

# Small sample correction methods for cluster-robust variance estimators and hypothesis tests

James E. Pustejovsky\*

Department of Educational Psychology  
University of Texas at Austin

and

Elizabeth Tipton

Department of Human Development  
Teachers College, Columbia University

August 11, 2015

## **Abstract**

The text of your abstract. 200 or fewer words.

*Keywords:* 3 to 6 keywords, that do not appear in the title

---

\*The authors gratefully acknowledge *please remember to list all relevant funding sources in the unblinded version*

# 1 Introduction

Cluster-robust variance estimators (CRVE) and hypothesis tests based upon such estimators are ubiquitous in applied econometric work. Nearly every respectable paper in the past 15 years uses cluster-robust variance estimators because to do otherwise would be to risk being seen as insufficiently rigorous (or anti-conservative....ughh....how gauche!).

There's been a lot of fretting recently that even CRVE may actually not be rigorous enough. Cite the following people so as not to get their ire up:

- Brewer et al. (2013)
- Cameron et al. (2008)
- Cameron & Miller (2015)
- Carter et al. (2013)
- Ibragimov & Müller (2010)
- Imbens & Kolesar (2012)
- Kezdi (2004)
- McCaffrey et al. (2001)
- McCaffrey & Bell (2006)
- Webb & MacKinnon (2013)
- Kline & Santos (2012)

## 1.1 Econometric framework

We will consider linear regression models in which the errors within a cluster have an unknown variance structure. The model is

$$\mathbf{Y}_j = \mathbf{X}_j\boldsymbol{\beta} + \boldsymbol{\epsilon}_j, \tag{1}$$

for  $j = 1, \dots, m$ , where  $\mathbf{Y}_j$  is  $n_j \times 1$ ,  $\mathbf{X}_j$  is an  $n_j \times p$  matrix of regressors for cluster  $j$ ,  $\beta$  is a  $p \times 1$  vector, and  $\epsilon_j$  is an  $n_j \times 1$  vector of errors. Assume that  $E(\epsilon_j | \mathbf{X}_j) = \mathbf{0}$  and  $\text{Var}(\epsilon_j | \mathbf{X}_j) = \Sigma_j$ , for  $j = 1, \dots, m$ , where  $\Sigma_1, \dots, \Sigma_m$  may be unknown, and the errors are independent across clusters. Let  $\mathbf{X} = (\mathbf{X}'_1, \mathbf{X}'_2, \dots, \mathbf{X}'_m)'$  and  $\Sigma = \bigoplus_{j=1}^m \Sigma_j$ .

The vector of regression coefficients is estimated by weighted least squares (WLS). Given a set of  $m$  symmetric weighting matrices  $\mathbf{W}_1, \dots, \mathbf{W}_m$ , the WLS estimator is

$$\hat{\beta} = \mathbf{M} \sum_{j=1}^m \mathbf{X}'_j \mathbf{W}_j \mathbf{Y}_j, \quad (2)$$

where  $\mathbf{M} = \left( \sum_{j=1}^m \mathbf{X}'_j \mathbf{W}_j \mathbf{X}_j \right)^{-1}$ . Let  $\mathbf{W} = \bigoplus_{j=1}^m \mathbf{W}_j$ .

Common choices for weighting include the unweighted case, in which  $\mathbf{W}_j$  is an identity matrix of dimension  $n_j \times n_j$ , and inverse-variance weighting under a working model. In the latter case, the errors are assumed to follow some known structure,  $\text{Var}(\epsilon_j | \mathbf{X}_j) = \Phi_j$ , where  $\Phi_j$  is a known function of a low-dimensional parameter. The weighting matrices are then taken to be  $\mathbf{W}_j = \hat{\Phi}_j^{-1}$ , where the  $\hat{\Phi}_j$  are constructed from estimates of the variance parameters.

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

## 2 Cluster-robust variance estimators

The variance of the WLS estimator is

$$\text{Var}(\hat{\beta}) = \mathbf{M} \left( \sum_{j=1}^m \mathbf{X}'_j \mathbf{W}_j \Sigma_j \mathbf{W}_j \mathbf{X}_j \right) \mathbf{M}, \quad (3)$$

which includes the unknown variance matrices. One approach to estimating this variance would be to posit a working model—potentially the same working model used to construct weights—and substitute estimates of the working variance structure in place of  $\Sigma$ . Under a

working model  $\Phi$ , denote this "model-based" variance estimator as

$$\mathbf{V}^M = \mathbf{M} \left( \sum_{j=1}^m \mathbf{X}_j' \mathbf{W}_j \hat{\Phi}_j \mathbf{W}_j \mathbf{X}_j \right) \mathbf{M}. \quad (4)$$

If  $\beta$  is estimated using inverse-variance weights defined under the same working model, then the model-based variance estimator simplifies to  $\mathbf{V}^M = \mathbf{M}$ .

Cluster-robust variance estimators provide a means of estimating  $\text{Var}(\hat{\beta})$  and test hypotheses regarding  $\hat{\beta}$  in the absence of a valid working model for the error structure, or if the working variance model used to develop weights is mis-specified. They are thus a generalization of heteroskedasticity-consistent (HC) variance estimators. Like the HC estimators, several different variants have been proposed, with different rationales and different finite-sample properties

The most widely used estimator is

$$\mathbf{V}^{CR0} = \mathbf{M} \left( \sum_{j=1}^m \mathbf{X}_j' \mathbf{W}_j \mathbf{e}_j \mathbf{e}_j' \mathbf{W}_j \mathbf{X}_j \right) \mathbf{M}, \quad (5)$$

where  $\mathbf{e}_j = \mathbf{Y}_j - \mathbf{X}_j \hat{\beta}$ . Following the naming conventions used by Cameron & Miller (2015), we will refer to this estimator as CR0. Note that CR0 is constructed by substituting  $\mathbf{e}_j \mathbf{e}_j'$  in place of  $\Sigma_j$  in (3). Although the individual squared residuals provide only very crude estimates of the unknown variance matrices, the resulting estimator is asymptotically consistent for the variance of  $\hat{\beta}$  as  $m$  increases (CITE). However, CR0 is known to have a downward bias when the number of independent clusters is small (CITE).

The small-sample bias in CR0 can be seen as analogous to that arising by estimating the variance of a sample of  $m$  observations using a denominator of  $m$  rather than  $m - 1$ . One approach to correcting this bias is to scale CR0 by a factor of  $m/(m - 1)$ . Thus, define the CR1 estimator as  $\mathbf{V}^{CR1} = [m/(m - 1)] \mathbf{V}^{CR0}$ . Some software implementations use a slightly different correction factor. For example, the Stata command `regress` scales CR0 by the factor  $m(N - 1)/[(m - 1)(N - p)]$ , where  $N = \sum_{j=1}^m n_j$ .

## 2.1 Correction based on a working-model

McCaffrey et al. (2001) proposed to correct the small-sample bias of CR0 so that it is exactly unbiased under a specified working model for the variance structure. In their

implementation, the residuals from each cluster are multiplied by an  $n_j \times n_j$  adjustment matrix  $\mathbf{A}_j$  that creates the unbiasedness property. The variance estimator, which we will call CR2a, is then

$$\mathbf{V}^{CR2a} = \mathbf{M} \left( \sum_{j=1}^m \mathbf{X}_j' \mathbf{W}_j \mathbf{A}_j \mathbf{e}_j \mathbf{e}_j' \mathbf{A}_j' \mathbf{W}_j \mathbf{X}_j \right) \mathbf{M}, \quad (6)$$

The adjustment matrices are chosen to satisfy

$$\mathbf{A}_j' (\mathbf{I} - \mathbf{H})_j \boldsymbol{\Phi} (\mathbf{I} - \mathbf{H})_j' \mathbf{A}_j = \boldsymbol{\Phi}_j, \quad (7)$$

where  $\boldsymbol{\Phi} = \bigoplus_{j=1}^m \boldsymbol{\Phi}_j$ ,  $\mathbf{H} = \mathbf{X} \mathbf{M} \mathbf{X}' \mathbf{W}$ , and  $(\mathbf{I} - \mathbf{H})_j$  denotes the rows of  $\mathbf{I} - \mathbf{H}$  corresponding to cluster  $j$ . However, this criterion does not uniquely define  $\mathbf{A}_j$ . Based on extensive simulations, McCaffrey et al. (2001) found that a symmetric solution worked well, with

$$\mathbf{A}_j = (\boldsymbol{\Phi}_j^C)' \mathbf{B}_j^{-1/2} \boldsymbol{\Phi}_j^C,$$

where  $\boldsymbol{\Phi}_j^C$  is the upper triangular Cholesky factorization of  $\boldsymbol{\Phi}_j$ ,

$$\mathbf{B}_j = \boldsymbol{\Phi}_j^C (\mathbf{I} - \mathbf{H})_j \boldsymbol{\Phi} (\mathbf{I} - \mathbf{H})_j' (\boldsymbol{\Phi}_j^C)',$$

and  $\mathbf{B}_j^{-1/2}$  is the inverse of the symmetric square root of  $\mathbf{B}$ . If ordinary (unweighted) least squares is used to estimate  $\hat{\boldsymbol{\beta}}$  and the working variance model posits that the errors are all independent and homoskedastic, then  $\mathbf{W} = \boldsymbol{\Phi} = \mathbf{I}_N$  and  $\mathbf{A}_j = (\mathbf{I}_j - \mathbf{H}_{jj})^{-1/2}$ , where  $\mathbf{I}_j$  is an  $n_j \times n_j$  identity matrix and  $\mathbf{H}_{jj} = \mathbf{X}_j (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}_j'$ .

Two difficulties arise in the implementation of CR2a. First, the matrices  $\mathbf{B}_1, \dots, \mathbf{B}_m$  may not be positive definite, so that  $\mathbf{B}_j^{-1/2}$  cannot be calculated for every cluster. This occurs, for instance, in balanced panel models when the specification includes fixed effects for each unit and each timepoint and clustering is over the units (Angrist & Pischke 2009, p. ??). However, this problem is easily overcome by using a generalized inverse of  $\mathbf{B}_j$ . A second, computational difficulty with CR2a is that it requires the inversion (or pseudo-inversion) of  $m$  matrices, each of dimension  $n_j \times n_j$ . Consequently, computation of CR2a will be slow if some clusters contain a large number of individual units.

## 2.2 Another working-model correction

The criterion (7) is not the only way to obtain a variance estimator that is precisely unbiased under a working model. An alternative approach, which to our knowledge is novel, is to

use

$$\mathbf{V}^{CR2b} = \mathbf{M} \left( \sum_{j=1}^m \mathbf{D}_j \mathbf{X}_j' \mathbf{W}_j \mathbf{e}_j \mathbf{e}_j' \mathbf{W}_j \mathbf{X}_j \mathbf{D}_j' \right) \mathbf{M}, \quad (8)$$

where the adjustment matrices  $\mathbf{D}_1, \dots, \mathbf{D}_m$  are chosen so that

$$\mathbf{D}_j \mathbf{X}_j' \mathbf{W}_j (\mathbf{I} - \mathbf{H})_j \boldsymbol{\Phi} (\mathbf{I} - \mathbf{H})_j' \mathbf{W}_j \mathbf{X}_j \mathbf{D}_j' = \mathbf{X}_j' \mathbf{W}_j \boldsymbol{\Phi} \mathbf{W}_j \mathbf{X}_j. \quad (9)$$

Just as with CR2a, there are several different forms of adjustment matrices that satisfy (9). A symmetric solution is to take

$$\mathbf{D}_j = (\mathbf{F}_j^C)' \left[ \mathbf{F}_j^C \mathbf{G}_j (\mathbf{F}_j^C)' \right]^{-1/2} \mathbf{F}_j^C \quad (10)$$

where  $\mathbf{F}_j = \mathbf{X}_j' \mathbf{W}_j \boldsymbol{\Phi} \mathbf{W}_j \mathbf{X}_j$  and  $\mathbf{G}_j = \mathbf{X}_j' \mathbf{W}_j (\mathbf{I} - \mathbf{H})_j \boldsymbol{\Phi} (\mathbf{I} - \mathbf{H})_j' \mathbf{W}_j \mathbf{X}_j$ . Note that  $\mathbf{F}_j$  and  $\mathbf{G}_j$  might not be positive definite, and so a generalized form of the Cholesky decomposition and symmetric inverse-square root must be used. The  $\mathbf{D}_j$  correction matrices will be defined (I think) so long as  $\text{rank}(\mathbf{F}_j) \geq \text{rank}(\mathbf{G}_j)$  for  $j = 1, \dots, m$ . In datasets where clusters have a large number of individual units, CR2b will be less computationally intensive than CR2a because the adjustment matrices are all of dimension  $p \times p$ .

## 2.3 Jackknife correction

## 2.4 Considerations with panel models

# 3 Hypothesis testing

## 3.1 Single-constraint tests

## 3.2 Multiple-constraint tests

# 4 Examples

# 5 Simulation evidence

# 6 Discussion

## References

- Angrist, J. D. & Pischke, J. (2009), *Mostly harmless econometrics: An empiricist's companion*, Princeton University Press, Princeton, NJ.
- Brewer, M., Crossley, T. F. & Joyce, R. (2013), Inference with difference-in-differences revisited.
- Cameron, A. C., Gelbach, J. B. & Miller, D. (2008), ‘Bootstrap-based improvements for inference with clustered errors’, *The Review of Economics and Statistics* **90**(3), 414–427.
- Cameron, A. C. & Miller, D. L. (2015), A practitioner’s guide to cluster-robust inference.
- Carter, A. V., Schnepel, K. T. & Steigerwald, D. G. (2013), ‘Asymptotic Behavior of a t Test Robust to Cluster Heterogeneity’, pp. 1–32.
- Ibragimov, R. & Müller, U. K. (2010), ‘t-Statistic based correlation and heterogeneity robust inference’, *Journal of Business & Economic Statistics* **28**(4), 453–468.

Imbens, G. W. & Kolesar, M. (2012), Robust standard errors in small samples: some practical advice.

**URL:** <http://www.nber.org/papers/w18478>

Kezdi, G. (2004), Robust standard error estimation in fixed-effects panel models.

**URL:** <http://papers.ssrn.com/sol3/Delivery.cfm?abstractid=596988>

Kline, P. & Santos, A. (2012), ‘A score based approach to wild bootstrap inference’, *Journal of Econometric Methods* **1**(1).

McCaffrey, D. F. & Bell, R. M. (2006), ‘Improved hypothesis testing for coefficients in generalized estimating equations with small samples of clusters.’, *Statistics in medicine* **25**(23), 4081–98.

McCaffrey, D. F., Bell, R. M. & Botts, C. H. (2001), Generalizations of biased reduced linearization, *in* ‘Proceedings of the Annual Meeting of the American Statistical Association’, number 1994.

Webb, M. & MacKinnon, J. G. (2013), Wild Bootstrap Inference for Wildly Different Cluster Sizes.