ownership), we obtain much larger standard errors, and we may encounter the problem of few treated clusters (in this case, just one).

7 Simulations

As we stressed in Section 4, choosing the correct level at which to cluster is extremely important. To investigate the issues of over-clustering and under-clustering, we now perform a set of Monte Carlo simulations in which clustering may occur at one of three levels. There are 60 zones, each with 100 observations, which are grouped into 20 cities and 10 states. Seven cities each contain one zone, six cities each contain three zones, and seven cities each contain five zones. Each state contains two cities. Thus states have between 200 and 1000 observations.

The model is

$$y_{izcs} = \beta_1 + \beta_2 x_{izcs} + u_{izcs}, \tag{18}$$

where the regressor x_{izcs} is equal to $x_s \sim N(0, 1)$ plus $x_i \sim N(0, 1)$. As the notation implies, x_s takes the same value for every observation in state s , and the x_i are independent across observations. Thus the correlation of the regressor between any pair of observations in the same state (or zone, or city) is $1/2$.

The disturbances may or may not be correlated within zones, cities, and states. For all $i, z, c,$ and s ,

$$u_{izcs} = \phi_z u_z + \phi_c u_c + \phi_s u_s + u_i,$$

where $u_i, u_z, u_c,$ and u_s are distributed as $N(0, 1)$, one of $\phi_z, \phi_c,$ and ϕ_s is equal to ϕ , and the others are equal to 0. Thus if, for example, there is clustering at the city level, u_{izcs} is equal to ϕ times a city-level random component u_c plus an individual component u_i .

In the experiments, we vary ϕ between 0 and 0.5, and we vary the level at which the disturbances are clustered. We then test the null hypothesis that $\beta_2 = 0$. Tables 3 and 4 report 5% rejection rates, as percentages, for experiments with 400,000 replications. In Table 3, the rejection rates are based on cluster-robust t -statistics and the $t(G - 1)$ distribution. In Table 4, they are based on the WCRB using the Rademacher distribution and 399 bootstrap replications.

When all correlation is within zones and we cluster at the zone level, inference based on the $t(G - 1)$ distribution is very good, and inference based on the WCRB appears to be perfect; see the first column of results in Tables 3 and 4, respectively. This is not surprising, of course, because 60 is a fairly large number of clusters, and all zones are the same size. The results become much more interesting when there is correlation at the city or state levels.

Both tables show that there is severe over-rejection whenever we under-cluster. This happens in columns 2, 3, and 6, where the DGP is marked with a “u”. The over-rejection increases with ϕ , initially very rapidly. The most severe over-rejection occurs in column 3, where the DGP has state-level clustering but standard errors are calculated at the zone level. Thus the standard errors are calculated two levels below the correct one. The over-rejection is less severe, but still very substantial, when the standard errors are calculated at the city level, only one level below the correct one.

The effects of under-clustering necessarily become worse as the sample size increases. We repeated the above experiments with 1000 observations per zone instead of 100. For large values of ϕ , the results did not change very much. For example, the rejection rate for $\phi = 0.50$ in the third column of Table 4 increased from 52.46 to 55.43. For smaller values, however, the rejection frequencies increased much more. For example, the rate for $\phi = 0.10$ in the third column of Table 4 increased from 25.69 to 48.73. Thus the consequences of ignoring small amounts of correlation by clustering at too fine a level become much more severe as the sample size increases.

The consequences of over-clustering are much less severe than those of under-clustering. Inference based on the $t(G - 1)$ distribution is somewhat unreliable when clustering by city and quite unreliable when clustering by state. In contrast, there is very little size distortion from over-clustering when the WCRB is used. This happens in columns 4, 7, and 8 of Table 4, where the DGP is marked with an “o”. Even though there are only 10 clusters, and they vary considerably in size, the combination of state-level clustering and the WCRB always works quite well, with rejection rates never much greater than 6%.

Table 3: CRVE Rejection Percentages for Three Levels of Clustering

ϕ /DGP	Zone $t(59)$			City $t(19)$			State $t(9)$		
	zone	city ^u	state ^u	zone ^o	city	state ^u	zone ^o	city ^o	state
0.00	5.24	5.28	5.23	6.26	6.28	6.24	7.17	7.21	7.15
0.05	5.29	8.66	12.98	6.42	6.98	11.22	7.43	8.39	9.08
0.10	5.48	15.65	26.81	6.74	7.71	17.94	8.08	9.80	10.44
0.15	5.50	21.62	36.79	6.98	8.10	22.16	8.52	10.49	10.90
0.20	5.55	25.80	43.02	7.15	8.30	24.60	8.79	10.93	11.18
0.30	5.70	30.34	49.61	7.48	8.47	27.08	9.37	11.26	11.51
0.40	5.65	32.53	52.55	7.51	8.59	28.13	9.48	11.44	11.59
0.50	5.75	33.78	54.19	7.65	8.57	28.67	9.69	11.45	11.67

The table shows rejection percentages for tests at the 5% level based on Monte Carlo experiments with 400,000 replications. An “o” superscript indicates over-clustering, and a “u” superscript indicates under-clustering.

The results in Table 4 suggest that, provided we use the WCRB, the cost of over-clustering is quite modest when the null hypothesis is true. But what if that hypothesis is false? Figure 1 shows six power functions for tests of $\beta_2 = 0$ in the model (18). For the three blue curves, clustering is actually at the zone level. The solid curve shows power when the CRVE correctly clusters at the zone level, the dashed curve shows power when the CRVE over-clusters at the city level, and the dotted curve shows power when the CRVE over-clusters at the state

Table 4: WCR Bootstrap Rejection Percentages for Three Levels of Clustering

ϕ/DGP	Zone WCRB			City WCRB			State WCRB		
	zone	city ^u	state ^u	zone ^o	city	state ^u	zone ^o	city ^o	state
0.00	4.99	5.02	4.95	5.10	5.11	5.07	5.28	5.38	5.30
0.05	4.98	8.19	12.39	5.12	5.21	8.25	5.35	5.50	5.62
0.10	5.03	14.79	25.69	5.21	5.27	12.51	5.46	5.71	5.85
0.15	5.00	20.44	35.50	5.15	5.31	15.21	5.57	5.83	5.88
0.20	4.98	24.35	41.48	5.16	5.31	16.83	5.60	5.93	5.97
0.30	5.03	28.66	48.00	5.26	5.37	18.52	5.76	6.05	6.02
0.40	4.98	30.77	50.84	5.21	5.40	19.32	5.74	6.06	6.13
0.50	5.02	31.95	52.46	5.27	5.35	19.80	5.84	6.01	6.17

The table shows rejection percentages for tests at the 5% level based on Monte Carlo experiments with 400,000 replications. The WCRB used the Rademacher distribution with 399 bootstrap replications. An “o” superscript indicates over-clustering, and a “u” superscript indicates under-clustering.

level. There is evidently some power loss due to over-clustering, which is more severe when we cluster at the state level than at the city level.

For the two red curves, clustering is actually at the city level. Power is very much less than it was with clustering at the zone level, because there are far more non-zero off-diagonal elements in the Ω matrix. Oddly, except for very large values of β_2 , power seems to be slightly higher when the CRVE over-clusters at the state level than when it correctly clusters at the city level. This apparent gain in power is spurious, of course. It arises because the test over-rejects a bit more in the former case than in the latter.

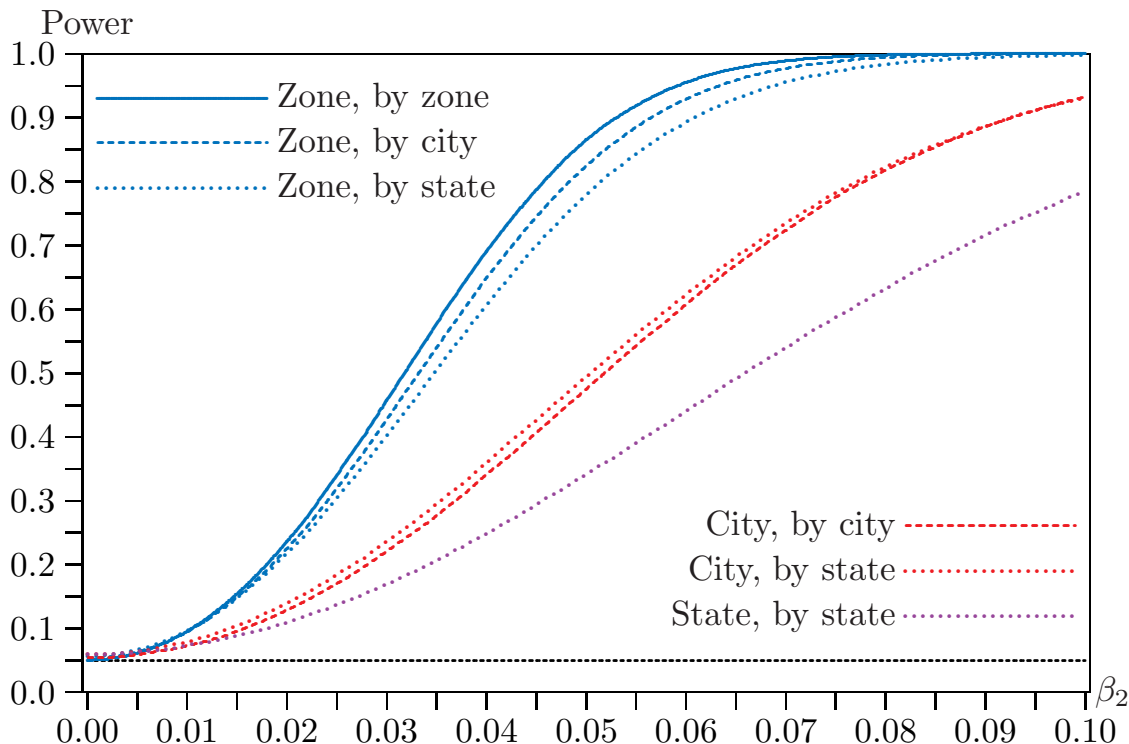
For the purple curve, clustering is actually at the state level. The only valid way for the CRVE to cluster is also at the state level, so only one curve is shown. The power loss, relative to clustering at the city level, is quite severe. Once again, this is not at all surprising. With only 10 states, there are a great many non-zero off-diagonal elements in the Ω matrix. Thus the sample contains much less information than it did with clustering at a finer level.

8 Conclusion

Ignoring the possibility that error terms may be correlated within clusters can cause severe problems for inference, especially when the sample size is large. The standard solution is to use t -statistics based on cluster-robust standard errors from a CRVE, together with the $t(G - 1)$ distribution. This approach generally works well when there is a large number of clusters that are balanced in terms of cluster sizes, the variances of the error terms, and the distributions of the key regressor(s). When these conditions are not met, however, the standard approach can yield very unreliable inferences.

The restricted wild cluster bootstrap (WCRB) works well in a far broader set of circumstances than the standard approach. Accordingly, we advocate using it as the default method of statistical inference. Computational tricks pioneered in the `boottest` package for Stata make computing the WCRB very fast and easy in most cases. There are still

Figure 1: Power of WCR Bootstrap Tests



Notes: In these experiments, $N = 6000$ and $\phi = 0.2$. The WCRB used the Rademacher distribution with 399 bootstrap replications. The label “ X , by y ” means that the disturbances are actually clustered at level X , and the CRVE is clustered at level y .

situations in which the WCRB will not be particularly reliable, most notably when there are very few clusters (say, five or less), when clusters are very heterogeneous, or when there are very few treated (or control) clusters. In such cases, it is important to check whether the WCRB is reliable. One approach is to compare WCRB P values or confidence intervals with those from alternative procedures such as randomization inference or the procedures discussed near the end of Section 5.1. If the results of the WCRB broadly agree with those from other procedures, then it seems reasonable to accept the former. A second approach is to conduct a Monte Carlo experiment that mimics the model and dataset on hand. This approach can be particularly useful in unusual settings with severely unbalanced clusters.

In Section 4, we suggested a few guidelines for how to cluster. In general, one should cluster at the broadest level possible. The simulation experiments in Section 7 suggest that there can be large size distortions from under-clustering compared with much smaller power losses from over-clustering. This is especially true in large samples. When studying the effects of policy changes, one should always cluster at least at the level of the policy change, and perhaps at a more aggregate level. When working with panel data, it is desirable to cluster by the cross-section dimension (perhaps in addition to other dimensions) in order to capture any serial correlation.

The validity of the asymptotic approximations that underlie cluster-robust inference is

driven by the number of clusters rather than the sample size. Unfortunately, and contrary to popular belief, there is no “golden number” of clusters beyond which CRVE-based inference becomes reliable. The requisite number of clusters depends on many factors. These include how much the cluster sizes vary, how unbalanced the clusters are in other respects, and, in many cases, how many clusters are treated, how many clusters are not treated, and the numbers of treated observations in each of the treated clusters. Because the WCRB conditions on all these aspects of the sample, which no “golden number” could ever do, our recommendation is to use the restricted wild cluster bootstrap essentially all the time.

Much more is known about cluster-robust inference than was the case even ten years ago. Nevertheless, there are still many unanswered questions. We do not know how to control for both spatial correlation and conventional within-cluster correlation at the same time. We do not know how to obtain reliable inferences when there are few treated clusters and the treated clusters are atypical. We do not really know how to determine the right level at which to cluster, especially when over-clustering means having few clusters or few treated clusters. Finally, we do not know how to conduct inference when there is a very small number of clusters, say, five or less.

References

- Abadie, A., S. Athey, G. W. Imbens, and J. Wooldridge (2017). When should you adjust standard errors for clustering? NBER Working Papers 24003, National Bureau of Economic Research, Inc.
- Abadie, A., S. Athey, G. W. Imbens, and J. Wooldridge (2019). Sampling-based vs. design-based uncertainty in regression analysis. *Econometrica* (to appear).
- Angrist, J. D. and J.-S. Pischke (2008). *Mostly Harmless Econometrics: An Empiricist's Companion* (1 ed.). Princeton University Press.
- Arellano, M. (1987). Computing robust standard errors for within-groups estimators. *Oxford Bulletin of Economics and Statistics* 49(4), 431–434.
- Barrios, T., R. Diamond, G. W. Imbens, and M. Kolesár (2012). Clustering, spatial correlations, and randomization inference. *Journal of the American Statistical Association* 107(498), 578–591.
- Bell, R. M. and D. F. McCaffrey (2002). Bias reduction in standard errors for linear regression with multi-stage samples. *Survey Methodology* 28(2), 169–181.
- Bertrand, M., E. Duflo, and S. Mullainathan (2004). How much should we trust differences-in-differences estimates? *The Quarterly Journal of Economics* 119(1), 249–275.
- Bester, C. A., T. G. Conley, and C. B. Hansen (2011). Inference with dependent data using cluster covariance estimators. *Journal of Econometrics* 165(1), 137–151.
- Bloom, N., R. Sadun, and J. Van Reenen (2012). Americans do it better: Us multinationals and the productivity miracle. *American Economic Review* 102(1), 167–201.
- Brewer, M., T. F. Crossley, and R. Joyce (2018). Inference with difference-in-differences revisited. *Journal of Econometric Methods* 7(1), 1–16.
- Cameron, A. C., J. B. Gelbach, and D. L. Miller (2008). Bootstrap-based improvements for inference with clustered errors. *The Review of Economics and Statistics* 90(3), 414–427.
- Cameron, A. C., J. B. Gelbach, and D. L. Miller (2011). Robust inference with multiway clustering. *Journal of Business & Economic Statistics* 29(2), 238–249.
- Cameron, A. C. and D. L. Miller (2015). A practitioner's guide to cluster-robust inference. *Journal of Human Resources* 50(2), 317–372.
- Canay, I. A., J. P. Romano, and A. M. Shaikh (2017). Randomization tests under an approximate symmetry assumption. *Econometrica* 85(3), 1013–1030.
- Canay, I. A., A. Santos, and A. Shaikh (2018). The wild bootstrap with a 'small' number of 'large' clusters. Technical Report 2019-17.
- Carter, A. V., K. T. Schnepel, and D. G. Steigerwald (2017). Asymptotic behavior of a t -test robust to cluster heterogeneity. *The Review of Economics and Statistics* 99(4), 698–709.
- Conley, T. (1999). GMM estimation with cross sectional dependence. *Journal of Econometrics* 92(1), 1–45.
- Conley, T. G. and C. R. Taber (2011). Inference with “Difference in Differences” with a small number of policy changes. *The Review of Economics and Statistics* 93(1), 113–125.

- Davezies, L., X. D’Haultfoeulle, and Y. Guyonvarch (2018). Asymptotic results under multiway clustering. ArXiv e-prints, CREST.
- Davidson, R. and E. Flachaire (2008). The wild bootstrap, tamed at last. *Journal of Econometrics* 146(1), 162–169.
- Davidson, R. and J. G. MacKinnon (2000). Bootstrap tests: How many bootstraps? *Econometric Reviews* 19(1), 55–68.
- Davidson, R. and J. G. MacKinnon (2006). Bootstrap methods in econometrics. In T. C. Mills and K. D. Patterson (Eds.), *Palgrave Handbook of Econometrics: Volume 1 Econometric Theory*, pp. 812–838. Palgrave Macmillan.
- Djogbenou, A. A., J. G. MacKinnon, and M. Ø. Nielsen (2019). Asymptotic theory and wild bootstrap inference with clustered errors. *Journal of Econometrics* 213(to appear).
- Donald, S. G. and K. Lang (2007). Inference with difference-in-differences and other panel data. *The Review of Economics and Statistics* 89(2), 221–233.
- Esarey, J. (2018). clusterSEs: Calculate Cluster-Robust p-Values and Confidence Intervals. *R package version 2.6.1*.
- Esarey, J. and A. Menger (2019). Practical and effective approaches to dealing with clustered data. *Political Science Research and Methods* 7(3), 541–559.
- Ferman, B. (2019). Inference in differences-in-differences: How much should we trust in independent clusters? MPRA Paper 93746, University Library of Munich, Germany.
- Ferman, B. and C. Pinto (2019). Inference in differences-in-differences with few treated groups and heteroskedasticity. *The Review of Economics and Statistics* 101, 452–467.
- Fisher, R. A. (1935). *The Design of Experiments*. Edinburgh: Oliver and Boyd.
- Giné, X. and G. Mansuri (2018). Together we will: Experimental evidence on female voting behavior in pakistan. *American Economic Journal: Applied Economics* 10(1), 207–35.
- Hagemann, A. (2019a). Placebo inference on treatment effects when the number of clusters is small. *Journal of Econometrics* 213(to appear).
- Hagemann, A. (2019b). Permutation inference with a finite number of heterogeneous clusters. ArXiv e-prints 1907.01049 [econ.EM].
- Hansen, B. E. and S. Lee (2019). Asymptotic theory for clustered samples. *Journal of Econometrics* 210(2), 268–290.
- Hess, S. (2017). Randomization inference with Stata: A guide and software. *Stata Journal* 17(3), 630–651.
- Horowitz, J. L. (2019). Bootstrap methods in econometrics. *Annual Review of Economics* 11(1), 193–224.
- Ibragimov, R. and U. K. Müller (2010). t-statistic based correlation and heterogeneity robust inference. *Journal of Business & Economic Statistics* 28(4), 453–468.
- Ibragimov, R. and U. K. Müller (2016). Inference with few heterogeneous clusters. *The Review of Economics and Statistics* 98(1), 83–96.
- Imbens, G. W. and M. Kolesár (2016). Robust standard errors in small samples: Some

- practical advice. *The Review of Economics and Statistics* 98(4), 701–712.
- Imbens, G. W. and D. B. Rubin (2015). *Causal Inference in Statistics, Social, and Biomedical Sciences*. New York: Cambridge University Press.
- Jackson, J. E. (2019). Corrected standard errors with clustered data. *Political Analysis*, to appear.
- Kelly, M. (2019, 06). The standard errors of persistence. Technical report.
- Kloek, T. (1981). Ols estimation in a model where a microvariable is explained by aggregates and contemporaneous disturbances are equicorrelated. *Econometrica* 49(1), 205–207.
- Kolenikov, S. (2010). Resampling variance estimation for complex survey data. *The Stata Journal* 10(2), 165–199.
- Lee, C. H. and D. G. Steigerwald (2018). Inference for clustered data. *Stata Journal* 18(2), 447–460.
- Lehmann, E. L. and J. P. Romano (2008). *Testing Statistical Hypotheses*. New York: Springer.
- Liang, K.-Y. and S. L. Zeger (1986). Longitudinal data analysis using generalized linear models. *Biometrika* 73(1), 13–22.
- MacKinnon, J. G. (2002). Bootstrap inference in econometrics. *Canadian Journal of Economics* 35(4), 615–645.
- 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*, pp. 437–461. Springer.
- MacKinnon, J. G. (2015). Wild cluster bootstrap confidence intervals. *L'Actualité Économique* 91, 11–33.
- MacKinnon, J. G. (2016). Inference with large clustered datasets. *L'Actualité Économique* 92, 649–665.
- MacKinnon, J. G. (2019). How cluster-robust inference is changing applied econometrics. *Canadian Journal of Economics* 52(3), to appear.
- MacKinnon, J. G., M. Ø. Nielsen, and M. D. Webb (2019). Wild bootstrap and asymptotic inference with multiway clustering. QED Working Paper 1415, Queen’s University, Department of Economics.
- MacKinnon, J. G. and M. D. Webb (2017a). Wild bootstrap inference for wildly different cluster sizes. *Journal of Applied Econometrics* 32(2), 233–254.
- MacKinnon, J. G. and M. D. Webb (2017b). Pitfalls when estimating treatment effects using clustered data. *The Political Methodologist* 24(2), 20–31.
- MacKinnon, J. G. and M. D. Webb (2018). The wild bootstrap for few (treated) clusters. *Econometrics Journal* 21(2), 114–135.
- MacKinnon, J. G. and M. D. Webb (2019). Wild bootstrap randomization inference for few treated clusters. In K. P. Huynh, D. T. Jacho-Chávez, and G. Tripathi (Eds.), *The Econometrics of Complex Survey Data: Theory and Applications*, Volume 39 of *Advances in Econometrics*, Chapter 3, pp. 61–85. Emerald Group.

- MacKinnon, J. G. and M. D. Webb (2020). Randomization inference for difference-in-differences with few treated clusters. *Journal of Econometrics* (to appear).
- MacKinnon, J. G. and H. White (1985). Some heteroskedasticity consistent covariance matrix estimators with improved finite sample properties. *Journal of Econometrics* 29(3), 305–325.
- Menzel, K. (2018). Bootstrap with cluster-dependence in two or more dimensions. ArXiv e-prints, New York University.
- Moulton, B. R. (1986). Random group effects and the precision of regression estimates. *Journal of Econometrics* 32(3), 385–397.
- Moulton, B. R. (1990). An illustration of a pitfall in estimating the effects of aggregate variables on micro units. *The Review of Economics and Statistics* 72(2), 334–338.
- Nickell, S. (1981). Biases in dynamic models with fixed effects. *Econometrica* 49, 1417–1426.
- Pustejovsky, J. (2017). clubsandwich: Cluster-robust (sandwich) variance estimators with small-sample corrections. *R package version 0.3.5*.
- Riddell, W. C. (1979). The empirical foundations of the Phillips curve: Evidence from Canadian wage contract data. *Econometrica* 47(1), 1–24.
- Rogers, W. H. (1993). Regression standard errors in clustered samples. *Stata Technical Bulletin* 13, 19–23.
- Roodman, D., J. G. MacKinnon, M. Ø. Nielsen, and M. D. Webb (2019). Fast and wild: bootstrap inference in Stata using boottest. *Stata Journal* 19(1), 4–60.
- Thompson, S. B. (2011). Simple formulas for standard errors that cluster by both firm and time. *Journal of Financial Economics* 99, 1–10.
- Webb, M. D. (2014). Reworking wild bootstrap based inference for clustered errors. QED Working Paper 1315, Queen’s University, Department of Economics.
- White, H. (1980). A heteroskedasticity-consistent covariance matrix estimator and a direct test for heteroskedasticity. *Econometrica* 48(4), 817–838.
- Young, A. (2016). Improved, nearly exact, statistical inference with robust and clustered covariance matrices using effective degrees of freedom corrections. Working paper, London School of Economics.