

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

James E. Pustejovsky <sup>\*</sup>  
University of Wisconsin - Madison  
and  
Elizabeth Tipton <sup>†</sup>  
Northwestern University

January 15, 2023

## Abstract

Pustejovsky and Tipton (2018) considered how to implement cluster-robust variance estimators for fixed effects models estimated by weighted (or unweighted) least squares. Theorem 2 of the paper concerns a computational short cut for a certain cluster-robust variance estimator in models with cluster-specific fixed effects. It claimed that this short cut works for models estimated by generalized least squares, as long as the weights are taken to be inverse of the working model. However, the theorem is incorrect. In this corrigendum, we review the CR2 variance estimator, describe the assertion of the theorem as originally stated, and demonstrate the error with a counter-example. We then provide a revised version of the theorem, which holds for the more limited set of models estimated by ordinary least squares.

*Keywords:*

---

<sup>\*</sup>Department of Educational Psychology, University of Wisconsin - Madison, 1025 West Johnson Street, Madison, WI 53706. Email: [pustejovsky@wisc.edu](mailto:pustejovsky@wisc.edu)

<sup>†</sup>Department of Statistics, Northwestern University. Email: [tipton@northwestern.edu](mailto:tipton@northwestern.edu)

# 1 A fixed effects model

For data that can be grouped into  $m$  clusters of observations, Pustejovsky and Tipton (2018) considered the model

$$\mathbf{y}_i = \mathbf{R}_i\boldsymbol{\beta} + \mathbf{S}_i\boldsymbol{\gamma} + \mathbf{T}_i\boldsymbol{\mu} + \boldsymbol{\epsilon}_i, \quad (1)$$

where  $\mathbf{y}_i$  is an  $n_i \times 1$  vector of responses for cluster  $i$ ,  $\mathbf{R}_i$  is an  $n_i \times r$  matrix of focal predictors,  $\mathbf{S}_i$  is an  $n_i \times s$  matrix of additional covariates that vary across multiple clusters, and  $\mathbf{T}_i$  is an  $n_i \times t$  matrix encoding cluster-specific fixed effects, all for  $i = 1, \dots, m$ . The cluster-specific fixed effects satisfy  $\mathbf{T}_h\mathbf{T}_i' = \mathbf{0}$  for  $h \neq i$ . Interest centers on inference for the coefficients on the focal predictors  $\boldsymbol{\beta}$ .

Pustejovsky and Tipton (2018) considered estimation of Model 1 by generalized or weighted least squares (WLS). Let  $\mathbf{W}_1, \dots, \mathbf{W}_m$  be a set of symmetric weight matrices used for WLS estimation, which may include off-diagonal elements. The CR2 variance estimator involves specifying a working model for the structure of the errors. Consider a working model  $\text{Var}(\boldsymbol{\epsilon}_i | \mathbf{R}_i, \mathbf{S}_i, \mathbf{T}_i) = \sigma^2 \boldsymbol{\Phi}_i$  where  $\boldsymbol{\Phi}_i$  is a symmetric  $n_i \times n_i$  matrix that may be a function of a low-dimensional, estimable parameter. In some applications, the weight matrices might be taken as  $\mathbf{W}_i = \hat{\boldsymbol{\Phi}}_i^{-1}$ , where  $\hat{\boldsymbol{\Phi}}_i$  is an estimate of  $\boldsymbol{\Phi}_i$ . In other applications, the weight matrices may be something else, such as diagonal matrices consisting of sampling weights or identity matrices (i.e., ordinary least squares).

# 2 The CR2 variance estimator

Pustejovsky and Tipton (2018) provided a generalization of the bias-reduced linearization estimator introduced by McCaffrey, Bell, and Botts (2001) and Bell and McCaffrey (2002) that can be applied to Model 1, referred to as the CR2 variance estimator. We follow the same notation as Pustejovsky and Tipton (2018) to define CR2. Let  $N = \sum_{i=1}^m n_i$  be the total sample size. Let  $\mathbf{U}_i = [\mathbf{R}_i \ \mathbf{S}_i]$  be the set of predictors that vary across clusters and  $\mathbf{X}_i = [\mathbf{R}_i \ \mathbf{S}_i \ \mathbf{T}_i]$  be the full set of predictors. Let  $\mathbf{R}$ ,  $\mathbf{S}$ ,  $\mathbf{T}$ ,  $\mathbf{U}$ ,  $\mathbf{X}$ , and  $\mathbf{y}$  denote the stacked versions of the cluster-specific matrices (i.e.,  $\mathbf{R} = [\mathbf{R}_1' \ \mathbf{R}_2' \ \dots \ \mathbf{R}_m']'$ , etc.). Let  $\mathbf{W} = \bigoplus_{i=1}^m \mathbf{W}_i$  and  $\boldsymbol{\Phi} = \bigoplus_{i=1}^m \boldsymbol{\Phi}_i$ . For a generic matrix  $\mathbf{Z}$ , let  $\mathbf{M}_Z = (\mathbf{Z}'\mathbf{W}\mathbf{Z})^{-1}$  and

$\mathbf{H}_Z = \mathbf{Z}\mathbf{M}_Z\mathbf{Z}'\mathbf{W}$ . Let  $\mathbf{C}_i$  be the  $n_i \times N$  matrix that selects the rows of cluster  $i$  from the full set of observations, such that  $\mathbf{X}_i = \mathbf{C}_i\mathbf{X}$ . Finally, let  $\mathbf{D}_i$  be the upper-right Cholesky factorization of  $\Phi_i$ .

These operators provide a means to define absorbed versions of the predictors and the outcome. Let  $\ddot{\mathbf{S}} = (\mathbf{I} - \mathbf{H}_T)\mathbf{S}$  be the covariates after absorbing the cluster-specific effects, let  $\ddot{\mathbf{U}} = (\mathbf{I} - \mathbf{H}_T)\mathbf{U}$  be an absorbed version of the focal predictors and the covariates, let  $\ddot{\mathbf{R}} = (\mathbf{I} - \mathbf{H}_{\ddot{\mathbf{S}}})(\mathbf{I} - \mathbf{H}_T)\mathbf{R}$  be the focal predictors after absorbing the covariates and the cluster-specific fixed effects, and let  $\mathbf{e} = (\mathbf{I} - \mathbf{H}_{\ddot{\mathbf{R}}})(\mathbf{I} - \mathbf{H}_{\ddot{\mathbf{S}}})(\mathbf{I} - \mathbf{H}_T)\mathbf{y}$  denote the vector of residuals, with  $\mathbf{e}_i = \mathbf{C}_i\mathbf{e}$ .

With this notation established, the CR2 variance estimator has the form

$$\mathbf{V}^{CR2} = \mathbf{M}_{\ddot{\mathbf{R}}} \left( \sum_{i=1}^m \ddot{\mathbf{R}}_i' \mathbf{W}_i \mathbf{A}_i \mathbf{e}_i \mathbf{e}_i' \mathbf{A}_i' \mathbf{W}_i \ddot{\mathbf{R}}_i \right) \mathbf{M}_{\ddot{\mathbf{R}}}, \quad (2)$$

where  $\ddot{\mathbf{R}}_i = \mathbf{C}_i\ddot{\mathbf{R}}$  is the cluster-specific matrix of absorbed focal predictors,  $\mathbf{e}_i$  is the vector of weighted least squares residuals from cluster  $i$ , and  $\mathbf{A}_1, \dots, \mathbf{A}_m$  are a set of adjustment matrices that correct the bias of the residual cross-products.

The adjustment matrices are calculated as follows. Define the matrices

$$\mathbf{B}_i = \mathbf{D}_i \mathbf{C}_i (\mathbf{I} - \mathbf{H}_X) \Phi (\mathbf{I} - \mathbf{H}_X)' \mathbf{C}_i' \mathbf{D}_i' \quad (3)$$

for  $i = 1, \dots, m$ . The adjustment matrices are then calculated as

$$\mathbf{A}_i = \mathbf{D}_i' \mathbf{B}_i^{+1/2} \mathbf{D}_i, \quad (4)$$

where  $\mathbf{B}_i^{+1/2}$  is the symmetric square root of the Moore-Penrose inverse of  $\mathbf{B}_i$ . Theorem 1 of Pustejovsky and Tipton (2018) shows that, if the working model  $\Phi$  is correctly specified and some conditions on the rank of  $\mathbf{U}$  are satisfied, then the CR2 estimator is exactly unbiased for the sampling variance of the weighted least squares estimator of  $\beta$ . Moreover, although the CR2 estimator is defined based on a working model, it remains close to unbiased and outperforms alternative sandwich estimators even when the working model is not correctly specified.

### 3 The original statement of Theorem 2

The adjustment matrices given in (4) can be expensive to compute directly because the  $\mathbf{B}_i$  matrices involve computing a “residualized” version of the  $N \times N$  matrix  $\Phi$  involving the full set of predictors  $\mathbf{X}$ —including the cluster-specific fixed effects  $\mathbf{T}_1, \dots, \mathbf{T}_m$ . Theorem 2 considered whether one can take a computational short cut by omitting the cluster-specific fixed effects from the calculation of the  $\mathbf{B}_i$  matrices. Specifically, define the modified matrices

$$\tilde{\mathbf{B}}_i = \mathbf{D}_i \mathbf{C}_i (\mathbf{I} - \mathbf{H}_{\tilde{\mathbf{U}}}) \Phi (\mathbf{I} - \mathbf{H}_{\tilde{\mathbf{U}}})' \mathbf{C}_i' \mathbf{D}_i' \quad (5)$$

and

$$\tilde{\mathbf{A}}_i = \mathbf{D}_i' \tilde{\mathbf{B}}_i^{+1/2} \mathbf{D}_i. \quad (6)$$

Theorem 2 claimed that if the weight matrices are inverse of the working model, such that  $\mathbf{W}_i = \Phi_i^{-1}$  for  $i = 1, \dots, m$ , then  $\tilde{\mathbf{B}}_i^{+1/2} = \mathbf{B}_i^{+1/2}$  and hence  $\tilde{\mathbf{A}}_i = \mathbf{A}_i$ . The implication is that the cluster-specific fixed effects can be ignored when calculating the adjustment matrices. However, the claimed equivalence does not hold in general. The proof of Theorem 2 as given in the supplementary materials of Pustejovsky and Tipton (2018) relied on a Woodbury identity for generalized inverses that does not hold for  $\mathbf{B}_i$  because necessary rank conditions are not satisfied.

We describe a simple numerical example that contradicts the original statement of Theorem 2. Consider a design with  $m = 3$  clusters, of sizes  $n_1 = 2$ ,  $n_2 = 3$ , and  $n_3 = 5$  for which we have the model

$$y_{it} = \beta_0 \times t + \mu_i + \epsilon_{it},$$

where the errors are heteroskedastic with  $\text{Var}(\epsilon_{it}) = \alpha \times t$ . We then have  $\mathbf{y}_i = [y_{i1} \cdots y_{in_i}]'$ ,  $\mathbf{R}_i = [1 \ 2 \ \cdots \ n_i]'$ ,  $\mathbf{T} = \bigoplus_{i=1}^3 \mathbf{1}_i$ , and  $\Phi_i = \text{diag}(1, 2, \dots, n_i)$ .

If the model is estimated using inverse variance weights, so that  $\mathbf{W}_i = \text{diag}\left(1, \frac{1}{2}, \dots, \frac{1}{n_i}\right)$ , then the CR2 adjustment matrices differ depending on whether they are calculated from the full model or from the model after absorbing the fixed effects. Table 1 reports the product of the adjustment matrices and the absorbed design matrix for the weighted least squares estimator. The column labelled Full uses the adjustment matrices based on the full design (i.e.,  $\mathbf{A}_i \mathbf{W}_i \ddot{\mathbf{R}}_i$ , with  $\mathbf{A}_i$  calculated from Equation 4) or based on the absorbed design

Table 1: Adjustment matrices based on weighted or unweighted least squares, calculated with or without absorbing fixed effects

Cluster	Y	Weighted		Unweighted (Hom.)	Unweighted (Het.)	
		Full	Absorbed	Full/Absorbed	Full	Absorbed
A	1.6	-0.110	-0.342	-0.510	-0.334	-0.497
	4.1	0.441	0.345	0.510	0.669	0.504
B	2.6	-0.174	-0.689	-1.091	-0.658	-0.997
	1.0	0.409	0.203	0.000	0.358	0.000
	7.6	0.647	0.518	1.091	1.437	1.072
C	6.7	-0.353	-1.954	-4.472	-2.786	-4.089
	5.0	0.483	-0.176	-2.236	-0.770	-2.342
	3.1	0.926	0.532	0.000	1.757	0.000
	3.7	1.193	0.925	2.236	4.487	2.564
	5.8	1.372	1.178	4.472	7.318	5.236
$V^{CR2}$		0.828	1.019	1.173	1.248	1.050

(i.e.,  $\tilde{\mathbf{A}}_i \mathbf{W}_i \ddot{\mathbf{R}}_i$ , with  $\tilde{\mathbf{A}}_i$  calculated from Equation 6). The values differ, contradicting the original statement of Theorem 2. The final row of the table reports the value of  $V^{CR2}$  based on fitting the model using the outcomes reported in the second column of the table. The difference in adjustment matrices leads to differences in the value of the variance estimator.

## 4 A revised Theorem 2

The implication of the original Theorem 2 was that using the modified adjustment matrices  $\tilde{\mathbf{A}}_i$  to calculate the CR2 estimator yields the same result as using the full adjustment matrices  $\mathbf{A}_i$ . Although this does not hold under the general conditions given above, a modified version of the theorem does hold for the more limited case of ordinary (unweighted) least squares regression with a homoskedastic working model. The precise conditions are given in the following theorem, with proof given in Section 4.2.

**Theorem.** Let  $\mathbf{L}_i = (\ddot{\mathbf{U}}'\ddot{\mathbf{U}} - \ddot{\mathbf{U}}_i'\ddot{\mathbf{U}}_i)$  and assume that  $\mathbf{L}_1, \dots, \mathbf{L}_m$  have full rank  $r + s$ . If  $\mathbf{W}_i = \mathbf{I}_i$  and  $\Phi_i = \mathbf{I}_i$  for  $i = 1, \dots, m$  and  $\mathbf{T}_i\mathbf{T}_k' = \mathbf{0}$  for  $i \neq k$ , then  $\mathbf{A}_i\ddot{\mathbf{R}}_i = \tilde{\mathbf{A}}_i\ddot{\mathbf{R}}_i$ , where  $\mathbf{A}_i$  and  $\tilde{\mathbf{A}}_i$  are as defined in (4) and (6), respectively.

The implication of the revised theorem is that, for ordinary least squares regression with a homoskedastic working model, calculating the CR2 with the modified adjustment matrices  $\tilde{\mathbf{A}}_i$  leads to the same result as using the full adjustment matrices  $\mathbf{A}_i$ . The equality between the full and absorbed adjustment matrices does not hold for weighted or generalized least squares, nor for ordinary least squares with working models other than  $\Phi_i = \mathbf{I}_i$ .

Continuing the example described in the previous section, Table 1 reports the product of the adjustment matrices and the absorbed design matrix for the ordinary least squares estimator with a homoskedastic working model in the column labelled Unweighted (Hom.). The values based on the full design matrix  $(\mathbf{A}_i\ddot{\mathbf{R}}_i)$  are numerically identical to the values based on the absorbed design matrix  $(\tilde{\mathbf{A}}_i\ddot{\mathbf{R}}_i)$ . In the columns labeled Unweighted (Het.), the same quantities are computed using the heteroskedastic working model described in the previous section. The values based on the full design matrix differ from the values based on the absorbed design matrix, leading to differences in the cluster-robust variance estimator reported in the final row of the table.

## 4.1 Remarks

The revised version of Theorem 2 holds for the class of linear models estimated using ordinary least squares, with the cluster-robust variance estimator constructed based on a homoskedastic working model. Considering the ubiquity of ordinary least squares in applied data analysis, this is clearly an important class of models—perhaps even the most important for application. However, other important classes fall outside the scope of Theorem 2. For instance, an analyst might prefer to estimate a linear model using weighted least squares based on a posited heteroskedastic variance structure, paired with heteroskedasticity- or cluster-robust variance estimators to buttress against misspecification of that variance structure (Romano and Wolf 2017). In other applications, an analyst might use generalized estimating equation with a compound symmetric or auto-regressive working model for the errors (Wang and Carey 2003). In still other applications, analysts

might need to estimate a linear model using survey weights or inverse propensity weights, while maintaining a homoskedastic working model. For such applications, efficient computation of small-sample adjusted cluster-robust variance estimators remains a topic for further research.

## 4.2 Proof

Setting  $\Phi_i = \mathbf{I}_i$  and observing that  $\ddot{\mathbf{U}}_i' \mathbf{T}_i = \mathbf{0}$  for  $i = 1, \dots, m$ , it follows that

$$\begin{aligned} \mathbf{B}_i &= \mathbf{D}_i \mathbf{C}_i (\mathbf{I} - \mathbf{H}_{\ddot{\mathbf{U}}}) (\mathbf{I} - \mathbf{H}_{\mathbf{T}}) \Phi (\mathbf{I} - \mathbf{H}_{\mathbf{T}})' (\mathbf{I} - \mathbf{H}_{\ddot{\mathbf{U}}})' \mathbf{C}_i' \mathbf{D}_i' \\ &= \mathbf{C}_i (\mathbf{I} - \mathbf{H}_{\ddot{\mathbf{U}}} - \mathbf{H}_{\mathbf{T}}) (\mathbf{I} - \mathbf{H}_{\ddot{\mathbf{U}}} - \mathbf{H}_{\mathbf{T}}) \mathbf{C}_i' \\ &= (\mathbf{I}_i - \ddot{\mathbf{U}}_i \mathbf{M}_{\ddot{\mathbf{U}}} \ddot{\mathbf{U}}_i' - \mathbf{T}_i \mathbf{M}_{\mathbf{T}} \mathbf{T}_i') \end{aligned} \quad (7)$$

and similarly,

$$\tilde{\mathbf{B}}_i = (\mathbf{I}_i - \ddot{\mathbf{U}}_i \mathbf{M}_{\ddot{\mathbf{U}}} \ddot{\mathbf{U}}_i'). \quad (8)$$

It is apparent that  $\tilde{\mathbf{B}}_i \mathbf{T}_i = \mathbf{T}_i$ . We now show that  $\tilde{\mathbf{A}}_i \mathbf{T}_i = \mathbf{T}_i$  as well. Denote the rank of  $\ddot{\mathbf{U}}_i$  as  $u_i \leq \min\{n_i, r + s\}$  and take the thin QR decomposition of  $\ddot{\mathbf{U}}_i$  as  $\ddot{\mathbf{U}}_i = \mathbf{Q}_i \mathbf{R}_i$ , where  $\mathbf{Q}_i$  is an  $n_i \times u_i$  semi-orthonormal matrix and  $\mathbf{R}_i$  is a  $u_i \times r + s$  matrix of rank  $u_i$ , with  $\mathbf{Q}_i' \mathbf{Q}_i = \mathbf{I}$ . Note that  $\mathbf{Q}_i' \mathbf{T}_i = \mathbf{0}$ . From the observation that  $\tilde{\mathbf{B}}_i$  can be written as

$$\tilde{\mathbf{B}}_i = \mathbf{I}_i - \mathbf{Q}_i \mathbf{Q}_i' + \mathbf{Q}_i (\mathbf{I} - \mathbf{R}_i \mathbf{M}_{\ddot{\mathbf{U}}} \mathbf{R}_i') \mathbf{Q}_i',$$

it can be seen that

$$\tilde{\mathbf{A}}_i = \tilde{\mathbf{B}}_i^{+1/2} = \mathbf{I}_i - \mathbf{Q}_i \mathbf{Q}_i' + \mathbf{Q}_i (\mathbf{I} - \mathbf{R}_i \mathbf{M}_{\ddot{\mathbf{U}}} \mathbf{R}_i')^{+1/2} \mathbf{Q}_i'. \quad (9)$$

It follows that  $\tilde{\mathbf{A}}_i \mathbf{T}_i = \mathbf{T}_i$ .

Define the matrices

$$\mathbf{A}_i = \tilde{\mathbf{A}}_i - \mathbf{T}_i \mathbf{M}_{\mathbf{T}} \mathbf{T}_i', \quad (10)$$

for  $i = 1, \dots, m$ . We claim that  $\mathbf{A}_i = \mathbf{B}_i^{+1/2}$ . This can be seen by observing that  $\mathbf{A}_i \mathbf{A}_i$  is equal to  $\mathbf{B}_i^+$ , the Moore-Penrose inverse of  $\mathbf{B}_i$ . Note that

$$\mathbf{A}_i \mathbf{A}_i = \tilde{\mathbf{A}}_i \tilde{\mathbf{A}}_i - \mathbf{T}_i \mathbf{M}_{\mathbf{T}} \mathbf{T}_i' = \tilde{\mathbf{B}}_i^+ - \mathbf{T}_i \mathbf{M}_{\mathbf{T}} \mathbf{T}_i'. \quad (11)$$

It can then readily be verified that i)  $\mathbf{B}_i \mathbf{A}_i \mathbf{A}_i \mathbf{B}_i = \mathbf{B}_i$ , ii)  $\mathbf{A}_i \mathbf{A}_i \mathbf{B}_i \mathbf{A}_i \mathbf{A}_i = \mathbf{A}_i \mathbf{A}_i$ , and iii)  $\mathbf{A}_i \mathbf{A}_i \mathbf{B}_i = \mathbf{B}_i \mathbf{A}_i \mathbf{A}_i$ . Thus,  $\mathbf{A}_i \mathbf{A}_i$  satisfies the definition of the Moore-Penrose inverse of  $\mathbf{B}_i$  (Rao and Mitra 1971).

Finally, because  $\mathbf{T}_i' \ddot{\mathbf{R}}_i = \mathbf{0}$ , it can be seen that  $\mathbf{A}_i \ddot{\mathbf{R}}_i = \left( \tilde{\mathbf{A}}_i - \mathbf{T}_i \mathbf{M}_T \mathbf{T}_i' \right) \ddot{\mathbf{R}}_i = \tilde{\mathbf{A}}_i \ddot{\mathbf{R}}_i$ .

## References

- Bell, Robert M, and Daniel F McCaffrey. 2002. “Bias reduction in standard errors for linear regression with multi-stage samples.” *Survey Methodology* 28 (2): 169–81.
- McCaffrey, Daniel F, Robert M Bell, and Carsten H Botts. 2001. “Generalizations of biased reduced linearization.” In *Proceedings of the Annual Meeting of the American Statistical Association*. 1994.
- Pustejovsky, James E, and Elizabeth Tipton. 2018. “Small-Sample Methods for Cluster-Robust Variance Estimation and Hypothesis Testing in Fixed Effects Models.” *Journal of Business & Economic Statistics* 36 (4): 672–83. <https://doi.org/10.1080/07350015.2016.1247004>.
- Rao, C. Radhakrishna, and Sujit Kumar Mitra. 1971. *Generalized Inverse of Matrices and its Applications*. New York, NY: John Wiley & Sons.
- Romano, Joseph P., and Michael Wolf. 2017. “Resurrecting Weighted Least Squares.” *Journal of Econometrics* 197 (1): 1–19. <https://doi.org/10.1016/j.jeconom.2016.10.003>.
- Wang, Y.-G., and Vincent J Carey. 2003. “Working Correlation Structure Misspecification, Estimation and Covariate Design: Implications for Generalised Estimating Equations Performance.” *Biometrika* 90 (1): 29–41. <https://doi.org/10.1093/biomet/90.1.29>.