

Supplementary information for: Estimators for longitudinal latent exposure models: examining measurement model assumptions

Brisa N. Sánchez^{*,1}, Sehee Kim¹, Mary D. Sammel²

A. Score equations for MLE

Let $\theta = (\theta_0, \theta_1, \theta_2, \theta_3)$ represent all model parameters, with $\theta_0 = (\alpha_y, \gamma^\top, \sigma_\epsilon^2)^\top$, and $\theta_1 = (\alpha_{\xi_i}, \alpha_x)$, $\theta_2 = (\Lambda)$, and $\theta_3 = (\Sigma_\xi, \Sigma_\delta)$. Assuming normality of ϵ_i , δ_i , and ξ_i , the observed marginal likelihood is $L(\theta) = \prod_{i=1}^N f_{Y|X}(Y_i|X_i; \theta) f_X(X_i; \theta_1, \theta_2, \theta_3)$, where $f_{Y|X}$, f_X are normal densities. Letting $\ell(\theta) = \log L(\theta)$, and θ_{jr} represent the r^{th} unique parameter contained in θ_j , $j = 1, 2, 3$, the i^{th} 's subject contribution to the likelihood score equations is

$$\frac{\partial \ell_i}{\partial \alpha_y} = (Y_i - \mu_{Y|X}^i) / \sigma_{Y|X}^2 \quad (1)$$

$$\frac{\partial \ell_i}{\partial \gamma} = \widehat{\xi}_i^{\text{LV}} (Y_i - \mu_{Y|X}^i) / \sigma_{Y|X}^2 + \frac{\Sigma_{\xi|X} \gamma}{2(\sigma_{Y|X}^2)^2} [(Y_i - \mu_{Y|X}^i)^2 - \sigma_{Y|X}^2] \quad (2)$$

$$\frac{\partial \ell_i}{\partial \sigma_\epsilon^2} = \frac{1}{(\sigma_{Y|X}^2)^2} [(Y_i - \mu_{Y|X}^i)^2 - \sigma_{Y|X}^2] \quad (3)$$

$$\frac{\partial \ell_i}{\partial \theta_{1r}} = \gamma^\top \frac{\partial \widehat{\xi}_i^{\text{LV}}}{\partial \theta_{1r}} (Y_i - \mu_{Y|X}^i) / \sigma_{Y|X}^2 + \left(\frac{\partial \mu_x^i}{\partial \theta_{1r}} \right)^\top \Sigma_x^{-1} (X_i - \mu_x^i) \quad (4)$$

$$\frac{\partial \ell_i}{\partial \theta_{2r}} = \frac{1}{2(\sigma_{Y|X}^2)^2} [(Y_i - \mu_{Y|X}^i)^2 - \sigma_{Y|X}^2] + \gamma^\top \frac{\partial \widehat{\xi}_i^{\text{LV}}}{\partial \theta_{2r}} (Y_i - \mu_{Y|X}^i) / \sigma_{Y|X}^2 \quad (5)$$

$$+ \frac{1}{2} \text{tr} \left\{ \frac{\partial \Sigma_x}{\partial \theta_{2r}} \Sigma_x^{-1} [(X_i - \mu_x^i)(X_i - \mu_x^i)^\top - \Sigma_x] \Sigma_x^{-1} \right\} + \left(\frac{\partial \mu_x^i}{\partial \theta_{2r}} \right)^\top \Sigma_x^{-1} (X_i - \mu_x^i) \quad (6)$$

$$\frac{\partial \ell_i}{\partial \theta_{3r}} = \frac{1}{2(\sigma_{Y|X}^2)^2} [(Y_i - \mu_{Y|X}^i)^2 - \sigma_{Y|X}^2] + \gamma^\top \frac{\partial \widehat{\xi}_i^{\text{LV}}}{\partial \theta_{3r}} (Y_i - \mu_{Y|X}^i) / \sigma_{Y|X}^2 \quad (7)$$

$$+ \frac{1}{2} \text{tr} \left\{ \frac{\partial \Sigma_x}{\partial \theta_{3r}} \Sigma_x^{-1} [(X_i - \mu_x^i)(X_i - \mu_x^i)^\top - \Sigma_x] \Sigma_x^{-1} \right\} \quad (8)$$

¹Department of Biostatistics, University of Michigan, Ann Arbor, MI USA 48109

²Center for Clinical Epidemiology and Biostatistics, University of Pennsylvania School of Medicine, Philadelphia, USA

*Correspondence to: brisa@umich.edu, Department of Biostatistics, University of Michigan, Ann Arbor, MI USA 48109

where $\widehat{\xi}_i^{LV} = E(\xi_i|X_i) = \alpha_{\xi_i} + \Sigma_{\xi}\Lambda_x^{\top}\Sigma_x^{-1}(X_i - \mu_x^i)$ is the expected value of the latent variable ξ_i given the items; $\mu_{y|x}^i = \alpha_y + \gamma^{\top}\widehat{\xi}_i^{LV}$; $\sigma_{y|x}^2 = \gamma^{\top}\Sigma_{\xi|x}\gamma + \sigma_{\epsilon}^2$; $\Sigma_{\xi|x} = \text{var}(\xi_i|X_i) = \Sigma_{\xi} - \Sigma_{\xi}\Lambda_x^{\top}\Sigma_x^{-1}\Lambda_x\Sigma_{\xi}$; $\mu_x^i = E(X_i) = \alpha_x + \Lambda_x\alpha_{\xi_i}$; $\Sigma_x = \text{Var}(X_i) = \Lambda_x\Sigma_{\xi}\Lambda_x^{\top} + \Sigma_{\delta}$.

B. Estimating equations for Regression Calibration approach

The estimating equations for IV1 are:

$$S_{\alpha_y} = \sum_{i=1}^N (Y_i - \mu_{y,IV}^i) \quad (9)$$

$$S_{\gamma} = \sum_{i=1}^N \widehat{X}_{i1} (Y_i - \mu_{y,IV}^i) \quad (10)$$

$$S_{G^*} = \sum_{i=1}^N \mathbf{V}_i^{\top} (\widehat{X}_{i1} - \mathbf{V}_i G^*) \quad (11)$$

where $\mu_{y,IV}^i = \alpha_y + \gamma^{\top}\widehat{X}_{i1}$, $\widehat{X}_{i1} = (\mathbf{V}_i \widehat{G}^*)^{\top}$, G is the matrix of regression coefficients in the first stage regression, G^* is a vector, the stacked columns of G , $\mathbf{V}_i = I_{\ell} \otimes V_i$ where V_i is the i th row of \mathbf{V} , and I_{ℓ} is an identity matrix of dimension ℓ . The variance of $\widehat{\theta}_{IV} = (\widehat{\alpha}_y, \widehat{\gamma}^{\top}, (\widehat{G}^*)^{\top})^{\top}$ can be obtained from the corresponding diagonal elements of $\widehat{\text{var}}(\widehat{\theta}_{IV}) = B^{-1}AB^{-\top}$, where $A = 1/N \sum_{i=1}^N S_i S_i^{\top}$ and $B = 1/N \sum_{i=1}^N \partial S_i / \partial \theta_{IV}$, and $S_i = (S_{\alpha_y}, S_{\gamma}^{\top}, S_{G^*}^{\top})^{\top}$. In particular, variances for $\widehat{\gamma}$ obtained using the estimating equations approach can also incorporate the variability inherent in the estimation of \widehat{Z}_i .

C. Constraints on G

Enforcing constraints on G is straightforward by modifying the design matrix \mathbf{V}_i in (11). Specifically, let V_{i1} be the elements of V_i that have constrained coefficients, and V_{i2} be the remaining coefficients. Then, instead of having

$$\mathbf{V}_i = \begin{pmatrix} V_i & 0 & 0 \\ 0 & V_i & 0 \\ 0 & 0 & V_i \end{pmatrix}$$

as above, we have

$$\mathbf{V}_i^c = \begin{pmatrix} V_{i1} & V_{i2} & 0 & 0 \\ V_{i1} & 0 & V_{i2} & 0 \\ V_{i1} & 0 & 0 & V_{i2} \end{pmatrix}.$$

In this case, $G_c^* = (G_1^*, G_2^*)$ is the vector of coefficients with G_1^* being the coefficients for V_{i1} that are constrained to be equal across scaling items (e.g., block diagonal elements of the matrix G), and G_2^* being the unconstrained coefficients.

Alternative ways to impose restrictions on G is to regress scaling items only on items measured at the same visit (e.g., x_{i11} regressed only on x_{i12}, \dots, x_{i1k}). This effectively constraints G to be a block diagonal matrix. One can further constrain the blocks of G to be equal, similar to the time invariance assumption for Λ . However although these constraints seem intuitive, they yielded biased estimates when serial independence is present, since constraining G to be block diagonal effectively assumes items measured at other time points are marginally uncorrelated (simulations not shown).

Table S1. Descriptive Statistics for ELEMENT dataset on lead exposure during pregnancy and child mental development

	Time	N	Mean	Std
Child's mental development	24mpp	318	91.4	11.30
Mother's log2 Plasma lead	T1	153	3.8	0.96
	T2	169	3.4	0.86
	T3	157	3.5	0.79
Mother's whole blood lead Laboratory 1	T1	155	6.9	4.70
	T2	173	6.2	3.10
	T3	159	6.7	3.40
Laboratory 2	T1	172	7.6	4.60
	T2	198	6.6	3.20
	T3	304	6.9	3.60
Child's log2 Cord blood lead	Birth	238	2.1	0.90
Child's Blood Lead	24mpp	318	4.8	3.40
Covariates				
Mother's IQ	24mpp	341	89.6	18.60
Age	Scr	341	25.9	5.10
Child's gender	Birth	341	0.5	

Abbreviations: T*t*=Trimester *t*; *x*mpp=*x* months post partum;
Scr=Screening

Figure S1. Percent bias in $\hat{\gamma}_1$ (top), $\hat{\gamma}_2$ (middle), and bias in $\hat{\gamma}_3$ (bottom) for instrumental variable and maximum likelihood estimators with various working models(x-axis) when data are generated under a variety of true models (denoted by different symbols in the legend) and true $\gamma = (-2, -2, 0)$.

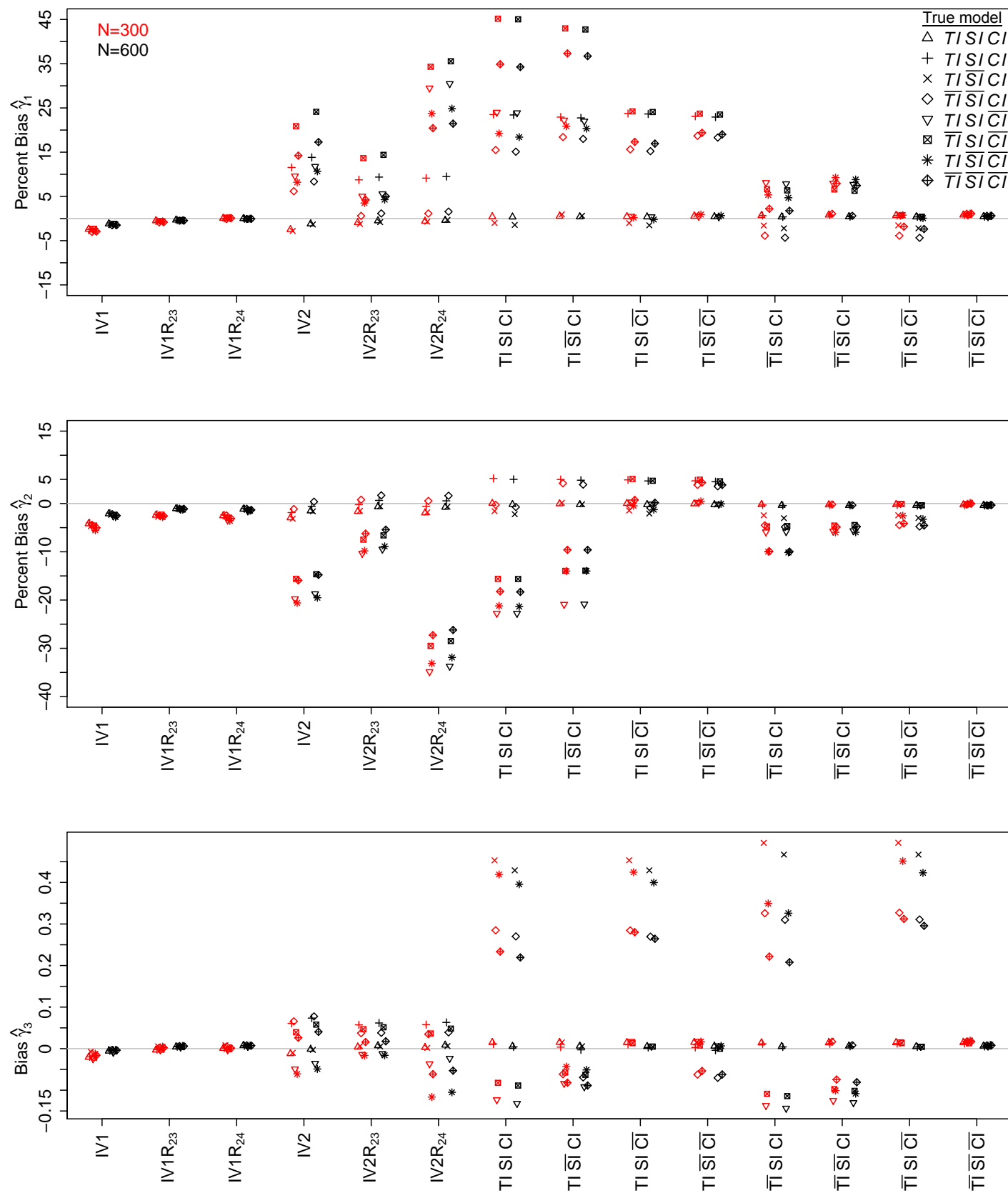


Figure S2. MSE of $\hat{\gamma}_1$ (top), $\hat{\gamma}_2$ (middle), $\hat{\gamma}_3$ (bottom) for instrumental variable and maximum likelihood estimators when data are generated from true $\gamma = (-2, -2, -2)$ and various degrees of lack of time invariance (left); various degrees of serial dependence (middle); various degrees of lack of conditional independence (right).

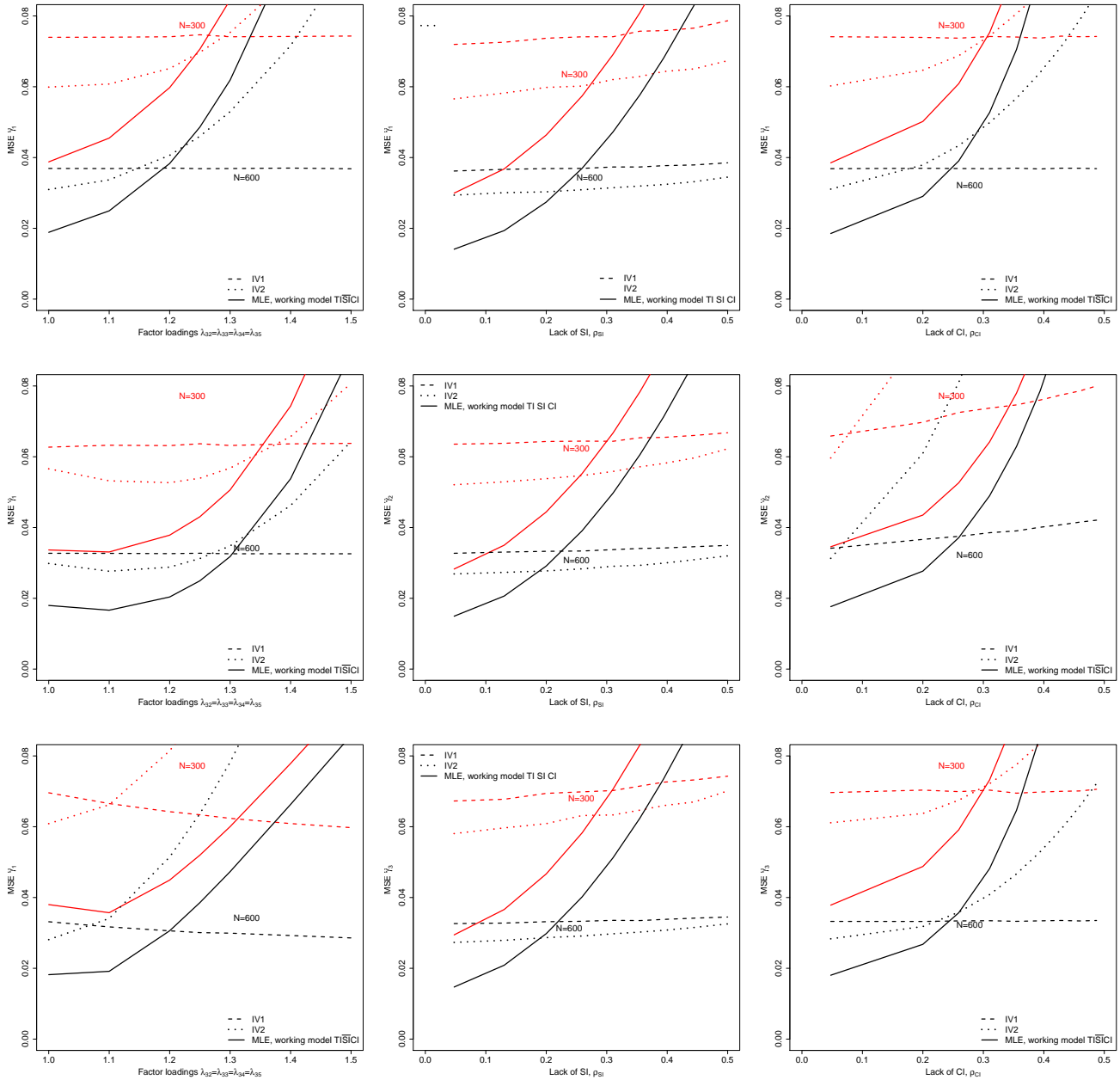


Figure S3. MSE of $\hat{\gamma}_1$ (top), $\hat{\gamma}_2$ (middle), $\hat{\gamma}_3$ (bottom) for instrumental variable and maximum likelihood estimators when data are generated from true $\gamma = (-2, -2, 0)$ and various degrees of lack of time invariance (left); various degrees of serial dependence (middle); various degrees of lack of conditional independence (right).

