

Multiple Imputation of Missing Data in Multilevel Designs: A Comparison of Different Strategies

Oliver Lüdtke, Alexander Robitzsch, and Simon Grund

Leibniz Institute for Science and Mathematics Education (IPN), Kiel, Germany, and Centre for International Student Assessment, Frankfurt, Germany

Multiple imputation is a widely recommended means of addressing the problem of missing data in psychological research. An often-neglected requirement of this approach is that the imputation model used to generate the imputed values must be at least as general as the analysis model. For multilevel designs in which lower level units (e.g., students) are nested within higher level units (e.g., classrooms), this means that the multilevel structure must be taken into account in the imputation model. In the present article, we compare different strategies for multiply imputing incomplete multilevel data using mathematical derivations and computer simulations. We show that ignoring the multilevel structure in the imputation may lead to substantial negative bias in estimates of intraclass correlations as well as biased estimates of regression coefficients in multilevel models. We also demonstrate that an ad hoc strategy that includes dummy indicators in the imputation model to represent the multilevel structure may be problematic under certain conditions (e.g., small groups, low intraclass correlations). Imputation based on a multivariate linear mixed effects model was the only strategy to produce valid inferences under most of the conditions investigated in the simulation study. Data from an educational psychology research project are also used to illustrate the impact of the various multiple imputation strategies.

Keywords: missing data, multiple imputation, multilevel modeling, multilevel data, intraclass correlation

Supplemental materials: <http://dx.doi.org/10.1037/met0000096.supp>

The pervasive problem of missing data has received considerable attention in psychological research during the last two decades (Enders, 2010; Graham, 2009; Schafer & Graham, 2002; see also West, 2001). There is consensus in the methodological literature that modern methods such as multiple imputation (MI) and model-based maximum likelihood procedures are much more effective at addressing missing data problems than traditional approaches such as listwise or pairwise deletion (Carpenter & Kenward, 2013; Little & Rubin, 2002). Although much has been published recently in the applied missing-data literature about these modern methods, less attention has been paid to the problem of missing values in multilevel designs. In such designs, lower level units (e.g., students, employees; Level 1) are typically nested within higher level units (e.g., classrooms, working units; Level 2). Multilevel modeling is a highly recommended statistical technique for analyzing these data structures, as it accounts for the dependence in the data as well as allowing researchers to estimate

relationships among variables located at different levels (Goldstein, 2010; Raudenbush & Bryk, 2002; Snijders & Bosker, 2012).

The purpose of this article is to evaluate several strategies for applying MI to incomplete multilevel data. The basic idea of MI is to draw a number of replacements for the missing values using the observed data and an imputation model (Rubin, 1987). One central feature of MI, however, is that the imputation model must be at least as general as the model of interest in order to preserve the relationships among the variables. For example, assume that the researcher is interested in testing an interaction effect of two variables in a multiple regression model in the presence of incomplete data. In that case, it is crucial that the interaction effect (i.e., product term) is also incorporated in the imputation model (e.g., Enders, Baraldi, & Cham, 2014; von Hippel, 2009).

Similarly, in the case of incomplete multilevel data, it is important to take the multilevel structure of the data into account in the imputation model in order to ensure valid statistical inferences in subsequent multilevel analyses (Enders, Mistler, & Keller, 2016). In the present article, we investigate how MI conducted with a single-level normal imputation model affects the estimation of variance components and regression coefficients in a multilevel analysis. We also discuss a strategy that is based on including dummy indicator (DI) variables in order to preserve the multilevel structure in a single-level imputation model. These ad hoc procedures will be compared with a multivariate linear mixed-effects imputation model that was developed by Schafer (Schafer, 2001; Schafer & Yucel, 2002).

Our article makes three main contributions to the literature. First, in contrast to previous research that mostly relied on simu-

This article was published Online First September 8, 2016.

Oliver Lüdtke, Alexander Robitzsch, and Simon Grund, Leibniz Institute for Science and Mathematics Education (IPN), Kiel, Germany and Centre for International Student Assessment, Frankfurt, Germany.

Data for the example were provided by the Research Data Centre at the Institute for Educational Quality Improvement, Berlin, and collected under the direction of Wilfried Bos.

Correspondence concerning this article should be addressed to Oliver Lüdtke, Leibniz Institute for Science and Mathematics Education (IPN), Olshausenstraße 62, 24118 Kiel, Germany. E-mail: oluedtke@ipn.uni-kiel.de

lations to illustrate the problems of not adequately accommodating the multilevel structure in the imputation model, we derive the asymptotic bias for the estimators of the intraclass correlation and the regression coefficients of a multilevel random-intercept model, when the two ad hoc strategies are used to deal with incomplete data. Second, we conducted a comprehensive simulation study that provides a thorough comparison of the ad hoc procedures with the multivariate linear mixed-effects imputation model, manipulating more factors in the simulation design than most previous studies. Third, our focus is on a multilevel random-intercept model in which the between part of the Level 1 predictor is treated as a latent variable. This model has been recommended in the methodological literature for assessing the group-level effects of Level 1 predictors in contextual studies (e.g., Croon & van Veldhoven, 2007; Lüdtke et al., 2008; Preacher, Zyphur, & Zhang, 2010; Shin & Raudenbush, 2010).

The article is organized as follows. We start by briefly describing the missing data mechanisms as defined by Rubin (1976) and introducing the basic idea of MI. We then describe the multilevel random-intercept model and motivate the analysis models we are interested in. We discuss two ad hoc procedures that have been used for imputing multilevel missing data and analytically investigate their asymptotic bias. A discussion of the multivariate linear mixed effects imputation model follows. We then use simulation methods to examine different strategies for dealing with incomplete multilevel data. Next, an empirical example from educational psychology is used to illustrate the impact of these strategies on the estimation of an intraclass correlation. Finally, we offer suggestions for applied researchers and propose directions for further research.

Missing Data and Multiple Imputation

In his well-known classification of missing data, Rubin (1976) distinguished three mechanisms. Suppose one has a complete data matrix, which can be decomposed into observed and unobserved parts $\mathbf{Y} = (\mathbf{Y}_O, \mathbf{Y}_M)$ by an indicator matrix $\mathbf{R} = (r_{iv})$ denoting the missing data such that $r_{iv} = 1$ if the variable v for person i is observed and $r_{iv} = 0$ if it is missing. If values are missing as a random sample of the complete hypothetical data, that is, if $P(\mathbf{R}|\mathbf{Y}) = P(\mathbf{R})$, the data are missing completely at random (MCAR). If the missingness depends on other variables but the data are MCAR when such variables are partialled out, that is, if $P(\mathbf{R}|\mathbf{Y}) = P(\mathbf{R}|\mathbf{Y}_O)$, the data are missing at random (MAR). This is in contrast to data that are missing *not* at random (MNAR), where missingness is also dependent on the missing part of the data, that is, $P(\mathbf{R}|\mathbf{Y}) = P(\mathbf{R}|\mathbf{Y}_O, \mathbf{Y}_M)$. Most software implementations of MI rely on the assumption that missing data are MAR (see Carpenter & Kenward, 2013, for a discussion of MI under a MNAR mechanism). This is a reasonable assumption that holds at least approximately if the observed data provide sufficient information about the missing data mechanism (Collins, Schafer, & Kam, 2001).

The MI procedure consists of three steps (see Enders, 2010, for a clear exposition). In the imputation phase, m copies of the data set are generated by filling in replacements for the missing values. In the analysis phase, the m completed data sets are then analyzed using standard complete-data methods. In the pooling phase, the

parameter estimates are pooled according to the rules described by Rubin (1987) for final parameter estimates and inference.

In our discussion of MI strategies for incomplete multilevel data, we are particularly concerned with the imputation phase. The key idea is to draw replacements of the missing values from the conditional distribution $P(\mathbf{Y}_M|\mathbf{Y}_O)$ of the missing data, given the observed data. In the Bayesian context, this distribution is also called the posterior predictive distribution of the missing data given the observed data (Gelman, Carlin, Stern, & Rubin, 2003; Hoff, 2009). To generate the imputed values, it is necessary to specify a joint distribution $P(\mathbf{Y}_M, \mathbf{Y}_O, \boldsymbol{\theta})$ for the missing and observed data. In research practice, the multivariate normal distribution with $\boldsymbol{\theta} = (\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is often used as an imputation model, where $\boldsymbol{\mu}$ is a vector of means and $\boldsymbol{\Sigma}$ is a covariance matrix.¹

In practice, the posterior predictive distribution of the missing values $P(\mathbf{Y}_M|\mathbf{Y}_O)$ is difficult to evaluate; Markov chain Monte Carlo (MCMC) techniques are generally used to draw values from this distribution. One commonly used MCMC procedure (also called data augmentation; Tanner & Wong, 1987) uses the following stochastic iterative algorithm, which cycles between two consecutive steps. At the $(t + 1)$ th iteration with current values $(\mathbf{Y}_M^{(t)}, \mathbf{Y}_O, \boldsymbol{\theta}^{(t)})$, the imputation step (I-step) draws missing values from the conditional distribution of the missing values given the observed data

$$\mathbf{Y}_M^{(t+1)} \sim P(\mathbf{Y}_M|\mathbf{Y}_O, \boldsymbol{\theta}^{(t)}). \quad (1)$$

In the next step, the posterior step (P-step), the completed data $(\mathbf{Y}_M^{(t+1)}, \mathbf{Y}_O)$ are used to generate new values for the parameter vector $\boldsymbol{\theta}$

$$\boldsymbol{\theta}^{(t+1)} \sim P(\boldsymbol{\theta}|\mathbf{Y}_O, \mathbf{Y}_M^{(t+1)}). \quad (2)$$

When the algorithm converges, these values can be viewed as simulated draws from the posterior distribution of the parameters, given observed and filled-in data. Typically, the initial samples of the algorithm are discarded (burn-in period) because the initial draws are affected by the starting values (for a discussion of assessing convergence in the context of MI, see Schafer & Olsen, 1998; see also Enders, 2010).

The crucial decision for our discussion of MI strategies is the choice of an imputation model. In the present article, we show that ignoring the multilevel structure can result in distorted parameter estimates in subsequent multilevel analyses. In the next section, we introduce the specific multilevel models we are interested in.

Multilevel Models With Missing Data

In the following, we consider a scenario with two variables X and Y , where Y has missing values and X is fully observed. More specifically, we assume a two-level structure with two individual-level variables X_{ij} and Y_{ij} for persons i ($i = 1, \dots, n_j$) in groups j ($j = 1, \dots, K$). The variables X_{ij} and Y_{ij} are decomposed as follows (see Snijders & Bosker, 2012, p. 29):

¹ Two broad approaches to performing MI can be distinguished. In the *joint modeling* approach, a single statistical model is used for incomplete variables simultaneously. In the *sequential* (or chained equations) approach, each variable is imputed in turn using a sequence of models. In the present article, we focus on the joint modeling approach (see Carpenter & Kenward, 2013, for a discussion).

$$X_{ij} = \mu_X + X_{B,j} + X_{W,ij}, \quad Y_{ij} = \mu_Y + Y_{B,j} + Y_{W,ij}. \quad (3)$$

In this model, group j has specific main effects $X_{B,j}$ and $Y_{B,j}$ for variables X and Y , and the within-group deviations $X_{W,ij}$ and $Y_{W,ij}$ are associated with individual i . The covariance matrix of X and Y within and between groups can be written as

$$\Sigma_W = \begin{pmatrix} \sigma_X^2 & \rho_W \sigma_X \sigma_Y \\ \rho_W \sigma_X \sigma_Y & \sigma_Y^2 \end{pmatrix} \text{ and } \Sigma_B = \begin{pmatrix} \tau_X^2 & \rho_B \tau_X \tau_Y \\ \rho_B \tau_X \tau_Y & \tau_Y^2 \end{pmatrix}, \quad (4)$$

where ρ_W and ρ_B denote the correlation between the individual deviations $\rho(X_{W,ij}, Y_{W,ij})$ and the between-groups correlation $\rho(X_{B,j}, Y_{B,j})$, respectively. The intraclass correlations $\rho_{I,X}$ and $\rho_{I,Y}$ indicate the proportion of total variance that can be attributed to between-groups differences and are defined as follows:

$$\rho_{I,X} = \frac{\tau_X^2}{\tau_X^2 + \sigma_X^2} \text{ and } \rho_{I,Y} = \frac{\tau_Y^2}{\tau_Y^2 + \sigma_Y^2}. \quad (5)$$

The intraclass correlation of a variable provides important information about the multilevel structure, and its calculation is usually the first step in analyzing multilevel data. We are interested in how different strategies of dealing with incomplete multilevel data affect the estimation of the intraclass correlation of Y .

Furthermore, we are interested in the relationship between X and Y within and between groups (see Cronbach, 1976). A multilevel random-intercept model is used, in which the dependent variable Y is predicted by

$$Y_{ij} = \mu_Y + \beta_{W,YX} X_{W,ij} + \beta_{B,YX} X_{B,j} + \delta_j + \epsilon_{ij}, \quad (6)$$

where μ_Y is the regression intercept, $\beta_{W,YX}$ is the within-group (Level 1) regression slope describing the relationship between Y and X within groups, and $\beta_{B,YX}$ is the between-groups (Level 2) regression slope that reflects the relationships between the group means of Y and X . The group-level residual δ_j and the individual-level residual ϵ_{ij} are normally distributed with zero means. It should also be noted that the model in Equation 6 treats the group mean of the predictor variable X as a latent variable and corrects the group-level effect $\beta_{B,YX}$ for the unreliability of the manifest, observed group mean (e.g., Lüdtke et al., 2008; see also Croon & van Veldhoven, 2007; Shin & Raudenbush, 2010). This model has been used to estimate the individual and group-level effects of Level 1 predictors by researchers in various subdisciplines of psychology, such as educational psychology (e.g., Dettmers, Trautwein, Lüdtke, Kunter, & Baumert, 2010), health psychology (e.g., Henry, Stanley, Edwards, Harkabus, & Chapin, 2009), and organizational psychology (e.g., Walsh, Matthews, Tuller, Parks, & McDonald, 2010).

Additionally, researchers are often interested in estimating contextual effects (Raudenbush & Bryk, 2002). A contextual effect is present if $\beta_{B,YX}$ is different from $\beta_{W,YX}$, meaning that the relationship at the aggregated level (Level 2) is stronger or weaker than the relationship at the individual level (Level 1). Contextual effects are of great interest in educational psychology, for example, where several researchers have postulated that aggregated school socioeconomic status or mean ability has an effect on student outcomes (e.g., student achievement or academic achievement), even after controlling for the individual effects of the constructs at Level 1. Another important aspect of the existence of a contextual effect becomes evident when we write the regression coefficient $\beta_{total,YX}$

of Y on X in a single-level analysis (i.e., ignoring the clustering of persons into groups) as a function of $\beta_{B,YX}$ and $\beta_{W,YX}$ (Snijders & Bosker, 2012, p. 30):

$$\beta_{total,YX} = \rho_{I,X} \beta_{B,YX} + (1 - \rho_{I,X}) \beta_{W,YX}. \quad (7)$$

Thus, the total regression coefficient $\beta_{total,YX}$ in the regression of Y on X is a weighted mean of the within- and between-groups regression coefficients. If no contextual effect is present ($\beta_{B,YX} = \beta_{W,YX}$), the within-group and between-groups coefficients are equal to the total regression coefficient from the single-level analysis. This relationship will be of relevance for the bias derivations in the next section.

Alternatively, we also investigate the reversed relationship when X is the outcome variable and Y , the variable with missing values, is the predictor

$$X_{ij} = \mu_X + \beta_{W,XY} Y_{W,ij} + \beta_{B,XY} Y_{B,j} + \delta_j + \epsilon_{ij}, \quad (8)$$

where μ_X is the regression intercept, $\beta_{W,XY}$ is the within-group regression slope, and $\beta_{B,XY}$ is the between-groups regression slope. The assumptions about the residuals δ_j and ϵ_{ij} are the same as above. It is well known with regard to multiple regression models that missing values in the predictor variables can be more problematic than missing values in the dependent variable (see Carpenter & Kenward, 2013, p. 24, for a detailed discussion).

Two Ad Hoc Strategies for Dealing With Multilevel Missing Data

In this section, we discuss two strategies that have been used to deal with incomplete multilevel data. The first approach uses a single-level imputation model and ignores the multilevel structure of the data. The second approach includes a set of DI variables to represent the multilevel structure in the single-level imputation model.

Ignoring the Multilevel Structure in the Imputation Model

In research, the multivariate normal distribution with $\theta = (\mu, \Sigma)$ is often selected as an imputation model. To illustrate the multivariate normal imputation model and demonstrate how it ignores the multilevel structure, let us use our scenario with two variables, X and Y . The variable X is assumed to be fully observed and Y is missing for a subset of cases. For example, the i th person has the following data pattern²: $Y_{i(mis)}$ and $X_{i(mis)}$, where the subscript *mis* indicates cases for which Y is missing (see Drechsler, 2015). A replacement for a missing value on Y would then be generated by the following equation:

$$Y_{i(mis)}^{(t+1)} = \alpha^{(t)} + \beta_1^{(t)} X_{i(mis)} + \epsilon_i, \quad (9)$$

where the residual ϵ_i is normally and identically distributed across persons with constant variance $\sigma^{2(t)}$. The regression parameters

² For simplicity's sake, we use an example with two variables and only a single missing data pattern. A multivariate regression would be required if there were more than two variables and more than one missing observation for an individual (for an example, see Hoff, 2009, p. 119; Enders, 2010, p. 200).

$(\alpha^{(r)}, \beta_1^{(r)}, \sigma^{2(r)})$ are based on posterior draws of $\theta^{(r)} = (\boldsymbol{\mu}^{(r)}, \boldsymbol{\Sigma}^{(r)})$ from the previous P-step. It is evident that generating imputations using Equation 9 would not take into account a multilevel structure and that the dependencies of the data are not adequately represented in the imputed values. Furthermore, if a contextual effect is present (i.e., if the relationship between X and Y at the group level differs from the relationship within groups), the expectation of the regression coefficient β_1 equals the total regression coefficient $\beta_{total,YX}$ and will be a weighted average of the within- and between-groups regression coefficients (see Equation 7), which does not adequately represent the relationships at the various analysis levels. Thus, if the model of interest is a multilevel model, important relationships among the variables may be omitted from the imputation model, increasing the risk of distorted parameter estimates in subsequent multilevel analyses that are based on the filled-in data. In the following, we refer to the MI strategy that ignores the multilevel structure and specifies a single-level multivariate normal distribution for the imputations as the NORM approach. Previous research has shown with simulation studies that using the NORM approach for imputing incomplete multilevel data produces intraclass correlations that underestimate their true size (e.g., Black, Harel, & McCoach, 2011; Taljaard, Donner, & Klar, 2008; van Buuren, 2011; see also Snijders & Bosker, 2012). In the next section, we discuss an ad hoc procedure for incorporating group effects in a single-level imputation model.

DI Approach

In the DI approach, a set of dummy variables is created to represent the multilevel structure. The dummy variables are included in the single-level imputation model and a separate intercept (or fixed effect) is estimated for each group.³ Group effects are thereby incorporated in the imputation model. More specifically, in the DI approach, the K groups are represented by $K-1$ dummy variables—or K indicator variables when the overall intercept is excluded (see Allison, 2009). To illustrate how imputations are generated using this strategy, we return to our example with two variables X and Y . Assuming now that a two-level structure should be represented in the imputation model by adding separated intercepts (or fixed effects) for each group, the linear regression model for imputing Y is written as follows:

$$Y_{ij(mis)}^{(t+1)} = \sum_{c=1}^K \alpha_c^{(t)} I(c=j) + \beta_1^{(t)} X_{ij(mis)} + \varepsilon_{ij}, \quad (10)$$

where $I(\cdot)$ denotes an indicator function that takes on the value 1 when a person belongs to a group and 0 otherwise. The regression parameters are again based on posterior draws of $\theta^{(r)} = (\boldsymbol{\mu}^{(r)}, \boldsymbol{\Sigma}^{(r)})$ from the previous P-step. It can be shown that the expectation of the coefficient β_1 in the DI approach is the within-group coefficient $\beta_{w,YX}$, which describes the relationship between X and Y within groups (see Equation 6). The DI approach has been supported by Graham (2009; White, Royston, & Wood, 2011; see Graham, 2012, for a less positive view) when the model of interest is a random-intercept model and the number of groups is not too large. Andridge (2011) took a critical look at the DI approach in the context of cluster randomized trials and showed that it results in biased standard errors for the regression coefficients (see also van Buuren, 2011). In a recent evaluation of the DI approach,

Drechsler (2015) demonstrated analytically and through simulations that unless the missing data rate is large ($>10\%$), and/or the intraclass correlation is small ($<.10$) and the number of persons per group is small, the DI method produced approximately unbiased estimates of regression coefficients and their standard errors if the model of interest is a multilevel random-intercept model. However, that evaluation considered only the case of missing values on the dependent variable and complete data on the predictors. The DI approach might be expected to be more problematic when there are missing values on the predictor variables (see also Enders et al., 2016, for a critical discussion of the DI approach). In the next section, we show how the two ad hoc procedures (NORM and DI approach) can result in biased estimators of variance components and multilevel regression coefficients for incomplete multilevel data.

Asymptotic Bias for the Two Ad Hoc MI Strategies

For our scenario with two variables (X and Y), we now investigate how the treatment of missing values in the NORM and DI approach affects the parameter estimates in subsequent multilevel analyses. For the following derivations, it is assumed that the number of groups approaches infinity. Furthermore, we assume that the values in Y are MCAR and that the variable X is fully observed. Without loss of generality, both variables are assumed to be mean-centered (i.e., zero mean in the population). In each group, n_1 persons have observed values, and $n_0 = n - n_1$ values on Y are missing. For simplicity's sake, it is assumed that the missing data rate $p_0 = n_0/n$ is the same in each group.

We focus on two analysis models. First, we are interested in the intraclass correlation of Y . Second, we investigate the within- and between-groups regression coefficients for the regression of Y on X (i.e., $\beta_{w,YX}$ and $\beta_{B,YX}$; see Equation 6) as well as for the regression of X on Y (i.e., $\beta_{w,XY}$ and $\beta_{B,XY}$; see Equation 8). Again, note that we assume that the number of groups approaches infinity ($K \rightarrow \infty$). The details of the derivations are presented in the Appendix.

Intraclass Correlation of Y

In the following, we derive the asymptotic bias for the estimator of the intraclass correlation ρ_{LY} . In a first step, we investigate the bias for the estimators of the between-groups variance τ_Y^2 and the within-group variance σ_Y^2 .

NORM approach. If the multilevel structure is ignored in the imputation model, the asymptotic bias of the estimator of the between-groups variance τ_Y^2 can be expressed as follows:

$$\begin{aligned} \text{Bias}(\hat{\tau}_Y^2) = & -p_0 \cdot \tau_Y^2 \cdot \frac{n}{n-1} \cdot \{ \rho_B^2(1 - \rho_{LX})(\beta_{B,YX} - \beta_{w,YX})A_X \\ & + (1 - \rho_B^2)A_e \}, \end{aligned} \quad (11)$$

where the two terms A_X and A_e are introduced to simplify the expression. They are defined as $A_X \equiv \{2(1 - p_0) + (p_0 - 1/n)(\beta_{total,YX}/\beta_{B,YX} + 1)\}/\beta_{B,YX}$ and $A_e \equiv 2 - 1/n - p_0$. The first

³This approach is also sometimes called a fixed effects approach (Drechsler, 2015). Fixed effects models can be used to assess causal effects in the presence of unknown group-level confounders (see Allison, 2009).

part of Equation 11 shows that the bias becomes larger when the missing data rate per group p_0 rises and the between-groups variance τ_y^2 increases. The second part (involving the term A_X) shows that if a contextual effect exists in the regression of Y on X (i.e., $\beta_{B,YX} \neq \beta_{W,YX}$), the bias depends on both the between-groups correlation and the intraclass correlation of X . For example, if a positive contextual effect is present (i.e., $\beta_{B,YX} > \beta_{W,YX}$) and the relationship between Y and X is stronger at the group level than within groups, the estimator of the between-groups variance is negatively biased and the magnitude of the between-groups variance will be underestimated. The third part (involving the term A_e) indicates that the bias decreases when X and Y are more strongly correlated at the group level. The relationships among the bias, the missing data rate (25% and 50%), the group size, and the correlation at the group level ($\rho_B = .30$ and $.60$) are also depicted in Figure 1. As we see, the bias grows larger with a higher missing data rate and decreases with an increasing correlation at the group level. Furthermore, increasing group size has almost no effect on the bias.

Interestingly, the positive bias of the estimator of the within-group variance σ_y^2 is equal to the negative bias of the estimator of the between-groups variance,

$$Bias(\hat{\sigma}_y^2) = -Bias(\hat{\tau}_y^2). \tag{12}$$

This relationship between the biases of the two estimators can be explained by the fact that the NORM approach preserves the total variance of Y . Using Equations 11 and 12, the bias for the estimator of the intraclass correlation of Y is calculated as follows:

$$Bias(\hat{\rho}_{I,Y}) = \rho_{I,Y} \cdot \frac{Bias(\hat{\tau}_y^2)}{\tau_y^2}. \tag{13}$$

As we see, the absolute bias depends on the true size of the intraclass correlation as well as the bias of the estimator of the between-groups variance of Y . Thus, it can be concluded that the bias grows larger with an increase in true intraclass correlation and a higher missing data rate. Furthermore, increasing the group size has only a minimal effect on the bias. Even with very large groups, the NORM approach can be expected to produce substantially negatively biased estimates of the intraclass correlation.

DI approach. In the DI approach, the multilevel structure is taken into account by including an indicator variable for each group in the imputation model. The within-group variance σ_y^2 can be estimated without bias using the DI approach. However, the

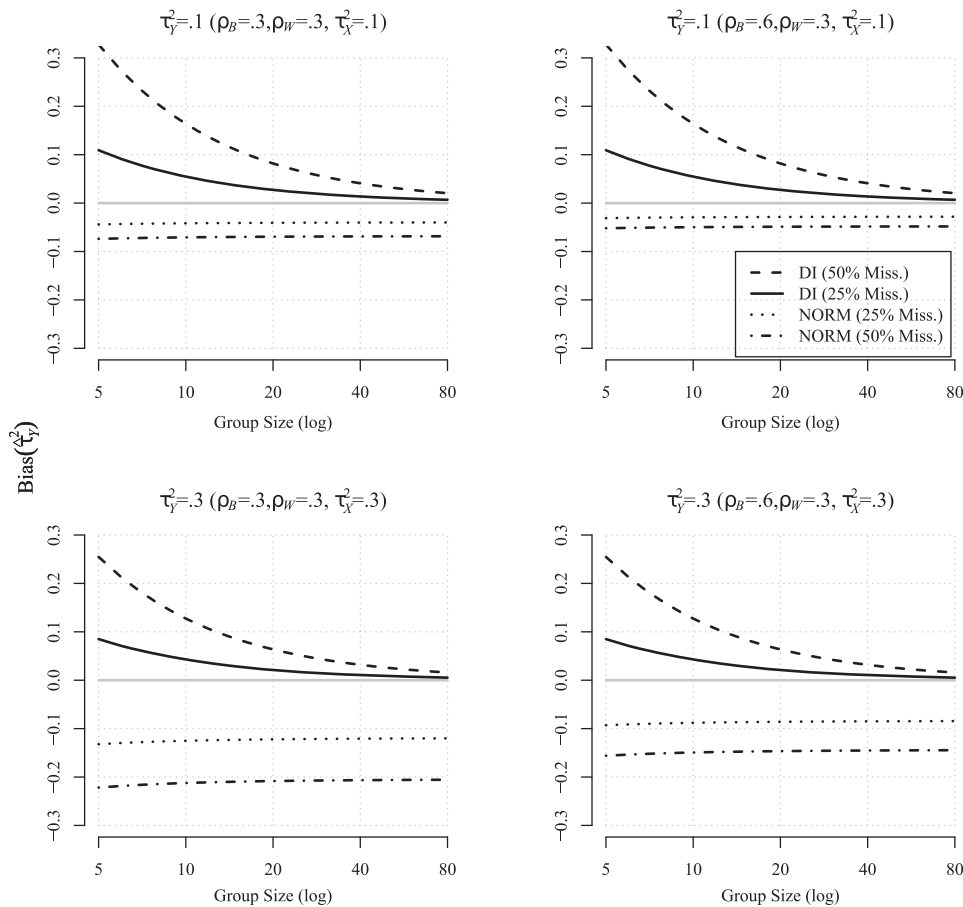


Figure 1. Asymptotic bias of the estimator of the between-groups variance τ_y^2 as a function of the logarithm of group size, missing data rate in groups, and multiple imputation strategy. Both variables X and Y are standardized with unit variance. It is assumed that the number of groups approaches infinity. NORM = normal model imputation; DI = dummy-indicator approach.

following relationship holds for the estimator of the between-groups variance $\hat{\tau}_Y^2$:

$$\text{Bias}(\hat{\tau}_Y^2) = \frac{\sigma_Y^2}{n} \cdot \frac{2p_0}{1-p_0} \cdot (1 - \rho_W^2). \quad (14)$$

As we see, the estimator of the between-groups variance is positively biased. The bias increases with an increase in the rate of missing data and as the group sizes decrease (see Figure 1). At the same time, the absolute bias also increases when the within-group variance grows larger. An intuitive explanation for the positively biased estimator of the between-groups variance is that the observed group means are an unreliable measure of the true group means when the group size is small and the within-group variance is large (Bliese, 2000). Thus, the DI approach, which estimates a separate intercept for each group, artificially inflates the true variation between groups. Graham (2012, p. 136) argues that the “dummy-coding strategy overcompensates” for the group structure. The bias also decreases when the within-group correlation (ρ_W) between Y and the fully observed variable X increases.

The positively biased estimator of the between-groups variance also affects the bias for the estimator of the intraclass correlation:

$$\text{Bias}(\hat{\rho}_{LY}) = (1 - \rho_{LY})^2 \cdot \frac{1}{n} \cdot \frac{2p_0}{1-p_0} \cdot (1 - \rho_W^2). \quad (15)$$

The bias for the intraclass correlation is positive, and a larger intraclass correlation as well as larger groups lead to a smaller absolute bias. As expected, the fraction of missing data within a group has a strong positive effect on the bias. Again, the bias decreases when the within-group correlation between X and Y increases.

Within- and Between-Group Regression Coefficients

In the following section, we investigate the asymptotic bias for the within- and between-groups regression coefficients that are estimated by a multilevel regression of Y on X ($\beta_{W,YX}$ and $\beta_{B,YX}$) and of X on Y ($\beta_{W,XY}$ and $\beta_{B,XY}$). As the estimators of these coefficients involve the within- and between-groups covariance of X and Y , we also investigate the estimator of the within-group covariance ($\sigma_{W,XY}$) and between-groups covariance ($\sigma_{B,XY}$).

NORM approach. If the multilevel structure is ignored in the imputation model, the asymptotic bias of the estimator of the within-group covariance $\sigma_{W,XY}$ can be expressed as follows:

$$\text{Bias}(\hat{\sigma}_{W,XY}) = p_0 \cdot \rho_{LX} \cdot (\beta_{B,YX} - \beta_{W,YX}) \cdot \sigma_X^2. \quad (16)$$

It is apparent that the bias depends on the rate of missing data as well as on the intraclass correlation and the within-group variance of X . More importantly, the estimator of the within-group covariance is only biased if a contextual effect exists in the population (i.e., $\beta_{B,YX} \neq \beta_{W,YX}$). In this case $\beta_{\text{total},XY}$ is different from $\beta_{B,XY}$ and $\beta_{W,XY}$ (see Equation 7) and the direction of the bias depends on whether a positive (i.e., $\beta_{B,YX} > \beta_{W,YX}$) or a negative (i.e., $\beta_{B,YX} < \beta_{W,YX}$) contextual effect of X on Y exists. In the case of a contextual effect in the regression of Y on X , therefore, the estimator of the within-group coefficient $\beta_{W,YX}$ is biased in the NORM approach:

$$\text{Bias}(\hat{\beta}_{W,YX}) = p_0 \cdot \rho_{LX} \cdot (\beta_{B,YX} - \beta_{W,YX}). \quad (17)$$

Similarly, the bias of the estimator of the within-group coefficient $\beta_{W,XY}$ in the regression of X on Y (i.e., when the predictor Y is MCAR) is given by

$$\text{Bias}(\hat{\beta}_{W,XY}) = \frac{\text{Bias}(\hat{\sigma}_{W,XY}) - \beta_{W,XY} \text{Bias}(\hat{\sigma}_Y^2)}{\sigma_Y^2 + \text{Bias}(\hat{\sigma}_Y^2)}. \quad (18)$$

As we see, the bias depends on the biases of the estimator of the within-group covariance (see Equation 16) and the estimator of the within-group variance of Y (see Equation 12).

The following relationship holds for the bias of the estimator of the between-groups covariance:

$$\text{Bias}(\hat{\sigma}_{B,XY}) = -p_0 \cdot (1 - \rho_{LX}) \cdot (\beta_{B,YX} - \beta_{W,YX}) \cdot \tau_X^2. \quad (19)$$

This relationship indicates that the bias depends on the fraction of missing data as well as on the intraclass correlation and the between-groups variance of X . Again, the presence of a contextual effect in the population is crucial for the existence and direction of the bias. For example, in case of a positive contextual effect, the difference $\beta_{B,YX} - \beta_{W,YX}$ is positive and the absolute magnitude of the between-groups covariance will be underestimated. If the estimator of the between-groups covariance is biased, the estimator of the between-groups regression coefficient $\beta_{B,XY}$ would also be biased:

$$\text{Bias}(\hat{\beta}_{B,XY}) = -p_0 \cdot (1 - \rho_{LX}) \cdot (\beta_{B,YX} - \beta_{W,YX}). \quad (20)$$

Based on Equations 11 and 19, the bias of the estimator for the between-groups regression coefficient $\beta_{B,XY}$ can be written as follows:

$$\text{Bias}(\hat{\beta}_{B,XY}) = \frac{\text{Bias}(\hat{\sigma}_{B,XY}) - \beta_{B,XY} \text{Bias}(\hat{\tau}_Y^2)}{\tau_Y^2 + \text{Bias}(\hat{\tau}_Y^2)}. \quad (21)$$

As we see, the bias depends on the bias of the estimator of the between-groups covariance and the estimator of the between-groups variance of Y . If no contextual effect is present, the bias depends primarily on the bias of the estimator for the between-groups variance because $\text{Bias}(\hat{\sigma}_{B,XY}) = 0$. Figure 2 shows that the bias grows larger with a higher missing data rate, and increasing the group size has only a very modest effect on the bias. However, raising the between-groups correlation (from .30 to .60) reduces the bias of the estimator of the between-groups coefficient. Interestingly, this positive effect of the larger between-groups correlation outweighs the bias in estimating the between-groups covariance that is introduced by the presence of a contextual effect. Overall, Equation 21 indicates that in the case of missing data in the predictor variable, the bias of the estimator of the between-groups regression in the NORM approach is a function of several different aspects of the multilevel structure of the data.

DI approach. Using the DI approach, both the within-group covariance $\sigma_{W,XY}$ and the between-groups covariance $\sigma_{B,XY}$ can be estimated without bias. Thus, the DI approach provides unbiased estimators of the within-group regression coefficient $\beta_{W,YX}$ and the between-groups regression $\beta_{B,YX}$, as these two estimators are based on the within- and between-groups variance of the fully observed variable X (i.e., σ_X^2 and τ_X^2) and the within- and between-groups covariances. In addition, the estimator of the within-group regression coefficient $\beta_{W,XY}$ is unbiased because the within-group

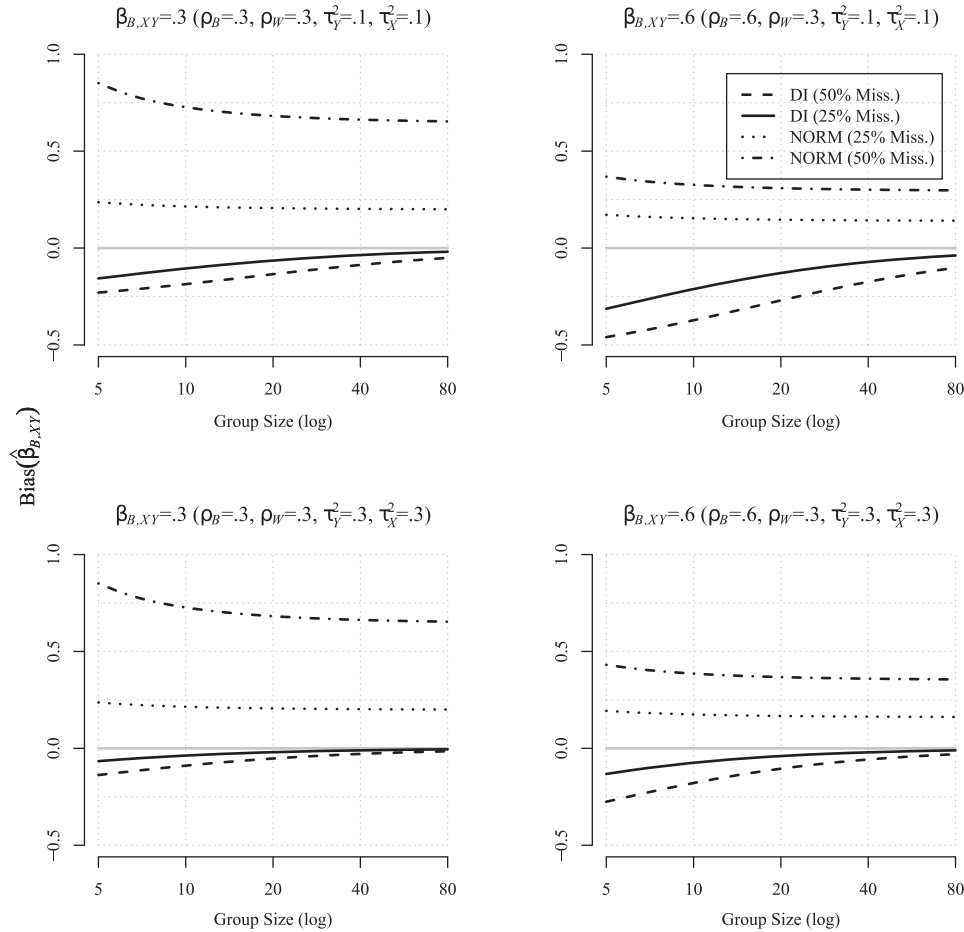


Figure 2. Asymptotic bias of the estimator of the between-groups coefficient $\beta_{B,XY}$ as a function of the size of the between-groups variance of Y , the between-groups correlation, the logarithm of group size, missing data rate in groups, and the multiple imputation strategy. Both variables X and Y are standardized with unit variance. It is assumed that the number of groups approaches infinity. NORM = normal model imputation; DI = dummy-indicator approach.

variance σ_Y^2 is also estimated without bias using the DI approach. However, the estimator of the between-groups regression coefficient $\beta_{B,XY}$ in the regression of X on Y is biased:

$$Bias(\hat{\beta}_{B,XY}) = -\beta_{B,XY} \cdot \frac{Bias(\hat{\tau}_Y^2)}{\tau_Y^2 + Bias(\hat{\tau}_Y^2)}. \quad (22)$$

It is evident that the bias becomes stronger when the bias of the estimator of the between-groups variance of Y increases. The relationship indicates that with missing values on the predictor variable, a positive between-groups regression coefficient ($\beta_{B,XY} > 0$) will be underestimated when the DI approach is used, particularly with small groups and a low ICC of the predictor Y (see Figure 2).

Overall, an examination of the asymptotic bias (i.e., the number of groups approaches infinity) for the two MI strategies showed that under the assumption of MCAR, the estimators of the between-groups variance and also of the intraclass correlation can be dramatically biased in certain data configurations. When the true intraclass correlation is not small, ignoring the multilevel

structure (NORM) can be problematic, particularly when the fraction of missing data is large. Additionally, in the NORM approach, all four estimators of regression coefficients are biased, and the magnitude and direction of the bias are a function of several different aspects of the multilevel structure. For the DI approach, the bias for the estimator of the between-groups variance approaches zero when the group size increases and/or the true intraclass correlation is large. This is because the DI approach relies on the information in the observed group means to approximate the true group effects of the multilevel structure. Furthermore, in the DI approach, only the estimator of the between-groups coefficient of the regression of X on Y is biased. In the next section, we present an imputation strategy that directly incorporates the true group effects and is based on a multivariate mixed effects model.

Multivariate Mixed Effects Imputation Model

A multivariate linear mixed effects model for imputing incomplete multilevel data has been developed by Schafer (2001; Schafer & Yucel, 2002). This model is used in the R package pan

(Schafer & Zhao, 2013), which several authors have identified as the method of choice for dealing with multilevel missing data (Andridge, 2011; Graham, 2012). In its general form, the model is written as

$$\mathbf{Y}_{ij} = \mathbf{X}_{ij}\boldsymbol{\beta} + \mathbf{Z}_{ij}\mathbf{b}_j + \boldsymbol{\varepsilon}_{ij}, \quad (23)$$

where \mathbf{Y}_{ij} is a $(1 \times r)$ vector of outcome variables for person i in group j , and \mathbf{X}_{ij} and \mathbf{Z}_{ij} are $(1 \times p)$ and $(1 \times q)$ vectors of covariate values (each containing a one for an intercept), $\boldsymbol{\beta}$ is a $(p \times r)$ matrix of regression coefficients, \mathbf{b}_j is a $(q \times r)$ matrix of random effects, and $\boldsymbol{\varepsilon}_{ij}$ is a $(1 \times r)$ vector of residuals. In most cases, the covariates in \mathbf{Z}_{ij} , which are allowed to have randomly varying effects across groups, are a subset of the values in \mathbf{X}_{ij} ($p \geq q$). The random effects matrix \mathbf{b}_j is assumed to follow a normal distribution with mean zero and covariance matrix $\boldsymbol{\Psi}$, and to be independently and identically distributed across groups. The residual vector $\boldsymbol{\varepsilon}_{ij}$ is independently and normally distributed across persons with mean zero and covariance matrix $\boldsymbol{\Sigma}$.

A limitation of the imputation model in Equation 23 is that only completely observed covariates can be included in \mathbf{X}_{ij} and \mathbf{Z}_{ij} . However, as the present study is interested only in random-intercept models, we simplify the right-hand side of the multivariate mixed effects model and write the model as an “empty model” without covariates:

$$\mathbf{Y}_{ij} = \boldsymbol{\mu} + \mathbf{Y}_{B,j} + \mathbf{Y}_{W,ij}, \quad (24)$$

where $\boldsymbol{\mu}$ is now a $(1 \times r)$ vector of means, $\mathbf{Y}_{B,j}$ is a $(1 \times r)$ vector of random effects between groups, and $\mathbf{Y}_{W,ij}$ is a $(1 \times r)$ vector of residuals within groups. The model is referred to by Hox (2010) as a multivariate multilevel model. It can also be interpreted as a variance decomposition model that decomposes the multivariate outcome \mathbf{Y}_{ij} into between-groups and within-group parts $\mathbf{Y}_{B,j}$ and $\mathbf{Y}_{W,ij}$ (Cronbach, 1976). We refer to the MI strategy that performs MIs based on this model as the PAN approach. The PAN approach is similar to a two-level imputation model that was proposed by Asparouhov and Muthén (2010) and is implemented in the software *Mplus* (Muthén & Muthén, 1998-2010; see H1 imputation model). The approach proposed by Asparouhov and Muthén (2010) has the further flexibility to generate imputations for Level 2 variables and categorical variables with missing values (see Enders et al., 2016).⁴

To illustrate the PAN approach, we use our example with two variables X and Y , and conclude that

$$(X_{ij(mis)}, Y_{ij(mis)}) = (\mu_X, \mu_Y) + (X_{B,j}, Y_{B,j}) + (X_{W,ij}, Y_{W,ij}), \quad (25)$$

where the random effects between and within groups are normally distributed as follows:

$$(X_{B,j}, Y_{B,j}) \sim N(0, \boldsymbol{\Sigma}_B), \quad (X_{W,ij}, Y_{W,ij}) \sim N(0, \boldsymbol{\Sigma}_W), \quad (26)$$

where $\boldsymbol{\Sigma}_B$ and $\boldsymbol{\Sigma}_W$ denote the between-groups and within-group covariance matrices. The algorithm for the PAN approach includes the random-effects step (RE-step) in addition to the I-step and P-step (Schafer, 2001; Schafer & Yucel, 2002). We again use a person i in group j with a missing value on Y and observations on X to illustrate the imputation of missing values in the I-step. In the $(t+1)$ iteration of the RE-step, the random effects $(X_{B,j}^{(t+1)}, Y_{B,j}^{(t+1)})$ are drawn from an appropriate multivariate normal distribution (e.g.,

Raudenbush & Bryk, 2002) that depends on values for the missing data and the parameters $\boldsymbol{\theta}^{(t)} = (\boldsymbol{\mu}^{(t)}, \boldsymbol{\Sigma}_B^{(t)}, \boldsymbol{\Sigma}_W^{(t)})$ from the previous I- and P-steps.

Based on the random effects, the within-group deviations are then calculated as follows:

$$X_{W,ij}^{(t+1)} = X_{ij(mis)} - \mu_X^{(t)} - X_{B,j}^{(t+1)}, \quad (27)$$

where $X_{W,ij}$ is a within-group deviation for a person i in group j who has a missing value on Y . In the I-step, an imputation for Y is then generated using a multilevel regression model

$$Y_{ij(mis)}^{(t+1)} = \mu_Y^{(t+1)} + \beta_1^{(t+1)}X_{W,ij}^{(t+1)} + Y_{B,j}^{(t+1)} + \boldsymbol{\varepsilon}_{ij}, \quad (28)$$

where the residual $\boldsymbol{\varepsilon}_{ij}$ is normally distributed across persons with constant variance $\sigma^{2(t+1)}$. The parameters $(\mu_Y^{(t+1)}, \beta_1^{(t+1)}, \sigma^{2(t+1)})$ are based on posterior draws of the parameters $\boldsymbol{\theta}^{(t+1)} = (\boldsymbol{\mu}^{(t+1)}, \boldsymbol{\Sigma}_B^{(t+1)}, \boldsymbol{\Sigma}_W^{(t+1)})$ from appropriate posterior distributions in the P-step (Schafer & Yucel, 2002). To demonstrate that the imputation model also takes into account the relationship among the variables between groups, we could rewrite Equation 28 (omitting index t) and replace $Y_{B,j}$ by $\beta_2 X_{B,j} + u_j$, where u_j is the part of $Y_{B,j}$ that is not explained by $X_{B,j}$, and β_2 describes the corresponding between-groups relation: $Y_{ij(mis)} = \mu_Y + \beta_1 X_{W,ij} + \beta_2 X_{B,j} + u_j + \boldsymbol{\varepsilon}_{ij}$. We now see that in the DI approach in Equation 10, the true group effects $\beta_2 X_{B,j} + u_j$ are approximated by the separate intercepts (or fixed effects), whereas the NORM approach in Equation 9 completely ignores the group effects. Furthermore, the expectations of β_1 and β_2 correspond to the within- and between-groups regression coefficients $\beta_{W,YX}$ and $\beta_{B,YX}$ in Equation 6, indicating that the empty model of the PAN approach (see Equation 24) adequately preserves the relationships within and between groups.

Simulation Study

We conducted a computer simulation to evaluate the statistical behavior of the three MI strategies (NORM, DI, and PAN) for dealing with incomplete multilevel data. The simulation study was designed to generate data that resembled data structures found in typical applications of multilevel analysis in psychological research (e.g., Level 1 individuals are nested within Level 2 units such as working groups or school classes). We also included listwise deletion, as this is still frequently used in research (Jelicic, Phelps, & Lerner, 2009; Peugh & Enders, 2004). The data-generating population model was a simple bivariate model with two normally distributed variables X and Y at Level 1 and Level 2. Missing values were limited to the variable Y , and X was assumed to be fully observed. We focused on two analysis models and on how the performance of their parameter estimates was influenced by the different MI strategies. First, we were interested in the intraclass correlation of Y that is routinely reported for multilevel data. Second, we were interested in a multilevel random-intercept model with latent group means (see Equation 6), which is often used to assess the relationship between two variables X and Y within and between groups. We investigated how the three MI strategies influenced estimation of the within-group and between-groups

⁴ For selected conditions of our simulation study, we compared the *Mplus* H1 imputation with the PAN approach. Both approaches yielded almost identical results.

regression coefficient in two scenarios: (a) when missing values occur only on the predictor variable (i.e., Y is the predictor variable and X is the dependent variable), and (b) when missing values occur only on the outcome variable (i.e., X is the predictor variable and Y is the dependent variable). Given the results from our mathematical derivations, our main focus is on the scenario when missing values occur on the predictor variable, particularly for estimating the between-groups regression coefficient.

Simulation Model and Conditions

Data were simulated based on a population with two standardized, bivariate, normally distributed variables X and Y (see Equations 3 and 4). The following population parameters were manipulated: the number of Level 2 groups, the number of observations per Level 2 group, the intraclass correlations of X and Y , and the correlation between X and Y .

Number of groups. The number of Level 2 groups was set to $K = 50$ and 150 . Although there are studies involving fewer than 50 groups, a sample size of 50 is commonly found in educational and organizational psychology (e.g., Maas & Hox, 2005; Mathieu, Aguinis, Culpepper, & Chen, 2012). However, because of the growing number of large-scale studies in the field, a condition with larger samples was also included.

Group size. The number of observations per Level 2 group was set to $n = 5, 15,$ and 30 . A group size of 5 is normal in small-group research, where multilevel modeling is also frequently applied (see Kenny, Mannetti, Pierro, Livi, & Kashy, 2002). Group sizes of 15 and 30 are typical of educational psychology research on class or school characteristics.

Intraclass correlation of X and Y . The intraclass correlations of X and Y (i.e., the amount of between-groups variance) were both varied and set to $\rho_{IX} = .10$ and $.30$, and $\rho_{IY} = .10$ and $.30$. Intraclass correlations rarely take on values greater than $.30$ in educational and organizational research (Bliese, 2000; Hedges & Hedberg, 2007). As the total variance of X and Y was assumed to be 1, the values of the intraclass correlations are equal to the between-groups variances of X and Y .

Correlation between X and Y . The correlation between X and Y at the group level ρ_B was set to $.35$ and $.60$, whereas the correlation of the individual deviations ρ_W was held constant at $.35$. The idea was to select conditions with medium-sized and large correlations in the sense of Cohen's (1988) classification. In addition, studies using multilevel data often show stronger correlations at the group-level compared with the individual level (see Ostroff, 1993), and we also expected that varying the between-groups correlation would be of importance. Given the bivariate two-level structure in Equation 4, the within- and between-groups correlations together with the intraclass correlations of X and Y completely determine the value of the regression coefficients within-group (i.e., $\beta_{W,YX}$ and $\beta_{W,XY}$) and between-groups (i.e., $\beta_{B,YX}$ and $\beta_{B,XY}$). For example, in the condition with $\rho_{IX} = .10$, $\rho_{IY} = .30$, and $\rho_B = .60$, the between-groups coefficient $\beta_{B,XY}$ in the regression of X on Y is $(\rho_{IX}/\rho_{IY})^{1/2} \cdot \rho_B = (.10/.30)^{1/2} \cdot .60 = .346$. In addition, the value of the corresponding within-group coefficient $\beta_{W,XY}$ is $(1 - \rho_{IX})^{1/2}/(1 - \rho_{IY})^{1/2} \cdot \rho_W = (.90/.70)^{1/2} \cdot .35 = .397$.

Missing Data Mechanism

For each simulated data set, missing values on Y were generated using two different missing data mechanisms (MCAR and MAR). More specifically, missing values were imposed on Y by defining a normally distributed response tendency r_{ij}^* , where an individual case on Y is missing if $r_{ij}^* > 0$. The probability of missingness on Y was modeled to be dependent on the within-group and between-groups portions of X by specifying the following relationship:

$$r_{ij}^* = \alpha + \lambda X_{B,j} + \lambda X_{W,ij} + \varepsilon_{ij}, \quad (29)$$

where α is a quantile of the standard normal distribution based on a missing data probability (i.e., $\alpha = -0.67$ for the condition 25% missing data on Y , and $\alpha = 0$ for 50%), and λ is used to control the missing data mechanism. The residual variance is set to $Var(\varepsilon_{ij}) = 1 - \lambda^2$. Note that specifying the same value of λ for the within and between effects of X on Y eliminates any contextual effects of X on the missingness of Y . In order for Y to be MCAR, we set $\lambda = 0$, and for the two MAR conditions, we set $\lambda = 0.4$ or $\lambda = 0.8$. The missing data rate was set to 25% and 50%. We decided to include such an extreme missing data condition. Missing data rates of up to 50% are common in the planned missingness designs that are used in educational research (see also our real-data example in the next section). However, our main focus was on the 25% missing data condition.

Missing Data Treatment

The software mice (van Buuren & Groothuis-Oudshoorn, 2011) was used to implement the two ad hoc procedures (NORM and DI). For the NORM approach, the "norm" method was specified in order to impute Y . The grouping variables were not included in the imputation model, thus ignoring the multilevel structure of the data. Through use of the "norm" method, the missing values were imputed assuming a normal distribution given the completely observed X . For the DI approach, a set of $K-1$ dummy variables was created and included in the imputation model. This resulted in the estimation of a separate intercept for each group to represent the multilevel structure.⁵ Schafer and Olsen (1998) suggested that $m = 10$ imputations are sufficient for most practical purposes. Following this recommendation, we generated for all ad hoc approaches $m = 10$ imputations for each data set (but see Bodner, 2008). We used the default value of 5 iterations in the software mice for each imputation. However, only a single iteration would suffice because the case of only one missing variable is a special case of a monotone missing data pattern (see Carpenter & Kenward, 2013, p. 77). The PAN approach is implemented in the R package pan (Schafer & Zhao, 2013). Variables X and Y were both specified as responses in the pan model to allow for variance and covariance at Level 1 and Level 2 (see Equation 24). The least informative inverse-Wishart priors were chosen for the covariance

⁵ Note that even though the mice package was used for implementing the NORM and DI approach, these approaches can still be considered as joint modeling. This is because for a single missing data pattern (or monotone patterns of missing data in general), the conditional imputations generated by mice are equivalent to imputations from the joint model with noninformative priors (Raghunathan, Lepkowski, Hoewyk, & Solenberger, 2001; see also Schafer, 1997).

matrices at Level 1 and Level 2, that is, $\Sigma_W \sim W^{-1}(I_2, 2)$ and $\Sigma_B \sim W^{-1}(I_2, 2)$. For the PAN approach, convergence behavior was assessed by inspecting the autocorrelation functions and trace plots of the different parameters. Applying the two criteria to a subsample of the replications of the simulation design, we concluded that the MCMC chains had reached convergence after the first 200 iterations. We let the software pan perform 200 burn-in iterations before drawing one imputed dataset for each 50 iterations, leading to $m = 10$ imputed data sets.

Analysis Models and Outcome Variables

For each of the $2 \times 3 \times 2 \times 2 \times 2 \times 3 \times 2 = 288$ conditions (five factors for the population model, two factors for the missing data mechanism), 1,000 simulated data sets were generated for each condition, which allowed for a precise estimation of bias, root mean square error (RMSE), and coverage rate. After imputing the missing values on Y using the three different MI strategies, all of the statistical analyses were conducted in *Mplus 6* (Muthén & Muthén, 1998–2010). We specified three different analysis models. First, the intraclass correlation of Y (see Equation 5) was estimated by specifying an empty two-level model with Y as the outcome variable. Second, a multilevel random-intercept model was specified in which Y was the predictor variable and X was the dependent variable (see Equation 8). This model produced estimates of the regression coefficients $\beta_{W,XY}$ and $\beta_{B,XY}$ under the scenario that missing values occur in the predictor variable. Third, we specified a multilevel random-intercept model with X as the predictor variable and Y as the dependent variable to estimate the within- and between-groups regression coefficients $\beta_{W,YX}$ and $\beta_{B,YX}$ (see Equation 6). Note that in both multilevel random-intercept models the group means of the predictor variable were treated as a latent variable (see Lüdtke et al., 2008).

We used bias, RMSE, and confidence interval coverage to evaluate the missing data strategies. The bias was estimated by calculating the difference between the mean parameter estimate from each design cell and the true population parameter. The overall accuracy of the parameter estimates was assessed using the RMSE, which was computed by taking the square root of the mean square difference of the estimate and the true parameter. When a parameter estimate is biased, the RMSE combines bias and variability (i.e., sampling variance) into an overall measure of accuracy. Furthermore, we analyzed the accuracy of the standard errors for the regression coefficients by determining the observed coverage of the 95% confidence interval (CI). Coverage was given a value of 1 if the true value was included in the confidence interval and a value of 0 if the true value was outside the confidence interval. To provide an additional benchmark for the results from the different MI strategies, we also show the results from the analysis of the complete data sets, that is, the results obtained from the data sets before the missing values have been induced.

Results

Intraclass Correlation of Y

Table 1 shows the estimated bias in the parameter estimates for the intraclass correlation of Y ($\rho_{I,Y}$) for selected conditions of the

simulation (number of groups is $K = 150$, missing rate = 25%, and MCAR and MAR; see the online supplemental materials for detailed information about all conditions). Using the PAN approach as an imputation model produced approximately unbiased estimates of the intraclass correlation. Only when the number of groups was small ($n = 5$) was there a slight tendency to overestimate the size of the intraclass correlation. In the worst condition ($n = 5$, $\rho_{I,X} = .30$, $\rho_{I,Y} = .10$, $\lambda = 0$), the estimated bias was 0.015, which is a relative percentage bias of 15%, given that the true intraclass correlation was .10. However, with larger group sizes ($n \geq 15$), the relative percentage bias was very small and ranged from -1.3% to 3.5% .

The intraclass correlation estimates of the NORM approach were negatively biased, and the estimates of the DI approach were positively biased. The magnitude of the absolute estimated bias for the NORM approach strongly depended on the size of the true intraclass correlation. It was much more pronounced for a large intraclass correlation ($\rho_{I,Y} = .30$), with values ranging from -0.127 to -0.096 , than for a small intraclass correlation ($\rho_{I,Y} = .10$), with values ranging from -0.041 to -0.026 . In contrast, the DI approach yielded *less* bias when the true intraclass correlation of Y was large, particularly for smaller groups. For example, when the group size was held at $n = 5$, the estimated bias ranged from 0.112 to 0.129 for a small intraclass correlation, and from 0.064 to 0.081 for a large intraclass correlation. However, with large groups ($n = 30$), the DI approach produced only slightly biased estimates of the intraclass correlation, with values ranging from 0.006 to 0.129. Listwise deletion provided approximately unbiased estimates when data were MCAR, but was negatively biased when the missingness in Y depended on X .

The main findings for the estimated bias of the intraclass correlation are also depicted in Figure 3. Differences between the approaches are particularly pronounced when the number of groups was small ($n = 5$). It is also apparent that a large group size and/or a high intraclass correlation of Y are needed to obtain acceptable estimates for the intraclass correlation with the DI approach. This reflects the fact that the observed group means yield more reliable estimates of the true group means when the group size is large and the intraclass correlation of Y is high (e.g., Bliese, 2000).

Next, we assessed the overall accuracy of the parameter estimates by estimating the RMSE. As expected, using the PAN approach resulted in the lowest estimated RMSE across most of the conditions. Overall, the differences and trends in the estimated RMSE values of the MI approaches were very similar to the results for bias.

Multilevel Regression of X on Y

Between-groups regression. The estimated bias for the estimator of the between-groups regression coefficient of X on Y ($\beta_{B,XY}$) is presented in Table 2. Again, only selected conditions are presented ($K = 150$, missing rate 25%, MCAR, and strong MAR). The PAN approach provided approximately unbiased estimates of the between-groups regression coefficient, except in conditions with a small number of groups. In the worst condition depicted in Table 2 ($\rho_{I,X} = .30$, $\rho_{I,Y} = .10$, $\rho_B = .60$), the estimated bias was -0.109 (corresponding to a relative percentage bias

Table 1
Bias of the Estimator of the Intraclass Correlation of Y for a Large Number of Groups (K = 150) and 25% Missing Data

Conditions	MCAR ($\lambda = 0$)					MAR ($\lambda = .8$)				
	NORM	DI	PAN	CD	LD	NORM	DI	PAN	CD	LD
$\rho_{I,X} = .10, \rho_{I,Y} = .10$										
$\rho_B = .35$										
$n = 5$	-.041	.112	.009	-.003	-.003	-.040	.118	.010	-.003	-.005
$n = 15$	-.039	.032	.002	-.001	-.001	-.039	.033	.001	-.001	-.004
$n = 30$	-.038	.017	.001	-.001	-.001	-.039	.016	.000	-.001	-.003
$\rho_B = .60$										
$n = 5$	-.034	.115	.014	.000	.000	-.036	.118	.011	-.001	-.010
$n = 15$	-.035	.032	.003	-.001	-.001	-.035	.034	.002	-.001	-.009
$n = 30$	-.035	.015	.001	-.001	-.001	-.035	.016	.001	-.001	-.009
$\rho_{I,X} = .10, \rho_{I,Y} = .30$										
$\rho_B = .35$										
$n = 5$	-.124	.070	-.001	-.002	-.002	-.124	.070	-.004	-.004	-.005
$n = 15$	-.124	.018	-.002	-.002	-.003	-.123	.019	-.003	-.003	-.003
$n = 30$	-.124	.008	-.002	-.002	-.002	-.123	.008	-.002	-.003	-.002
$\rho_B = .60$										
$n = 5$	-.118	.069	-.001	-.001	-.001	-.117	.070	-.004	-.003	-.011
$n = 15$	-.116	.019	-.001	-.003	-.003	-.116	.019	-.002	-.002	-.009
$n = 30$	-.116	.009	-.001	-.002	-.002	-.116	.007	-.003	-.004	-.009
$\rho_{I,X} = .30, \rho_{I,Y} = .10$										
$\rho_B = .35$										
$n = 5$	-.035	.114	.012	-.001	-.002	-.033	.129	.011	-.001	-.007
$n = 15$	-.034	.033	.002	-.001	-.001	-.033	.040	.002	-.001	-.007
$n = 30$	-.034	.016	.001	-.001	.000	-.032	.019	.001	-.001	-.007
$\rho_B = .60$										
$n = 5$	-.027	.115	.015	-.001	-.002	-.026	.128	.014	-.001	-.017
$n = 15$	-.027	.032	.003	-.002	-.002	-.026	.040	.003	-.001	-.016
$n = 30$	-.027	.015	.001	-.001	-.001	-.026	.017	.000	-.001	-.015
$\rho_{I,X} = .30, \rho_{I,Y} = .30$										
$\rho_B = .35$										
$n = 5$	-.120	.064	-.007	-.004	-.006	-.111	.081	.000	.000	-.006
$n = 15$	-.117	.017	-.003	-.004	-.004	-.111	.023	-.001	-.003	-.009
$n = 30$	-.118	.006	-.004	-.003	-.003	-.112	.009	-.003	-.003	-.008
$\rho_B = .60$										
$n = 5$	-.099	.071	.003	-.002	-.002	-.098	.076	-.003	-.004	-.024
$n = 15$	-.100	.019	-.001	-.002	-.002	-.096	.023	-.001	-.002	-.019
$n = 30$	-.101	.008	-.002	-.003	-.003	-.097	.010	-.002	-.003	-.020

Note. Biases larger than 10% are written in bold. n = group size; $\rho_{I,X}$ = intraclass correlation of X; $\rho_{I,Y}$ = intraclass correlation of Y; ρ_B = correlation at Level 2; λ = effect of X on missingness; NORM = normal model imputation; DI = dummy-indicator approach; PAN = two-level imputation; CD = complete data; LD = listwise deletion.

of -10.5%).⁶ However, with larger group sizes ($n \geq 15$), the negative bias disappeared, with values ranging from -0.040 to 0.005. The NORM approach, which ignores the multilevel structure, showed a positive bias and tended to overestimate the size of the true between-groups regression coefficient (range = 0.093 to 0.587). In contrast, the DI approach was negatively biased and underestimated the true value of the between-groups regression coefficient (range = -0.649 to -0.009). The magnitude of the estimated bias was particularly pronounced for a small group size ($n = 5$) and a low intraclass correlation ($\rho_{I,Y} = .10$). As the true value of the between-groups coefficient depends on the intraclass correlations of both X and Y, the effect of group size and the intraclass correlation can best be seen when comparing the conditions with the same true value of the between-groups regression coefficient in Table 2 (e.g., upper half; true value = .350). It is apparent that the estimated bias is strongly reduced when the group size and/or the intraclass correlation of Y

increase. Listwise deletion provided estimates that were strongly negatively biased under the MAR conditions, but were

⁶ The estimates produced by the PAN approach were slightly negatively biased in conditions with a small group size ($n = 5$) and a relatively low intraclass correlation of the predictor ($\rho_{I,Y} = .10$). We believe that this finding was due to the standard least-informative prior for the covariance matrix of the random effects, which induces bias into small variance components (see also Grund, Lüdtke, & Robitzsch, 2016). The software pan uses an inverse-Wishart prior distribution for the random effects that implies, in the case of two variables, a prior distribution for the variances that is loosely centered on a value of .50. With small group sizes and a relatively small true variance component, this could result in slightly positively biased estimates of variance components, which in turn would yield negatively biased between-groups regression coefficients. This explanation is also consistent with the finding that the estimator of the intraclass correlation of Y in the PAN approach is also slightly positively biased in conditions with a small number of groups and a low intraclass correlation of Y. Note that the complete data analysis produced approximately unbiased estimates in these conditions.

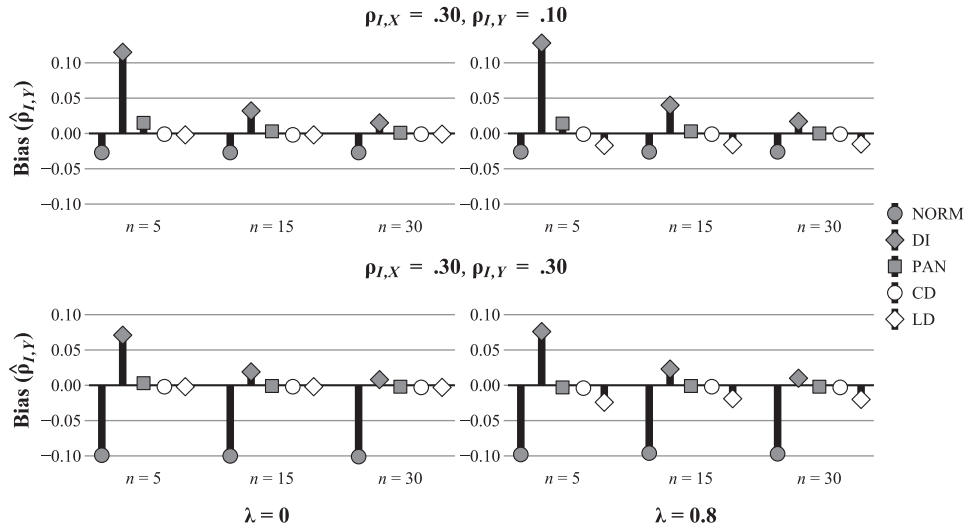


Figure 3. Bias of the estimator of the intraclass correlation of Y ($\rho_{t,Y} = .10$ and $.30$) for varying group size (n) and effect of X on missingness (λ). The intraclass correlation of X was fixed at $\rho_{t,X} = .30$, the correlation at Level 2 at $\rho_B = .60$, the number of groups at $K = 150$, and the missing data probability at 25%. NORM = normal model imputation; DI = dummy-indicator approach; PAN = two-level imputation; CD = complete data; LD = listwise deletion.

only slightly biased when the data were MCAR. The main findings are summarized in Figure 4.

In the next step, we estimated the RMSE for the between-groups regression coefficient. Figure 5 shows the main findings when the data were MAR. As we see, the PAN approach outperformed the other approaches and the NORM approach showed the largest estimated RMSE across the conditions. In conditions with a large correlation at Level 2 ($\rho_B = .30$), the performance of the DI approach improved considerably with larger group sizes and a large intraclass correlation of Y . It is also apparent that listwise deletion, which provides substantially biased estimates under MAR, showed the second largest RMSE of the compared methods.

The accuracy of the standard errors produced by the various MI strategies was evaluated in terms of the coverage rate, which was assessed using the 95% CIs. As shown in Figure 6, the coverage rate for the DI approach was not accurate and mirrored the pattern for the estimated bias. The probability that the CIs cover the true value was higher for large groups ($n \geq 15$) than for small groups ($n = 5$). In contrast, the coverage rate for the NORM approach dropped substantially when the group size increased, although the bias for the NORM approach was only slightly influenced by the group size. It is also evident that the coverage rates produced by listwise deletion were close to the nominal value when data were MCAR, but were poor in MAR conditions. Overall, the PAN approach provided acceptable coverage rates, with values ranging from 89.5 to 96.5. However, it should be added that with a high rate of missing data (50%), there were a few conditions in which the coverage values produced by the PAN approach were not acceptable. This was the case when the number of groups was small, the intraclass correlation of Y was low, the intraclass correlation of X was large, the correlation at Level 2 was large, and the number of groups was $K = 150$. In these conditions the coverage rates were 80.3, 80.3, and 82.6.

In addition to examining coverage rates, we looked at whether the estimated standard errors correctly describe the sampling distribution of the point estimates. Figure 7 shows box plots of the ratio of the estimated standard errors and the empirical standard deviation of the estimates for the between-groups regression coefficient in selected conditions. With low intraclass correlations ($\rho_{t,X} = \rho_{t,Y} = .10$) and small groups ($n = 5$), the individual standard errors were sometimes a poor estimate of the sampling variability of the between-groups regression coefficient—a result that is most pronounced for NORM and least pronounced for the DI approach. In general, however, the median standard errors for the between-groups regression coefficient were very close to the observed standard deviation of the point estimates.

Within-group regression. For the within-group coefficient of the regression of X on Y , we found no substantial estimated bias for the DI or PAN approach ($\beta_{w,XY}$). For example, the estimated bias observed for the PAN approach ranged from -0.008 to 0.008 , which represents 2.2% downward and upward bias. Both the NORM approach and listwise deletion had a tendency to underestimate the true within-group coefficient. The largest absolute bias was -0.057 (or -14.4%) for the NORM approach (range = -0.057 to -0.004), and -0.102 (or -25.8%) for listwise deletion (range = -0.102 to 0.004). The NORM approach was particularly biased when the intraclass correlation of Y was large with relative bias values ranging from -14.8% to -7.1% , whereas listwise deletion was especially biased when data were MAR (range = -26.2% to -19.6%). The estimated RMSE values were lowest overall for the DI and PAN approaches. For the NORM approach and listwise deletion, the RMSE was usually larger in conditions where estimates were biased. Coverage of the 95% CIs was satisfactory across all conditions for the DI approach (range = 91.9 to 96.6) and the PAN approach (range = 92.4 to 96.3). As was expected from our findings regarding bias, the NORM ap-

Table 2
Bias of the Estimator of the Between-Group Regression Coefficient (X on Y) for a Large Number of Groups (K = 150) and 25% Missing Data

Conditions	MCAR ($\lambda = 0$)					MAR ($\lambda = .8$)				
	NORM	DI	PAN	CD	LD	NORM	DI	PAN	CD	LD
Moderate correlation at Level 2 ($\rho_B = .35$)										
$\rho_{I,X} = .10, \rho_{I,Y} = .10$ (true value = .350)										
$n = 5$.279	-.209	-.043	-.007	-.009	.241	-.223	-.067	-.008	-.174
$n = 15$.231	-.094	-.010	.003	.004	.211	-.109	-.023	-.006	-.171
$n = 30$.216	-.056	-.009	.000	.001	.214	-.057	-.007	-.006	-.166
$\rho_{I,X} = .10, \rho_{I,Y} = .30$ (true value = .202)										
$n = 5$.094	-.060	-.013	-.003	-.005	.093	-.063	-.012	-.003	-.080
$n = 15$.103	-.019	-.001	-.001	-.002	.099	-.019	-.001	.001	-.076
$n = 30$.102	-.009	-.001	.000	.000	.099	-.010	-.001	.000	-.075
$\rho_{I,X} = .30, \rho_{I,Y} = .10$ (true value = .606)										
$n = 5$.587	-.361	-.051	.021	.052	.528	-.379	-.059	.014	-.361
$n = 15$.498	-.157	-.001	.018	.017	.445	-.199	-.021	.000	-.356
$n = 30$.459	-.096	-.010	.001	.000	.434	-.112	-.013	-.002	-.351
$\rho_{I,X} = .30, \rho_{I,Y} = .30$ (true value = .350)										
$n = 5$.231	-.097	.001	.001	.002	.198	-.110	-.008	-.004	-.156
$n = 15$.221	-.030	.002	.001	.001	.204	-.038	.000	.002	-.147
$n = 30$.215	-.019	-.003	-.002	-.002	.209	-.014	.005	.004	-.139
Large correlation at Level 2 ($\rho_B = .60$)										
$\rho_{I,X} = .10, \rho_{I,Y} = .10$ (true value = .600)										
$n = 5$.261	-.354	-.090	.012	.020	.275	-.366	-.091	.026	-.188
$n = 15$.233	-.171	-.037	-.001	-.002	.241	-.171	-.028	.010	-.193
$n = 30$.232	-.094	-.016	.003	.004	.229	-.096	-.014	.005	-.198
$\rho_{I,X} = .10, \rho_{I,Y} = .30$ (true value = .346)										
$n = 5$.126	-.100	-.018	-.001	-.001	.121	-.104	-.016	.000	-.110
$n = 15$.131	-.030	-.001	-.001	-.001	.123	-.034	-.004	-.001	-.109
$n = 30$.127	-.016	-.001	.001	.001	.120	-.019	-.004	-.002	-.108
$\rho_{I,X} = .30, \rho_{I,Y} = .10$ (true value = 1.039)										
$n = 5$.520	-.617	-.109	.060	.095	.485	-.649	-.107	.067	-.311
$n = 15$.427	-.287	-.040	.004	.009	.405	-.331	-.040	.013	-.364
$n = 30$.412	-.160	-.016	.003	.004	.379	-.191	-.025	.005	-.362
$\rho_{I,X} = .30, \rho_{I,Y} = .30$ (true value = .600)										
$n = 5$.245	-.161	-.001	.005	.007	.223	-.189	-.009	.005	-.190
$n = 15$.231	-.056	-.003	.000	.000	.212	-.068	-.003	-.002	-.185
$n = 30$.230	-.028	-.001	.001	.001	.212	-.034	-.002	.000	-.178

Note. Biases larger than 10% are written in bold. n = group size; $\rho_{I,X}$ = intraclass correlation of X; $\rho_{I,Y}$ = intraclass correlation of Y; λ = effect of X on missingness; NORM = normal model imputation; DI = dummy-indicator approach; PAN = two-level imputation; CD = complete data; LD = listwise deletion.

proach had low coverage rates when the intraclass correlation of Y was large (range = 10.8 to 94.7), whereas listwise deletion provided unsatisfactory coverage when data were MAR (range = 0 to 80.8).

Multilevel Regression of Y on X

We also investigated an analysis model in which Y was the dependent and X the predictor variable. In this case, missing values occurred only on the dependent variable. The DI and the PAN approach as well as listwise deletion produced approximately unbiased estimators of the within-group regression coefficient $\beta_{w,YX}$. For example, the estimated bias for the PAN approach ranged from -0.014 to 0.004, or from -3.5% to 1.1% in relative terms. Furthermore, the PAN approach led to approximately unbiased estimates of the between-groups coefficient across all conditions, with absolute bias ranging from -0.052 to 0.130 (or -14.7% to 12.5%). The moderate bias exhibited by PAN, however, was limited to conditions with small groups ($n = 5$) and vanished as soon as the groups grew larger ($n \geq$

15; range -0.024 to 0.035, or -6.8% to 3.4%). In contrast, the NORM approach produced estimates of the between-groups coefficient that were often biased, with absolute bias ranging from -0.163 to 0.037 (or -15.7% to 18.2%). The DI approach had a tendency to overestimate the between-groups coefficient, with bias ranging from -0.030 to 0.267 (or -8.7% to 25.7%). However, this substantial positive bias was only observed for small group sizes and was reduced with large numbers of groups. When the group sizes increased ($n \geq 15$), the bias disappeared with values ranging from -0.011 to 0.051 (or -5.3% to 4.9%). Finally, listwise deletion tended to overestimate the true between-groups coefficient, with absolute bias ranging from -0.074 to 0.466 (or -36.8% to 44.8%).

Summary

The main results of the simulation study can be summarized as follows. First, the simulation confirmed the findings of our mathematical derivations, namely, that the estimator of the intraclass correlation was negatively biased for the NORM approach and

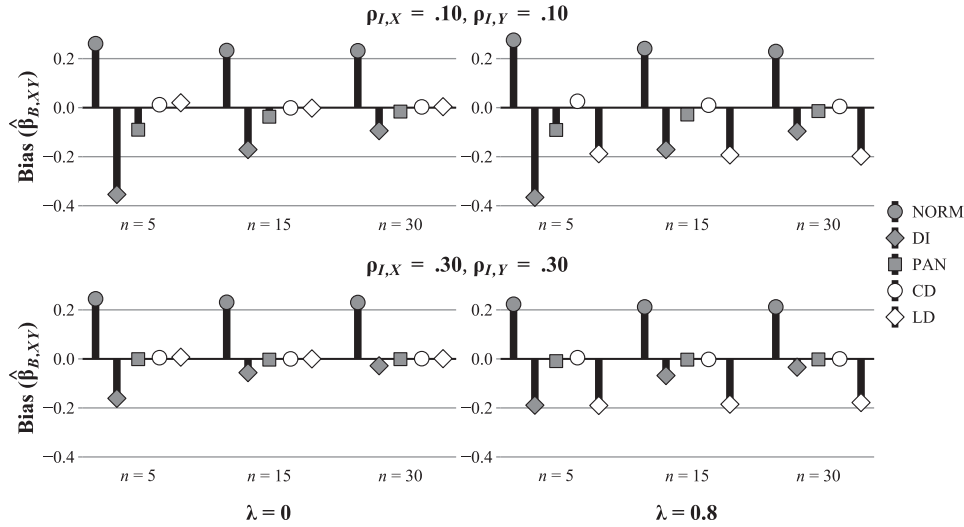


Figure 4. Bias of the estimator of the between-groups regression coefficient (Level 2, X regressed on Y , true value = 0.600) for varying group size (n), intraclass correlation ($\rho_{L,X}$ and $\rho_{L,Y}$), and effect of X on missingness (λ). The correlation at Level 2 was fixed at $\rho_B = .60$, the number of groups at $K = 150$, and the missing data probability at 25%. NORM = normal model imputation; DI = dummy-indicator approach; PAN = two-level imputation; CD = complete data; LD = listwise deletion.

positively biased for the DI approach. Second, for both the intraclass correlation and the between-groups regression coefficient, the performance of the DI approach was particularly problematic in data constellations with small group sizes and low intraclass correlations. In contrast, the performance of the NORM approach did not improve with larger group sizes and was even worse when the true intraclass correlations were large. Third, the PAN approach provided approximately unbiased estimates and accurate

standard errors (i.e., coverage values near the nominal value) across the simulated conditions. It was only in a few conditions with a small group size that the estimates of the intraclass correlation and the between-groups regression coefficient were slightly positively biased. Fourth, listwise deletion produced acceptable parameter estimates only under MCAR conditions. The NORM, DI, and PAN approaches were not strongly influenced by the missing data mechanism (MCAR or MAR). Fifth, increasing the

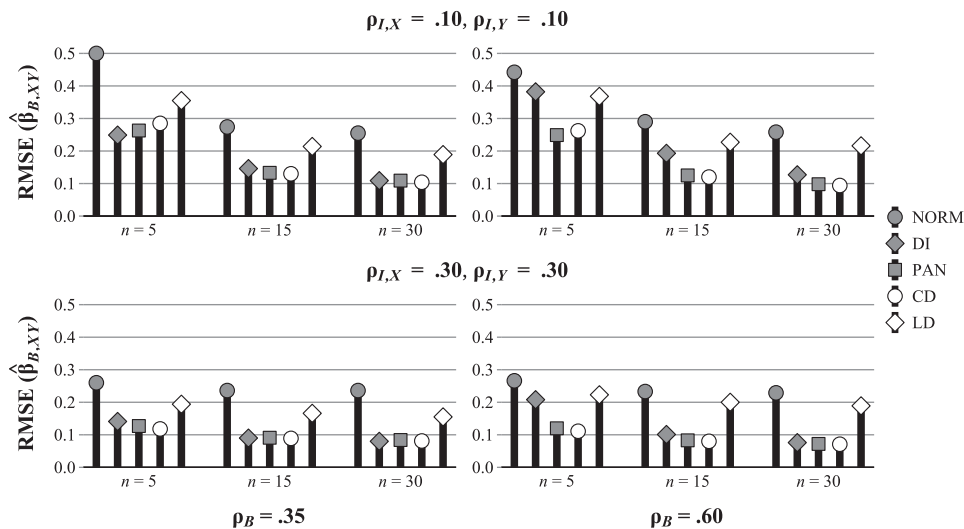


Figure 5. Root mean square error of the estimator of the between-groups regression coefficient (Level 2, X regressed on Y) for moderate ($\rho_B = .35$, true value = .350) and large correlation at Level 2 ($\rho_B = .60$, true value = .600), and varying group size (n) and intraclass correlation ($\rho_{L,X}$ and $\rho_{L,Y}$). The effect of X on missingness was fixed at $\lambda = 0.8$, the number of groups at $K = 150$, and the missing data probability at 25%. NORM = normal model imputation; DI = dummy-indicator approach; PAN = two-level imputation; CD = complete data; LD = listwise deletion.

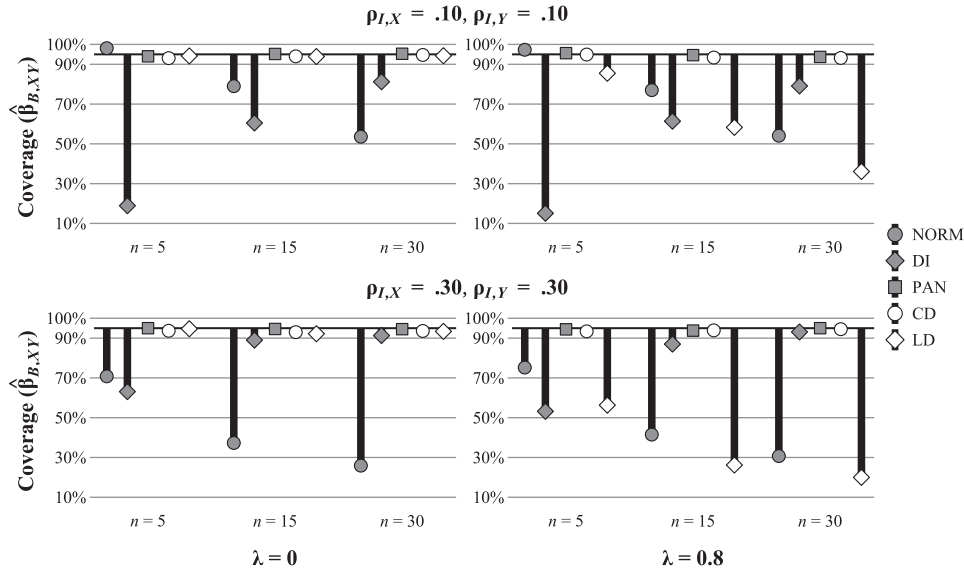


Figure 6. Coverage of the 95% confidence interval of the estimator of the between-groups regression coefficient (Level 2, X regressed on Y , true value = 0.600) for varying group size (n), intraclass correlation (ρ_{LX} and ρ_{LY}), and effect of X on missingness (λ). The correlation at Level 2 was fixed at $\rho_B = .60$, the number of groups at $K = 150$, and the missing data probability at 25%. NORM = normal model imputation; DI = dummy-indicator approach; PAN = two-level imputation; CD = complete data; LD = listwise deletion.

missing data rate from 25% to 50% generally increased bias but did not change the overall picture of the results, with the exception that in a very few conditions the coverage rates for the PAN approach were too low.

Illustrative Data Example

An example from educational psychology is used to illustrate the impact of various MI strategies when estimating the intraclass

correlation with incomplete multilevel data. The data were taken from the German sample of primary school students who participated in 2001 in the Progress in International Reading Literacy Study (Bos et al., 2003; Mullis, Martin, Gonzales, & Kennedy, 2003). In this study, students were asked to rate several specific aspects of their instruction in German and mathematics. However, owing to time constraints, the students in a class were randomly administered different versions of the student questionnaire (six

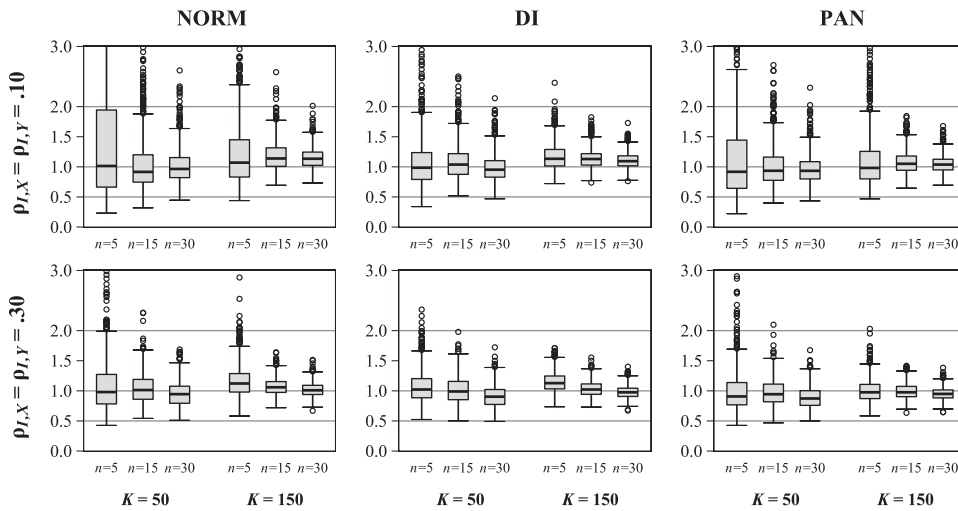


Figure 7. Standard errors divided by the standard deviation of the (point) estimates of the between-groups regression coefficient (X regressed on Y). The boxes indicate the median and the quartiles. Extreme individual values are shown as dots. The correlation at Level 2 was fixed at $\rho_B = .60$, the effect of X on missingness at $\lambda = 0$, and the missing data probability at 25%. NORM = normal model imputation; DI = dummy-indicator approach; PAN = two-level imputation.

different booklets). All students were asked questions that addressed basic background variables, but only three of the six booklets contained questions about the classroom environment in mathematics (planned missing data design; see [Graham, Taylor, Olchowski, & Cumsille, 2006](#)). As a result, approximately 50% of the items are missing by design and can be assumed to be MCAR. The data set contains $N = 8,828$ students nested within 476 classes (average cluster size = 18.5).

In the present example, we focused on two aspects of students' mathematics lessons. First, we were interested in disciplinary problems, which lead to class disruptions and wasted time (see [Kounin, 1970](#)). Students were asked to rate, on five items, how chaotic and unstructured they perceived their mathematics lessons to be (sample item: "The teacher has to wait a long time for students to quiet down"; Cronbach's $\alpha = .80$). Second, we examined students' teacher-related anxiety in mathematics, as assessed by a five-item scale (sample item: "In this teacher's class, I'm afraid that I might do something wrong"; Cronbach's $\alpha = .83$). For both scales the percentage of missing values was above 50% (disciplinary problems in mathematics, 61.5% missing; teacher-related anxiety, 61.7%). Only in the case of 50% of the data could it be assumed that they were MCAR. For the purpose of illustration, we considered three additional measures as auxiliary variables in the imputation model: disciplinary problems in German lessons (21.6% missing), reading achievement scores (0.6% missing), and student ratings of school climate (21.9% missing). Although the scales measuring disciplinary problems in German lessons and school climate were administered in all six booklets, a substantial percentage of the students failed to complete those items.

As in the simulation study, the R package mice was used for the NORM and DI approach. The PAN approach was specified in the pan software. Ten imputations were generated using each procedure. The R code that was used for the three MI approaches is provided in the online supplemental materials. We also used listwise deletion, which excluded 62.2% of the sample from the analyses because of the extreme pattern of missing data.

Table 3 shows the parameter estimates of the within-group variance, the between-groups variance, and the intraclass correlation for the three MI strategies and listwise deletion. The intraclass correlation estimates obtained by the PAN approach were .181 for disciplinary problems and .059 for anxiety. This indicates that 18% of the total variance in the student ratings was located at the class

level for disciplinary problems, but only about 6% for teacher-related anxiety. For disciplinary problems in mathematics, the intraclass correlation estimates of the DI approach were close to the estimates produced by PAN, but for teacher-related anxiety, they were substantially larger. This is in line with the simulation results and our mathematical derivations, which showed that the positive bias of the DI approach is particularly pronounced with a small intraclass correlation. Also consistent with the previous results, the NORM approach, which ignores the multilevel structure, led to smaller estimates of the intraclass correlation for both scales. Finally, as a large amount of the missing data in the two scales was MCAR by design, it is not surprising that estimates produced by listwise deletion deviated only slightly from the estimates obtained by PAN.

Discussion

MI approaches for dealing with missing data problems have received growing attention in psychological research in the last two decades. In this article, we presented mathematical derivations, a computer simulation, and a real-data example to demonstrate the importance of correctly specifying the dependence in the data when using MI for incomplete multilevel data. We showed that of three different MI strategies, only the approach that is based on a multilevel imputation model produced valid parameter estimates of intraclass correlations and regression coefficients in random-intercept models under most of the simulated conditions.

What are the consequences of our findings for dealing with incomplete multilevel data? It is difficult to give general recommendations for research practice, as the impact of the various MI strategies depends on the model of interest and the specific data constellation (e.g., percentage of missing data, intraclass correlations). However, we suggest the following. First, when researchers are not only viewing the multilevel structure as a nuisance factor that needs to be controlled for, but are interested in decomposing the variance of variables at different levels of analysis, there is no alternative to a multilevel imputation model. Even when the focus is on estimating regression coefficients and not on estimating variance components, however, we recommend the PAN approach as a way of obtaining appropriate standard errors.

Second, if the missing data rate is low and the intraclass correlations of the variables are small, the NORM approach that ignores the multilevel structure of the data may produce results that are close to those obtained using a multilevel imputation model. In addition, the NORM approach might be improved by including strong auxiliary variables (e.g., with low rates of missing data, substantial correlations with missingness) that are also associated with the between-groups portion of the missing variables. Furthermore, the DI approach could be a reasonable alternative in the case of large groups and a substantial intraclass correlation, particularly when the focus is on the regression coefficients (see [Drechsler, 2015](#)). However, in most settings, the true values of intraclass correlations are unknown, and with small intraclass correlations, the DI approach might substantially overestimate the variation between groups.

Third, it is important that researchers not only report the amount of missing data but also provide more details about the technique used to deal with that issue. As the present study has shown, results might differ dramatically depending on the MI strategy used. More

Table 3
Variance Components at Level 1 and Level 2, and Intraclass Correlation in the Example Data Set for Different Missing Data Strategies

Method	Disciplinary problems			Teacher-related anxiety		
	$\hat{\sigma}^2$	$\hat{\tau}^2$	$\hat{\rho}_I$	$\hat{\sigma}^2$	$\hat{\tau}^2$	$\hat{\rho}_I$
NORM	.553	.060	.098	.583	.013	.022
DI	.518	.118	.185	.570	.108	.156
PAN	.510	.113	.181	.567	.035	.059
LD	.521	.105	.167	.569	.027	.045

Note. $\hat{\sigma}^2$ = variance at Level 1; $\hat{\tau}^2$ = variance at Level 2; $\hat{\rho}_I$ = intraclass correlation; NORM = normal model imputation; DI = dummy-indicator approach; PAN = two-level imputation; LD = listwise deletion.

specifically, it is very important for researchers to report the variables that were used in the imputation model. This allows other researchers at least to infer how the results might differ if other MI strategies had been chosen (see [Díaz-Ordaz, Kenward, Cohen, Coleman, & Eldridge, 2014](#)).

Although MI is gaining popularity among applied researchers, multilevel imputation models are only rarely used in research practice. A main reason for this is that it can be challenging to apply software that is capable of performing MI using a multilevel imputation model, and documentation is rather technical. It is therefore important for methodologists to provide tutorials that familiarize applied researchers with these important methods. Other multilevel imputation routines are available in addition to the R-package *pan* ([Schafer & Zhao, 2013](#)) that was used in the present study. The REALCOM-IMPUTE software ([Carpenter, Goldstein, & Kenward, 2011](#)) is a standalone software that can handle missing data for both Level 1 and Level 2 variables, as well as categorical variables ([Goldstein, Carpenter, & Browne, 2014](#); see also the R package *jomo*, [Quartagno & Carpenter, 2016](#)). The software *Mplus* ([Muthén & Muthén, 1998-2010](#)) also provides an imputation model (H1 imputation; see [Asparouhov & Muthén, 2010](#)) that can deal with missing values on categorical as well as continuous variables at both Level 1 and Level 2. [Mistler \(2013\)](#) offers a SAS macro (MMI_IMPUTE) that performs multilevel imputation in SAS. It is also possible to use a multilevel model for imputing incomplete, continuous Level 1 variables in a chained equations approach with the function `mice.impute.2l.norm` in the software *mice* ([van Buuren, 2012](#); see also [Enders et al., 2016](#)).

As is true of any simulation study, the results of our study cannot be generalized beyond its specific conditions, for six reasons. First, we did not compare the performance of the various MI strategies under more extreme conditions. For example, a group size of two is common in research with dyads, which are studied in many different psychological disciplines ([Kenny, Kashy, & Cook, 2006](#)). Moreover, psychological variables often show intraclass correlations that are substantially lower (e.g., .05 or smaller) than under the conditions included in the simulation. The DI approach might be expected to be even more problematic in these conditions.

Second, the model of interest in the present study was a multilevel random-intercept model in which the between part of the Level 1 predictor was treated as a latent variable ([Lüdtke et al., 2008](#); [Preacher et al., 2010](#)). This model was used as an analysis model, and it was also assumed that it was the data-generating model in the population. Alternatively, a traditional multilevel model could be used to estimate the group-level effects of Level 1 predictors (e.g., [Raudenbush & Bryk, 2002](#)). The important difference is that in the traditional model, the observed group mean of the predictor is treated as a manifest variable rather than a latent variable. In the online supplement, we provide an analytical argument that the PAN approach is also an appropriate strategy for dealing with incomplete variables in multilevel models with manifest group means. This argument relies on the fact that a bivariate multilevel model can also be represented as a multivariate single-level model ([Mehta & Neale, 2005](#)). It can then be shown that the covariance structure implied by the model with manifest group means will be preserved by the PAN approach (see [Carpenter & Kenward, 2013](#), p. 221).⁷ The analytical argument was also confirmed by an additional simulation in which the PAN approach

produced an approximately unbiased estimator of the group-level effect of the model with manifest group means with coverage rates near the nominal value.

Third, we only considered missing values that occur at Level 1. The treatment of missing data at Level 2 has received less attention in the literature (see [Gibson & Olejnik, 2003](#); [van Buuren, 2011](#)), but can be very important when the model of interest includes Level 1 and Level 2 variables simultaneously. For example, in a study of teacher effects on students' motivation, the whole class of students would need to be excluded from the analysis if the teacher's data are missing. The R package *pan* that uses a multivariate linear mixed effects model (see [Equation 23](#)) is capable only of handling missing data in Level 1 variables (or Level 2 variables that result from aggregating Level 1 variables), but it cannot address missing data that occur at Level 2. [Yucel \(2008\)](#) and [Goldstein et al. \(2014\)](#) developed multilevel MI models that can be used for treating incomplete data at Level 2 (see also [Shin, 2013](#)). In terms of statistical software, *Mplus* and REALCOM-IMPUTE incorporate joint modeling procedures that can address missing data at Level 2 (see also [Enders et al., 2016](#)). The chained equation approach can also be used to impute Level 2 variables in the *mice* package using the `mice.impute.2lonly.norm` function (see also [Yucel, 2008](#)). Clearly, more simulation research is needed to evaluate the performance of these models.

Fourth, we focused only on random-intercept models, which assume that relationships between the variables do not vary across groups. For multilevel models including random slopes (i.e., slopes that are allowed to vary across groups), proper MI can be difficult when values of the covariate are missing. In the imputation model of the software *pan* (see [Equation 23](#)), missing values are allowed only in the multivariate outcome, and predictor variables must be completely observed. We conducted an additional simulation (see the online supplemental materials) in which we evaluated how the PAN approach performs for a random slope model with an incomplete predictor variable. The main finding was that the estimators of the within-group and between-groups regression coefficients are still approximately unbiased, whereas the size of the slope variance was underestimated (see also [Grund, Lüdtke, & Robitzsch, 2016](#)). [Enders et al. \(2016\)](#) discussed a chained equations approach for handling missing values in multilevel models with random slopes. In a simulation study, this approach outperformed the NORM, DI, and PAN approaches of the present study with regard to estimating the slope variance, but still provided negatively biased estimates of the true slope variance. [Yucel \(2011\)](#) presented an adaptation of the multivariate linear mixed effects model of the software *pan* that allows the within-group covariance matrix to vary across groups (see also [Carpenter & Kenward, 2013](#)). However, this approach is not implemented in standard software, and further research is needed to evaluate its performance. In addition, [Graham \(2009\)](#) suggested that MI for multilevel models with random slopes may be carried

⁷ It is worth mentioning that these results also hold for the opposite case—when the model with latent means is the analysis model and the multilevel with manifest group means is used as an imputation model (see the online supplemental materials). It can be concluded that the joint imputation approach in PAN and the chained equations approach (with manifest group means) generate imputations from the same distribution (see also empirical examples in [Enders et al., 2016](#)).

out separately *within* each group. However, as Graham (2012) pointed out, this approach requires that the groups be quite large. Evaluating and developing strategies for dealing with incomplete variables in multilevel models with random slopes is a subject for future research (see also Enders et al., 2016).

Fifth, a further limitation is that the performance of the different MI strategies was only explored with multivariate normally distributed data. It would be important to investigate how robust normal-distribution-based MI strategies are against violations of these assumptions. Previous research has shown that parameter estimates by MI can lead to serious errors of inferences when the assumption of normality is violated, particularly with small sample sizes and a nontrivial proportion of missing data (Demirtas, Freels, & Yucel, 2008; Yuan, Yang-Wallentin, & Bentler, 2012). The bias was particularly pronounced for estimates of variance parameters, and there is some evidence that this also holds for the estimates of variance parameters in multilevel models (see Yucel & Demirtas, 2010).

Sixth, it would also be important to compare the MI strategies with a model-based approach that produces maximum likelihood estimates with incomplete multilevel data in a structural equation modeling framework (Black et al., 2011; Enders, 2010). The software *Mplus* uses a full-information maximum likelihood approach to estimate two-level multilevel structural equation models with incomplete predictor variables (Muthén & Asparouhov, 2011; see also Hox, van Buuren, & Jolani, 2016). However, the model-based approach is limited in its flexibility to include broad sets of auxiliary variables, which are often needed to make the MAR assumption more plausible (see Enders, 2010). Alternatively, two-stage maximum likelihood approaches could be used to estimate two-level structural equation models with missing data (Yuan & Bentler, 2007). Two-stage approaches have the advantage that they can incorporate broad sets of auxiliary variables and also seem to be more robust against violations of the assumption of multivariate normality (e.g., Savalei & Falk, 2014; Yuan, Tong, & Zhang, 2015).

We conclude that although MI is a highly recommended technique for dealing with the issue of missing data, researchers must bear in mind that the imputation model needs to represent the structure of the data. Our comparison of MI strategies for multiply imputing incomplete multilevel data has shown that a multilevel imputation model would be a reasonable choice if one is interested in estimating multilevel random-intercept models with missing values at Level 1.

References

- Allison, P. D. (2009). *Fixed effects regression models*. Thousand Oaks, CA: Sage. <http://dx.doi.org/10.4135/9781412993869>
- Andridge, R. R. (2011). Quantifying the impact of fixed effects modeling of clusters in multiple imputation for cluster randomized trials. *Biometrical Journal*, 53, 57–74. <http://dx.doi.org/10.1002/bimj.201000140>
- Asparouhov, T., & Muthén, B. O. (2010). Multiple imputation with *Mplus*. Retrieved from <http://www.statmodel.com/download/Imputations7.pdf>
- Black, A. C., Harel, O., & McCoach, D. B. (2011). Missing data techniques for multilevel data: Implications of model misspecification. *Journal of Applied Statistics*, 38, 1845–1865. <http://dx.doi.org/10.1080/02664763.2010.529882>
- Bliese, P. D. (2000). Within-group agreement, non-independence, and reliability: Implications for data aggregation and analysis. In K. J. Klein & S. W. Kozlowski (Eds.), *Multilevel theory, research, and methods in organizations* (pp. 349–381). San Francisco, CA: Jossey-Bass.
- Bodner, T. E. (2008). What improves with increased missing data imputations? *Structural Equation Modeling*, 15, 651–675. <http://dx.doi.org/10.1080/10705510802339072>
- Bos, W., Lankes, E.-M., Prenzel, M., Schwippert, K., Walther, G., & Valtin, R. (Eds.). (2003). *Erste Ergebnisse aus IGLU. Schülerleistungen am Ende der vierten Jahrgangsstufe im internationalen Vergleich* [The first results from IGLU - An international comparison of student achievement at the end of fourth grade]. Münster, Germany: Waxmann.
- Carpenter, J. R., Goldstein, H., & Kenward, M. G. (2011). REALCOM-IMPUTE software for multilevel multiple imputation with mixed response types. *Journal of Statistical Software*, 45, 1–14. <http://dx.doi.org/10.18637/jss.v045.i05>
- Carpenter, J. R., & Kenward, M. G. (2013). *Multiple imputation and its application*. Chichester, UK: Wiley. <http://dx.doi.org/10.1002/9781119942283>
- Cohen, J. (1988). *Statistical power for the behavioral sciences*. Hillsdale, NJ: Erlbaum.
- Collins, L. M., Schafer, J. L., & Kam, C.-M. (2001). A comparison of inclusive and restrictive strategies in modern missing data procedures. *Psychological Methods*, 6, 330–351. <http://dx.doi.org/10.1037/1082-989X.6.4.330>
- Cronbach, L. J. (1976). *Research on classrooms and schools: Formulations of questions, design and analysis*. Stanford, CA: Stanford Evaluation Consortium.
- Croon, M. A., & van Veldhoven, M. J. P. M. (2007). Predicting group-level outcome variables from variables measured at the individual level: A latent variable multilevel model. *Psychological Methods*, 12, 45–57. <http://dx.doi.org/10.1037/1082-989X.12.1.45>
- Demirtas, H., Freels, S. A., & Yucel, R. M. (2008). Plausibility of multivariate normality assumption when multiply imputing non-Gaussian continuous outcomes: A simulation assessment. *Journal of Statistical Computation and Simulation*, 78, 69–84. <http://dx.doi.org/10.1080/10629360600903866>
- Dettmers, S., Trautwein, U., Lüdtke, O., Kunter, M., & Baumert, J. (2010). Homework works if homework quality is high: Using multilevel modeling to predict the development of achievement in mathematics. *Journal of Educational Psychology*, 102, 467–482. <http://dx.doi.org/10.1037/a0018453>
- Díaz-Ordaz, K., Kenward, M. G., Cohen, A., Coleman, C. L., & Eldridge, S. (2014). Are missing data adequately handled in cluster randomised trials? A systematic review and guidelines. *Clinical Trials*, 11, 590–600. <http://dx.doi.org/10.1177/1740774514537136>
- Drechsler, J. (2015). Multiple imputation of multilevel missing data—Rigor versus simplicity. *Journal of Educational and Behavioral Statistics*, 40, 69–95. <http://dx.doi.org/10.3102/1076998614563393>
- Enders, C. K. (2010). *Applied missing data analysis*. New York, NY: Guilford Press.
- Enders, C. K., Baraldi, A. N., & Cham, H. (2014). Estimating interaction effects with incomplete predictor variables. *Psychological Methods*, 19, 39–55. <http://dx.doi.org/10.1037/a0035314>
- Enders, C. K., Mistler, S. A., & Keller, B. T. (2016). Multilevel multiple imputation: A review and evaluation of joint modeling and chained equations imputation. *Psychological Methods*, 21, 222–240. <http://dx.doi.org/10.1037/met0000063>
- Gelman, A., Carlin, J. B., Stern, H. S., & Rubin, D. B. (2003). *Bayesian data analysis*. London, UK: CRC Press.
- Gibson, N. M., & Olejnik, S. (2003). Treatment of missing data at the second level of hierarchical linear models. *Educational and Psychological Measurement*, 63, 204–238. <http://dx.doi.org/10.1177/0013164402250987>

- Goldstein, H. (2010). *Multilevel statistical models*. London, UK: Edward Arnold. <http://dx.doi.org/10.1002/9780470973394>
- Goldstein, H., Carpenter, J. R., & Browne, W. J. (2014). Fitting multilevel multivariate models with missing data in responses and covariates that may include interactions and non-linear terms. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, *177*, 553–564. <http://dx.doi.org/10.1111/rssa.12022>
- Graham, J. W. (2009). Missing data analysis: Making it work in the real world. *Annual Review of Psychology*, *60*, 549–576. <http://dx.doi.org/10.1146/annurev.psych.58.110405.085530>
- Graham, J. W. (2012). *Missing data: Analysis and design*. New York, NY: Springer. <http://dx.doi.org/10.1007/978-1-4614-4018-5>
- Graham, J. W., Taylor, B. J., Olchowski, A. E., & Cumsille, P. E. (2006). Planned missing data designs in psychological research. *Psychological Methods*, *11*, 323–343. <http://dx.doi.org/10.1037/1082-989X.11.4.323>
- Grund, S., Lüdtke, O., & Robitzsch, A. (2016). Multiple imputation of missing covariate values in multilevel models with random slopes: A cautionary note. *Behavior Research Methods*, *48*, 640–649. <http://dx.doi.org/10.3758/s13428-015-0590-3>
- Hedges, L. V., & Hedberg, E. C. (2007). Intraclass correlation values for planning group-randomized trials in education. *Educational Evaluation and Policy Analysis*, *29*, 60–87. <http://dx.doi.org/10.3102/0162373707299706>
- Henry, K. L., Stanley, L. R., Edwards, R. W., Harkabus, L. C., & Chapin, L. A. (2009). Individual and contextual effects of school adjustment on adolescent alcohol use. *Prevention Science*, *10*, 236–247. <http://dx.doi.org/10.1007/s11121-009-0124-2>
- Hoff, P. D. (2009). *A first course in Bayesian statistical methods*. New York, NY: Springer. <http://dx.doi.org/10.1007/978-0-387-92407-6>
- Hox, J. J. (2010). *Multilevel analysis: Techniques and applications*. Mahwah, NJ: Erlbaum.
- Hox, J., van Buuren, S., & Jolani, S. (2016). Incomplete multilevel data. In J. Harring, L. M. Stapleton, & S. N. Beretvas (Eds.), *Multilevel modeling for educational research: Addressing practical issues found in real-world applications* (pp. 39–61). New York, NY: Information Age.
- Jelicic, H., Phelps, E., & Lerner, R. M. (2009). Use of missing data methods in longitudinal studies: The persistence of bad practices in developmental psychology. *Developmental Psychology*, *45*, 1195–1199. <http://dx.doi.org/10.1037/a0015665>
- Kenny, D. A., Kashy, D. A., & Cook, W. L. (2006). *The analysis of dyadic data*. New York, NY: Guilford Press.
- Kenny, D. A., Mannetti, L., Pierro, A., Livi, S., & Kashy, D. A. (2002). The statistical analysis of data from small groups. *Journal of Personality and Social Psychology*, *83*, 126–137. <http://dx.doi.org/10.1037/0022-3514.83.1.126>
- Kounin, J. S. (1970). *Discipline and group management in classrooms*. New York, NY: Holt, Rinehart & Winston.
- Little, R. J. A., & Rubin, D. B. (2002). *Statistical analysis with missing data*. New York, NY: Wiley. <http://dx.doi.org/10.1002/9781119013563>
- Lüdtke, O., Marsh, H. W., Robitzsch, A., Trautwein, U., Asparouhov, T., & Muthén, B. (2008). The multilevel latent covariate model: A new, more reliable approach to group-level effects in contextual studies. *Psychological Methods*, *13*, 203–229. <http://dx.doi.org/10.1037/a0012869>
- Maas, C. J. M., & Hox, J. J. (2005). Sufficient sample sizes for multilevel modeling. *Methodology: European Journal of Research Methods for the Behavioral and Social Sciences*, *1*, 86–92. <http://dx.doi.org/10.1027/1614-2241.1.3.86>
- Magnus, J. R., & Neudecker, H. (1988). *Matrix differential calculus with applications in statistics and econometrics*. Hoboken, NJ: Wiley.
- Mathieu, J. E., Aguinis, H., Culppepper, S. A., & Chen, G. (2012). Understanding and estimating the power to detect cross-level interaction effects in multilevel modeling. *Journal of Applied Psychology*, *97*, 951–966. <http://dx.doi.org/10.1037/a0028380>
- Mehta, P. D., & Neale, M. C. (2005). People are variables too: Multilevel structural equations modeling. *Psychological Methods*, *10*, 259–284. <http://dx.doi.org/10.1037/1082-989X.10.3.259>
- Mistler, S. A. (2013). A SAS macro for applying multiple imputation to multilevel data. In *Proceedings of the SAS Global Forum 2013*. Contributed paper (Statistics and Data Analysis) 438–2013. San Francisco, CA. Retrieved from <https://support.sas.com/resources/papers/proceedings13/438-2013.pdf>
- Mistler, S. A. (2015). *Multilevel multiple imputation: An examination of competing methods* (Unpublished doctoral dissertation). Arizona State University, Tempe, AZ.
- Mullis, I. V. S., Martin, M. O., Gonzales, E. J., & Kennedy, A. M. (2003). *PIRLS 2001 international report: IEA's study of reading literacy achievement in primary schools in 35 countries*. Boston, MA: International Study Center, Lynch School of Education, Boston College.
- Muthén, B. O., & Asparouhov, T. (2011). Beyond multilevel regression modeling: Multilevel analysis in a general latent variable framework. In J. J. Hox & J. K. Roberts, *The handbook of advanced multilevel analysis* (pp. 15–40). Milton Park, UK: Routledge.
- Muthén, L. K., & Muthén, B. O. (1998–2010). *Mplus user's guide* (6th ed.). Los Angeles, CA: Author.
- Ostroff, C. (1993). Comparing correlations based on individual-level and aggregated data. *Journal of Applied Psychology*, *78*, 569–582. <http://dx.doi.org/10.1037/0021-9010.78.4.569>
- Peugh, J. L., & Enders, C. K. (2004). Missing data in educational research: A review of reporting practices and suggestions for improvement. *Review of Educational Research*, *74*, 525–556. <http://dx.doi.org/10.3102/00346543074004525>
- Preacher, K. J., Zhang, Z., & Zyphur, M. J. (2016). Multilevel structural equation models for assessing moderation within and across levels of analysis. *Psychological Methods*, *21*, 189–205. <http://dx.doi.org/10.1037/met0000052>
- Preacher, K. J., Zyphur, M. J., & Zhang, Z. (2010). A general multilevel SEM framework for assessing multilevel mediation. *Psychological Methods*, *15*, 209–233. <http://dx.doi.org/10.1037/a0020141>
- Quartagno, M., & Carpenter, J. (2016). jomo: A package for multilevel joint modelling multiple imputation. R package version 2.1–2. Retrieved from <http://CRAN.R-project.org/package=jomo>
- Raghunathan, T. E., Lepkowski, J. E., Hoewyk, J. V., & Solenberger, P. (2001). A multivariate technique for multiply imputing missing values using a sequence of regression models. *Survey Methodology*, *27*, 85–96.
- Raudenbush, S. W., & Bryk, A. S. (2002). *Hierarchical linear models* (2nd ed.). Thousand Oaks, CA: Sage.
- Rubin, D. B. (1976). Inference and missing data. *Biometrika*, *63*, 581–592. <http://dx.doi.org/10.1093/biomet/63.3.581>
- Rubin, D. B. (1987). *Multiple imputation for nonresponse in surveys*. Hoboken, NJ: Wiley. <http://dx.doi.org/10.1002/9780470316696>
- Savalei, V., & Falk, C. (2014). Robust two-stage approach outperforms robust full information maximum likelihood with incomplete nonnormal data. *Structural Equation Modeling*, *21*, 280–302. <http://dx.doi.org/10.1080/10705511.2014.882692>
- Schafer, J. L. (1997). *Analysis of incomplete multivariate data*. Boca Raton, FL: CRC Press.
- Schafer, J. L. (2001). Multiple imputation with PAN. In L. M. Collins & A. G. Sayer (Eds.), *New methods for the analysis of change* (pp. 357–377). Washington, DC: American Psychological Association. <http://dx.doi.org/10.1037/10409-012>
- Schafer, J. L., & Graham, J. W. (2002). Missing data: Our view of the state of the art. *Psychological Methods*, *7*, 147–177.
- Schafer, J. L., & Olsen, M. K. (1998). Multiple imputation for multivariate missing-data problems: A data analyst's perspective. *Multivariate Behavioral Research*, *33*, 545–571. http://dx.doi.org/10.1207/s15327906mbr3304_5

- Schafer, J. L., & Yucel, R. M. (2002). Computational strategies for multivariate linear mixed-effects models with missing values. *Journal of Computational and Graphical Statistics, 11*, 437–457. <http://dx.doi.org/10.1198/106186002760180608>
- Schafer, J. L., & Zhao, J. H. (2013). pan: Multiple imputation for multivariate panel or clustered data (R package version 0.9) [Computer program]. Retrieved from <http://CRAN.R-project.org/package=pan>
- Shin, Y. (2013). Efficient handling of predictors and outcomes having missing values. In L. Rutkowski, M. von Davier, & D. Rutkowski (Eds.), *A handbook of international large-scale assessment data analysis: Background, technical issues, and methods of data analysis* (pp. 451–479). London, UK: Chapman & Hall/CRC Press.
- Shin, Y., & Raudenbush, S. W. (2010). A latent cluster-mean approach to the contextual effects model with missing data. *Journal of Educational and Behavioral Statistics, 35*, 26–53. <http://dx.doi.org/10.3102/1076998609345252>
- Snijders, T. A. B., & Bosker, R. J. (2012). *Multilevel analysis: An introduction to basic and advanced multilevel modeling*. London, UK: Sage.
- Taljaard, M., Donner, A., & Klar, N. (2008). Imputation strategies for missing continuous outcomes in cluster randomized trials. *Biometrical Journal, 50*, 329–345. <http://dx.doi.org/10.1002/bimj.200710423>
- Tanner, M. A., & Wong, W. H. (1987). The calculation of posterior distributions by data augmentation. *Journal of the American Statistical Association, 82*, 528–540. <http://dx.doi.org/10.1080/01621459.1987.10478458>
- van Buuren, S. (2011). Multiple imputation of multilevel data. In J. J. Hox & J. K. Roberts, *The handbook of advanced multilevel analysis* (pp. 173–196). Milton Park, UK: Routledge.
- van Buuren, S. (2012). *Flexible imputation of missing data*. Boca Raton, FL: Chapman & Hall/CRC Press. <http://dx.doi.org/10.1201/b11826>
- van Buuren, S., & Groothuis-Oudshoorn, K. (2011). MICE: Multivariate imputation by chained equations in R. *Journal of Statistical Software, 45*, 1–67.
- von Hippel, P. T. (2009). How to impute square, interactions, and other transformed variables. *Sociological Methodology, 39*, 265–291. <http://dx.doi.org/10.1111/j.1467-9531.2009.01215.x>
- Walsh, B. M., Matthews, R. A., Tuller, M. D., Parks, K. M., & McDonald, D. P. (2010). A multilevel model of the effects of equal opportunity climate on job satisfaction in the military. *Journal of Occupational Health Psychology, 15*, 191–207. <http://dx.doi.org/10.1037/a0018756>
- West, S. G. (2001). New approaches to missing data in psychological research: Introduction to the special section. *Psychological Methods, 6*, 315–316. <http://dx.doi.org/10.1037/1082-989X.6.4.315>
- White, I. R., Royston, P., & Wood, A. M. (2011). Multiple imputation using chained equations: Issues and guidance for practice. *Statistics in Medicine, 30*, 377–399. <http://dx.doi.org/10.1002/sim.4067>
- Yuan, K.-H., & Bentler, P. M. (2007). Multilevel covariance structure analysis by fitting multiple single-level models. *Sociological Methodology, 37*, 53–82. <http://dx.doi.org/10.1111/j.1467-9531.2007.00182.x>
- Yuan, K.-H., Tong, X., & Zhang, Z. (2015). Bias and efficiency for SEM with missing data and auxiliary variables: Two-stage robust method versus two-stage ML. *Structural Equation Modeling, 22*, 178–192. <http://dx.doi.org/10.1080/10705511.2014.935750>
- Yuan, K.-H., Yang-Wallentin, F., & Bentler, P. M. (2012). ML versus MI for missing data with violation of distribution conditions. *Sociological Methods & Research, 41*, 598–629. <http://dx.doi.org/10.1177/0049124112460373>
- Yucel, R. M. (2008). Multiple imputation inference for multivariate multilevel continuous data with ignorable non-response. *Philosophical Transactions. Series A, Mathematical, Physical, and Engineering Sciences, 366*, 2389–2403. <http://dx.doi.org/10.1098/rsta.2008.0038>
- Yucel, R. M. (2011). Random-covariances and mixed-effects models for imputing multivariate multilevel continuous data. *Statistical Modelling, 11*, 351–370. <http://dx.doi.org/10.1177/1471082X1001100404>
- Yucel, R. M., & Demirtas, H. (2010). Impact of non-normal random effects on inference by multiple imputation: A simulation assessment. *Computational Statistics & Data Analysis, 54*, 790–801. <http://dx.doi.org/10.1016/j.csda.2009.01.016>

Appendix

Derivation of Bias for Ad Hoc Multiple Imputation Strategies

In this Appendix, we derive the asymptotic bias for the NORM approach and the DI approach, when estimating the within-group variance, the between-groups variance, the intraclass correlation, and the within- and between-groups regression coefficients from incomplete multilevel data. We assume that the data have a two-level structure with two mean-centered variables X and Y (see Equations 3 and 4). The values in Y are MCAR and X is fully observed. In addition, we assume that in each group n_1 persons have observed values and $n_0 = n - n_1$ have missing values for Y . With no loss of generality, we assume that the first n_1 values in a group are observed and that the other n_0 values are missing. In the following, the given data (X_{ij}, Y_{ij}) are denoted as $(X_{ij(obs)}, Y_{ij(obs)})$, if Y_{ij} is observed and $(X_{ij(mis)}, Y_{ij(mis)})$, if Y_{ij} is missing. Missing values are replaced with the imputed values $Y_{ij(imp)}$ and the completed data are denoted as Y_{ij}^* . If Y_{ij} is observed, then $Y_{ij}^* = Y_{ij(obs)}$, otherwise $Y_{ij}^* = Y_{ij(imp)}$. Finally, the following derivations are based on the assumption that the number of groups approaches infinity ($K \rightarrow \infty$).

NORM Approach

The NORM approach ignores the multilevel structure of the data and uses a simple regression for imputing the missing values

$$Y_{ij} = \alpha + \beta_{total,YX}X_{ij} + e_{ij}, \quad Var(e_{ij}) = \sigma_e^2, \quad (A.1)$$

where the regression coefficient $\beta_{total,YX}$ indicates the total relationship between X and Y (i.e., ignoring the grouped data structure) and the residual variance σ_e^2 is assumed to be homogenous. As $E(X_{ij}) = 0$, the intercept α is estimated to be zero if the number of groups approaches infinity ($K \rightarrow \infty$). By using the MCAR assumption, the regression coefficient $\beta_{total,YX}$ is given by

$$\begin{aligned} E(\hat{\beta}_{total,YX}) &= \frac{Cov(X_{ij(obs)}, Y_{ij(obs)})}{Var(X_{ij(obs)})} \\ &= \frac{Cov(X_{ij}, Y_{ij})}{Var(X_{ij})} \\ &= \frac{Cov(X_{B,ij}, Y_{B,ij}) + Cov(X_{W,ij}, Y_{W,ij})}{Var(X_{B,ij}) + Var(X_{W,ij})} \\ &= \frac{\rho_B \tau_X \tau_Y + \rho_W \sigma_X \sigma_Y}{\tau_X^2 + \sigma_X^2} \\ &= \frac{\tau_X^2 \beta_{B,YX} + \sigma_X^2 \beta_{W,YX}}{\tau_X^2 + \sigma_X^2}. \end{aligned} \quad (A.2)$$

The residual variance is estimated by the following expression

$$\hat{\sigma}_e^2 = \frac{1}{Kn_1 - 2} \sum_{j=1}^K \sum_{i=1}^{n_1} (Y_{ij(obs)} - \hat{Y}_{ij(obs)})^2. \quad (A.3)$$

We can now derive

$$\begin{aligned} E(\hat{\sigma}_e^2) &= E(Y_{ij(obs)} - \beta_{total,YX}X_{ij(obs)})^2 \\ &= E(Y_{ij} - \beta_{total,YX}X_{ij})^2 \\ &= \tau_Y^2 + \sigma_Y^2 - \beta_{total,YX}^2(\tau_Y^2 + \sigma_Y^2). \end{aligned} \quad (A.4)$$

Thus, the imputed values are given by

$$Y_{ij(imp)} = \beta_{total,YX}X_{ij(mis)} + e_{ij}, \quad e_{ij} \sim N(0, \sigma_e^2). \quad (A.5)$$

The within-group variance σ_Y^2 can be estimated using the observed within-group variance of the completed data (Snijders & Bosker, 2012):

$$\hat{\sigma}_Y^2 = S_{within,Y}^2 = \frac{1}{K(n-1)} \sum_{j=1}^K \left(\sum_{i=1}^{n_1} (Y_{ij(obs)} - \bar{Y}_{\bullet j})^2 + \sum_{i=n_1+1}^{n_1+n_0} (Y_{ij(imp)} - \bar{Y}_{\bullet j})^2 \right). \quad (A.6)$$

For the estimator of the between-groups variance τ_Y^2 , we need the observed between-groups variance of the completed data

$$S_{between,Y}^2 = \frac{1}{K-1} \sum_{j=1}^K (\bar{Y}_{\bullet j}^* - \bar{Y}_{\bullet\bullet})^2. \quad (A.7)$$

The estimator of the between-groups variance is then given by

$$\hat{\tau}_Y^2 = S_{between,Y}^2 - S_{within,Y}^2/n. \quad (A.8)$$

We now show that for the NORM approach the following relation holds: $Bias(\hat{\tau}_Y^2) = -Bias(\hat{\sigma}_Y^2)$. We first write the total sum of squares for the completed data

$$\begin{aligned} SS_{total,Y} &= \sum_{j=1}^K \sum_{i=1}^n (Y_{ij}^* - \bar{Y}_{\bullet\bullet})^2 \\ &= \sum_{j=1}^K \sum_{i=1}^{n_1} (Y_{ij}^* - \bar{Y}_{\bullet\bullet})^2 + \sum_{j=1}^K \sum_{i=1}^{n_0} (Y_{ij}^* - \bar{Y}_{\bullet j})^2 \\ &= SS_{between,Y} + SS_{within,Y}. \end{aligned} \quad (A.9)$$

Using that $SS_{within,Y} = K(n-1)S_{within,Y}^2$ and $SS_{between,Y} = (K-1)nS_{between,Y}^2$, the following relationship holds:

$$\frac{E(SS_{total,Y})}{Kn} = E(Y_{ij}^*)^2 = \frac{K-1}{K} \cdot E(S_{between,Y}^2) + \frac{n-1}{n} E(S_{within,Y}^2). \quad (A.10)$$

For a large number of groups ($K \rightarrow \infty$), this reduces to

$$E(Y_{ij}^*)^2 = E(S_{between,Y}^2) + (1 - 1/n)E(S_{within,Y}^2). \quad (A.11)$$

(Appendix continues)

Note that the imputed values in the NORM approach preserve the total variance of Y :

$$\text{Var}(Y_{ij(\text{imp})}) = E(Y_{ij(\text{imp})}^*)^2 = \beta_{\text{total}, YX}^2(\tau_X^2 + \sigma_X^2) + \sigma_e^2 = \tau_Y^2 + \sigma_Y^2. \tag{A.12}$$

Thus, it follows that

$$\tau_Y^2 + \sigma_Y^2 = E(\hat{\tau}_Y^2) + (1/n)E(\hat{\sigma}_Y^2) + (1 - 1/n)E(\hat{\sigma}_Y^2) = E(\hat{\tau}_Y^2) + E(\hat{\sigma}_Y^2). \tag{A.13}$$

Based on this relation, it can be concluded that

$$\text{Bias}(\hat{\tau}_Y^2) = -\text{Bias}(\hat{\sigma}_Y^2). \tag{A.14}$$

We now derive the bias for the estimator of the between-groups variance $\hat{\tau}_Y^2$ using the following relationship:

$$E(S_{\text{between}, Y}^2) = E(\hat{\tau}_Y^2) + (1/n)E(\hat{\sigma}_Y^2) = \tau_Y^2 + (1/n)\sigma_Y^2 + (1 - 1/n)\text{Bias}(\hat{\tau}_Y^2). \tag{A.15}$$

Rearranging terms and solving for the bias term, we obtain

$$\text{Bias}(\hat{\tau}_Y^2) = n(n - 1)^{-1}(E(S_{\text{between}, Y}^2) - \tau_Y^2 - (1/n)\sigma_Y^2). \tag{A.16}$$

For taking the expectation of $S_{\text{between}, Y}^2$, we use the following relationship:

$$E(\bar{Y}_{\bullet j}^* - Y_{\bullet j}^*)^2 = E\left(\frac{K-1}{K}\bar{Y}_{\bullet j}^* - \frac{1}{K}\sum_{k \neq j} Y_{\bullet k}^*\right)^2 = \frac{K-1}{K}E(\bar{Y}_{\bullet j}^*)^2. \tag{A.17}$$

The expectation is now given by $E(S_{\text{between}, Y}^2) = E(\bar{Y}_{\bullet j}^*)^2$, and for the average of the completed data, we consider the sum

$$\sum_{i=1}^n Y_{ij}^* = n_1 Y_{B,j} + \sum_{i=1}^{n_1} Y_{W,ij(\text{obs})} + n_0 \beta_{\text{total}, YX} X_{B,j} + \sum_{i=n_1+1}^{n_1+n_0} (\beta_{\text{total}, YX} X_{W,ij(\text{mis})} + e_{ij}). \tag{A.18}$$

Then it can be shown that

$$\begin{aligned} n^2 E(\bar{Y}_{\bullet j}^*)^2 &= \tau_X^2 (n_1 \beta_{B, YX} + n_0 \beta_{\text{total}, YX})^2 + n_1^2 \tau_Y^2 (1 - \rho_B^2) \\ &\quad + n_1 \sigma_Y^2 + n_0 \beta_{\text{total}, YX}^2 \sigma_X^2 + n_0 \sigma_e^2 \\ &= \tau_X^2 (n_1 \beta_{B, YX} + n_0 \beta_{\text{total}, YX})^2 + n_1^2 \tau_Y^2 (1 - \rho_B^2) \\ &\quad + n \sigma_Y^2 + n_0 \tau_Y^2 - n_0 \beta_{\text{total}, YX}^2 \tau_X^2 \end{aligned} \tag{A.20}$$

This is used to show that the following relation holds:

$$\begin{aligned} n^2 E(\bar{Y}_{\bullet j}^*)^2 - n^2 \tau_Y^2 - n \sigma_Y^2 &= \tau_X^2 (n_1 \beta_{B, YX} + n_0 \beta_{\text{total}, YX})^2 \\ &\quad + n_1^2 \tau_Y^2 (1 - \rho_B^2) + n_0 \tau_Y^2 - n_0 \beta_{\text{total}, YX}^2 \tau_X^2 - n^2 \tau_Y^2 \end{aligned}$$

$$\begin{aligned} &= \tau_Y^2 \left\{ \rho_B^2 \left[\left(n_1 + n_0 \frac{\beta_{\text{total}, YX}}{\beta_{B, YX}} \right)^2 + n_0 - n_0 \left(\frac{\beta_{\text{total}, YX}}{\beta_{B, YX}} \right)^2 - n^2 \right] \right. \\ &\quad \left. + (1 - \rho_B^2)(n_1^2 + n_0 - n^2) \right\}. \end{aligned} \tag{A.21}$$

Inserting Equation A.21 into Equation A.16, and using $p_0 = n_0/n$, the bias of the estimator of the between-groups variance is given by

$$\begin{aligned} \text{Bias}(\hat{\tau}_Y^2) &= -p_0 \tau_Y^2 \frac{n}{n-1} \left\{ \rho_B^2 (1 - \rho_{l, X})(\beta_{B, YX} - \beta_{W, YX}) A_X \right. \\ &\quad \left. + (1 - \rho_B^2) A_e \right\}, \end{aligned} \tag{A.22}$$

where $A_X \equiv \{2(1 - p_0) + (p_0 - 1/n)(\beta_{\text{total}, YX}/\beta_{B, YX} + 1)\}/\beta_{B, YX}$ and $A_e \equiv 2 - 1/n - p_0$.

The bias for the estimator of the intraclass correlation of Y can be written as

$$\text{Bias}(\hat{\rho}_{l, Y}) = \frac{\text{Bias}(\hat{\tau}_Y^2)}{\tau_Y^2 + \sigma_Y^2} = \rho_{l, Y} \cdot \frac{\text{Bias}(\hat{\tau}_Y^2)}{\tau_Y^2}. \tag{A.23}$$

Then for the bias of the estimator of the within-group variance, we can use Equation A.14 to show that $\text{Bias}(\hat{\sigma}_Y^2) = -\text{Bias}(\hat{\tau}_Y^2)$.

In the next step, we investigate the bias of the estimator of the within-group regression coefficient. The estimator for the within-group covariance $\sigma_{W, XY}$ is given as follows:

$$\begin{aligned} C_{\text{within}} &= \frac{1}{K(n-1)} \sum_{j=1}^K \left(\sum_{i=1}^{n_1} (X_{ij(\text{obs})} - \bar{X}_{\bullet j})(Y_{ij(\text{obs})} - \bar{Y}_{\bullet j}^*) \right. \\ &\quad \left. + \sum_{i=n_1+1}^{n_1+n_0} (X_{ij(\text{mis})} - \bar{X}_{\bullet j})(Y_{ij(\text{mis})} - \bar{Y}_{\bullet j}^*) \right). \end{aligned} \tag{A.24}$$

For the expectations of the single cross-products of the observed values the following relationships hold:

$$\begin{aligned} E(X_{ij(\text{obs})} Y_{ij(\text{obs})}) &= E(X_{B,j} Y_{B,j}) + E(X_{W,ij(\text{obs})} Y_{Y,ij(\text{obs})}) \\ &= \beta_{B, YX} \tau_X^2 + \beta_{W, YX} \sigma_X^2 \end{aligned} \tag{A.25}$$

$$\begin{aligned} E(\bar{X}_{\bullet j} Y_{ij(\text{obs})}) &= E(X_{B,j} Y_{B,j}) + (1/n) E(X_{W,ij(\text{obs})} Y_{Y,ij(\text{obs})}) \\ &= \beta_{B, YX} \tau_X^2 + (1/n) \beta_{W, YX} \sigma_X^2 \end{aligned} \tag{A.26}$$

$$E(X_{ij(\text{obs})} \bar{Y}_{\bullet j}^*) = (1/n) \tau_X^2 (n_1 \beta_{B, YX} + n_0 \beta_{\text{total}, YX}) + (1/n) \sigma_X^2 \beta_{W, YX} \tag{A.27}$$

$$\begin{aligned} E(\bar{X}_{\bullet j} \bar{Y}_{\bullet j}^*) &= (1/n) \tau_X^2 (n_1 \beta_{B, YX} + n_0 \beta_{\text{total}, YX}) + (n_1/n^2) \sigma_X^2 \beta_{W, YX} \\ &\quad + (n_0/n^2) \beta_{\text{total}, YX} \sigma_X^2. \end{aligned} \tag{A.28}$$

Using these relationships, the expectation for the cross-product of the observed values is given by

$$\begin{aligned} E(X_{ij(\text{obs})} - \bar{X}_{\bullet j})(Y_{ij(\text{obs})} - \bar{Y}_{\bullet j}^*) &= (1 - 1/n) \beta_{W, YX} \sigma_X^2 \\ &\quad + (1/n^2) \sigma_X^2 n_0 (\beta_{\text{total}, YX} - \beta_{W, YX}). \end{aligned} \tag{A.29}$$

(Appendix continues)

In a similar way, the expectations of the single cross-product terms for the imputed values are given by

$$E(X_{ij(mis)}Y_{ij(imp)}) = \beta_{total,YX}\tau_X^2 + \beta_{total,YX}\sigma_X^2 \quad (A.30)$$

$$E(X_{ij(mis)}\bar{Y}_{\bullet j}^*) = (1/n)\tau_X^2(n_1\beta_{B,YX} + n_0\beta_{total,YX}) + (1/n)\sigma_X^2\beta_{total,YX} \quad (A.31)$$

$$E(\bar{X}_{\bullet j}Y_{ij(imp)}) = \beta_{total,YX}\tau_X^2 + (1/n)\beta_{total,YX}\sigma_X^2. \quad (A.32)$$

Using these relationships, together with Equation A.28 the expectation of C_{within} yields

$$\begin{aligned} E(C_{within}) &= \frac{n_1\beta_{W,YX} + n_0\beta_{total,YX}}{n}\sigma_X^2 \\ &= \beta_{W,YX}\sigma_X^2 + \frac{n_0}{n}(\beta_{total,YX} - \beta_{W,YX})\sigma_X^2. \end{aligned} \quad (A.33)$$

The bias of the estimator of the within-group covariance $\sigma_{W,XY}$ is then given by

$$Bias(\hat{\sigma}_{W,XY}) = p_0 \cdot \rho_{LX} \cdot (\beta_{B,YX} - \beta_{W,YX}) \cdot \sigma_X^2. \quad (A.34)$$

The bias of the estimator of the within-group coefficient $\beta_{W,YX}$ can now be expressed as follows:

$$Bias(\hat{\beta}_{W,YX}) = p_0 \cdot \rho_{LX} \cdot (\beta_{B,YX} - \beta_{W,YX}). \quad (A.35)$$

For the regression of X on Y , the bias of the estimator of the within-group coefficient $\beta_{W,XY}$ can be written as a function of the biases of the estimators of the within-group covariance and the within-group variance of Y :

$$Bias(\hat{\beta}_{W,XY}) = \frac{Bias(\hat{\sigma}_{W,XY}) - \beta_{W,XY}Bias(\hat{\sigma}_Y^2)}{\sigma_Y^2 + Bias(\hat{\sigma}_Y^2)}. \quad (A.36)$$

For investigating the between-groups covariance $\sigma_{B,XY}$, we first define the covariance of the observed group means

$$C_{between} = \frac{1}{K-1} \sum_{j=1}^K (\bar{X}_{\bullet j} - \bar{X}_{\bullet\bullet})(\bar{Y}_{\bullet j} - \bar{Y}_{\bullet\bullet}). \quad (A.37)$$

The estimator of the between-groups covariance $\sigma_{B,XY}$ is then given by

$$\hat{\sigma}_{B,XY} = C_{between} - C_{within}/n. \quad (A.38)$$

For the expectation of $C_{between}$, the following relation holds:

$$\begin{aligned} E(C_{between}) &= \frac{1}{K-1} \sum_{j=1}^K E(\bar{X}_{\bullet j} - \bar{X}_{\bullet\bullet})(\bar{Y}_{\bullet j} - \bar{Y}_{\bullet\bullet}) \\ &= \frac{K}{K-1} \cdot \left(1 - \frac{1}{K}\right) \cdot E(\bar{X}_{\bullet j}\bar{Y}_{\bullet j}) \\ &= E(\bar{X}_{\bullet j}\bar{Y}_{\bullet j}^*) \\ &= (1/n)\tau_X^2(n_1\beta_{B,YX} + n_0\beta_{total,YX}) + (n_1/n^2)\beta_{W,YX}\sigma_X^2 \\ &\quad + (n_0/n^2)\beta_{total,YX}\sigma_X^2 \end{aligned} \quad (A.39)$$

The expectation of the estimator of the between-groups covariance is then given by

$$\begin{aligned} E(\hat{\sigma}_{B,XY}) &= E(C_{between}) - (1/n)E(C_{within}) \\ &= \beta_{B,YX}\tau_X^2 + \tau_X^2 n_0(\beta_{total,YX} - \beta_{B,YX})(1/n). \end{aligned} \quad (A.40)$$

Thus, the bias can be calculated as

$$Bias(\hat{\sigma}_{B,XY}) = -p_0 \cdot (1 - \rho_{LX}) \cdot (\beta_{B,YX} - \beta_{W,YX}) \cdot \tau_X^2. \quad (A.41)$$

The bias of the estimator of the between-groups coefficient $\beta_{B,XY}$ can be written as

$$Bias(\hat{\beta}_{B,XY}) = -p_0 \cdot (1 - \rho_{LX}) \cdot (\beta_{B,YX} - \beta_{W,YX}). \quad (A.42)$$

We now derive the expectation of the estimator of the between-groups coefficient $\beta_{B,XY}$, using the bias for the estimator of the between-groups covariance and the between-groups variance

$$\begin{aligned} E(\hat{\beta}_{B,XY}) &= \frac{\sigma_{B,XY} + Bias(\hat{\sigma}_{B,XY})}{\tau_Y^2 + Bias(\hat{\tau}_Y^2)} \\ &= \beta_{B,XY} + \frac{Bias(\hat{\sigma}_{B,XY}) - \beta_{B,XY}Bias(\hat{\tau}_Y^2)}{\tau_Y^2 + Bias(\hat{\tau}_Y^2)}. \end{aligned} \quad (A.43)$$

The bias is then given by

$$Bias(\hat{\beta}_{B,XY}) = \frac{Bias(\hat{\sigma}_{B,XY}) - \beta_{B,XY} \cdot Bias(\hat{\tau}_Y^2)}{\tau_Y^2 + Bias(\hat{\tau}_Y^2)}. \quad (A.44)$$

DI Approach

In the DI approach, dummy variables for the groups are included in the imputation model. The regression $Y_{ij} = \alpha_j + \beta X_{ij} + e_{ij}$ is used for imputing missing values in Y , where α_j is a group-specific fixed effect. The regression coefficient β consistently estimates the within-group coefficient $\beta_{W,YX}$, when the number of groups approaches infinity. The residual variance is assumed to be homogeneous $Var(e_{ij}) = \sigma_e^2$. The group-specific fixed effects α_j are estimated as follows:

$$\begin{aligned} \hat{\alpha}_j &= \bar{Y}_{\bullet j(obs)} - \beta \bar{X}_{\bullet j(obs)} \\ &= Y_{B,j} + (1/n_1) \sum_{i=1}^{n_1} Y_{W,ij(obs)} - \beta_{W,YX} \left(X_{B,j} + (1/n_1) \sum_{i=1}^{n_1} X_{W,ij(obs)} \right) \\ &= Y_{B,j} - \beta_{W,YX} X_{B,j} + (1/n_1) \sum_{i=1}^{n_1} (Y_{W,ij(obs)} - \beta_{W,YX} X_{W,ij(obs)}) \end{aligned} \quad (A.45)$$

(Appendix continues)

Thus, the expected value of a group-specific effect is given by $E(\hat{\alpha}_j) = Y_{B,j} - \beta_{W,YX}X_{B,j}$, and its variance is $Var(\hat{\alpha}_j) = (1/n_1)\sigma_Y^2(1 - \rho_W^2)$. Furthermore, it can be shown that the estimator of the residual variance in the regression provides an unbiased estimate of the within-group variance $E(\hat{\sigma}_e^2) = \sigma_Y^2(1 - \rho_W^2) = \sigma_{W,YX,e}^2$, if the number of groups is large. The imputed values $Y_{ij(imp)}$ are then generated by a regression with sampled group-specific effects $\alpha_j^* = \hat{\alpha}_j + v_{\alpha_j}$ with $v_{\alpha_j} \sim N(0, Var(\hat{\alpha}_j))$ and normally distributed residuals $e_{ij} \sim N(0, \sigma_{W,YX,e}^2)$:

$$\begin{aligned} Y_{ij(imp)} &= \alpha_j^* + \beta X_{ij(mis)} + e_{ij} \\ &= Y_{B,j} + v_{\alpha_j} + \beta_{W,YX}X_{W,ij(mis)} \\ &\quad + (1/n_1) \sum_{i=1}^{n_1} (Y_{W,ij(obs)} - \beta_{W,YX}X_{W,ij(obs)}) + e_{ij}. \end{aligned} \quad (A.46)$$

Again, we first calculate the bias for the estimator of the within-group variance $S_{within,Y}^2$. The group average of the completed data $\bar{Y}_{\bullet j}$ is given by

$$\begin{aligned} \bar{Y}_{\bullet j} &= Y_{B,j} + \frac{n_0}{n} v_{\alpha_j} + \frac{1}{n} \sum_{i=1}^{n_1} Y_{W,ij(obs)} + \frac{1}{n} \sum_{i=n_1+1}^{n_1+n_0} \beta_{W,YX} X_{W,ij(mis)} \\ &\quad + \frac{1}{n} \sum_{i=n_1+1}^{n_1+n_0} e_{ij} + \frac{n_0}{nn_1} \sum_{i=1}^{n_1} (Y_{W,ij(obs)} - \beta_{W,YX} X_{W,ij(obs)}) \\ &= Y_{B,j} + (n_0/n) v_{\alpha_j} + (1/n) \sum_{i=1}^{n_1} \beta_{W,YX} X_{W,ij} \\ &\quad + (1/n_1) \sum_{i=1}^{n_1} \varepsilon_{W,Yij} + (1/n) \sum_{i=n_1+1}^{n_1+n_0} e_{ij} \end{aligned} \quad (A.47)$$

where $\varepsilon_{W,Yij}$ is the residual of the within-group regression of Y on X .

Then the variance of the average of the completed data can be calculated as

$$E(\bar{Y}_{\bullet j})^2 = \tau_Y^2 + (1/n)\sigma_Y^2 + (2n_0/nn_1)\sigma_Y^2(1 - \rho_W^2). \quad (A.48)$$

The squared deviation of an observed value $Y_{ij(obs)}$ in a group is given by

$$E(Y_{ij(obs)} - \bar{Y}_{\bullet j})^2 = E(Y_{ij(obs)}^2) + E(\bar{Y}_{\bullet j})^2 - 2E(Y_{ij(obs)}\bar{Y}_{\bullet j}). \quad (A.49)$$

For the variance of an observed value, we write

$$E(Y_{ij(obs)}^2) = Var(Y_{ij(obs)}) = Var(Y_{B,j}) + Var(Y_{W,ij}) = \tau_Y^2 + \sigma_Y^2. \quad (A.50)$$

The covariance is given by

$$\begin{aligned} E(Y_{ij(obs)}\bar{Y}_{\bullet j}) &= E(Y_{B,j}^2) + (1/n)E[(\beta_{W,YX}X_{W,ij(obs)})^2] \\ &\quad + (1/n_1)E(e_{W,Yij}^2) \\ &= \tau_Y^2 + (1/n)\sigma_Y^2\rho_W^2 + (1/n_1)\sigma_Y^2(1 - \rho_W^2). \end{aligned} \quad (A.51)$$

Using Equations A.48, A.50, and A.51, we can write

$$E(Y_{ij(obs)} - \bar{Y}_{\bullet j})^2 = (1 - 1/n)\sigma_Y^2. \quad (A.52)$$

The variance of an imputed value is given by

$$E(Y_{ij(imp)}^2) = Var(Y_{ij(imp)}) = \tau_Y^2 + \sigma_Y^2 + (2/n_1)\sigma_Y^2(1 - \rho_W^2). \quad (A.53)$$

The covariance of an imputed value with the mean of the completed data in a group can be calculated as follows:

$$\begin{aligned} E(Y_{ij(imp)}\bar{Y}_{\bullet j}) &= E(Y_{B,j}^2) + (n_0/n)E(V_{\alpha_j}^2) + (1/n)E(e_{ij}^2) \\ &\quad + (1/n)E[(\beta_{W,YX}X_{W,ij(mis)})^2] + (1/n_1)\sigma_Y^2(1 - \rho_W^2) \\ &= \tau_Y^2 + (2/n_1)\sigma_Y^2 - (1/n)\sigma_Y^2\rho_W^2 - (2n_0/nn_1)\sigma_Y^2\rho_W^2 \end{aligned} \quad (A.54)$$

Using Equations A.48, A.53, and A.54, we can write for the squared deviation of an imputed value

$$E(Y_{ij(imp)} - \bar{Y}_{\bullet j})^2 = (1 - 1/n)\sigma_Y^2. \quad (A.55)$$

Now the sum for all squared deviations is calculated by combining Equations A.52 and A.55:

$$\sum_{i=1}^n E(Y_{ij}^* - \bar{Y}_{\bullet j})^2 = (n - 1)\sigma_Y^2. \quad (A.56)$$

Thus, the DI approach provides an unbiased estimator of the within-group variance $E(S_{within,Y}^2) = \sigma_Y^2$. For the expectation of the observed between-groups variance $S_{between,Y}^2$, we use Equations A.17 and A.48. The bias for the estimator of the between-groups variance can now be calculated as follows:

$$Bias(\hat{\tau}_Y^2) = E\left(S_{between,Y}^2 - \frac{S_{within,Y}^2}{n}\right) - \hat{\tau}_Y^2 = \frac{\sigma_Y^2}{n} \cdot \frac{2p_0}{1 - p_0} \cdot (1 - \rho_W^2). \quad (A.57)$$

In order to obtain the bias for the intraclass correlation of Y , we first write

$$\begin{aligned} Bias(\hat{\rho}_{I,Y}) &= \frac{\tau_Y^2 + \sigma_Y^2 \cdot r}{\tau_Y^2 + \sigma_Y^2 + \sigma_Y^2 \cdot r} - \frac{\tau_Y^2}{\tau_Y^2 + \sigma_Y^2} \\ &= (1 - \rho_{I,Y})^2 \cdot \left[1 - \frac{\sigma_Y^2 \cdot r}{\tau_Y^2 + \sigma_Y^2 + \sigma_Y^2 \cdot r}\right] \cdot r, \end{aligned} \quad (A.58)$$

(Appendix continues)

where we define $r = \frac{1}{n} \cdot \frac{2p_0}{1-p_0} \cdot (1 - \rho_W^2)$. Expanding the factors in squared brackets in Equation A.58 and neglecting the second term because it is of power n^2 leads to

$$Bias(\hat{\rho}_{LY}) \approx (1 - \rho_{LY})^2 \cdot r = (1 - \rho_{LY})^2 \cdot \frac{1}{n} \cdot \frac{2p_0}{1-p_0} \cdot (1 - \rho_W^2). \tag{A.59}$$

For the bias of the within- and between-groups regression coefficients, we start again with C_{within} . First, we show for the cross-product $(X_{ij(obs)} - \bar{X}_{\bullet j})(Y_{ij(obs)} - \bar{Y}_{\bullet j})$ involving observed values that the following relationships hold:

$$\begin{aligned} E(X_{ij(obs)}Y_{ij(obs)}) &= E(X_{B,j}Y_{B,j}) + E(X_{W,ij(obs)}Y_{W,ij(obs)}) \\ &= \rho_B\tau_X\tau_Y + \rho_W\sigma_X\sigma_Y \end{aligned} \tag{A.60}$$

$$\begin{aligned} E(X_{ij(obs)}\bar{Y}_{\bullet j}^*) &= E(X_{B,j}Y_{B,j}) + (1/n)E(X_{W,ij(obs)}X_{W,ij(obs)})\beta_{W,YX} \\ &= \rho_B\tau_X\tau_Y + (1/n)\sigma_X^2\rho_W\sigma_X\sigma_Y \end{aligned} \tag{A.61}$$

$$E(\bar{X}_{\bullet j}Y_{ij(obs)}) = (1/n)E(X_{ij}Y_{ij}) = \rho_B\tau_X\tau_Y + (1/n)\rho_W\sigma_X\sigma_Y \tag{A.62}$$

$$\begin{aligned} E(\bar{X}_{\bullet j}\bar{Y}_{\bullet j}^*) &= E(X_{B,j}Y_{B,j}) + (1/n)E(X_{W,ij(obs)}X_{W,ij(obs)})\beta_{W,YX} \\ &= \rho_B\tau_X\tau_Y + (1/n)\rho_W\sigma_X\sigma_Y. \end{aligned} \tag{A.63}$$

Combining the Equations A.60 to A.63, the expectation for the deviations of the observed values within groups is given by

$$E[(X_{ij(obs)} - \bar{X}_{\bullet j})(Y_{ij(obs)} - \bar{Y}_{\bullet j})] = (1 - 1/n)\rho_W\sigma_X\sigma_Y. \tag{A.64}$$

In a similar manner, this relation can be shown to hold for cross-product terms involving imputed values. It follows that

C_{within} is an unbiased estimator of the within-group covariance $\sigma_{W,XY}$:

$$\begin{aligned} E(\hat{\sigma}_{W,XY}) &= \frac{1}{K(n-1)} \sum_{j=1}^K \sum_{i=1}^n E[(X_{ij} - \bar{X}_{\bullet j})(Y_{ij}^* - Y_{\bullet j}^*)] \\ &= \rho_W\sigma_X\sigma_Y = \sigma_{W,XY}. \end{aligned} \tag{A.65}$$

For the between-groups covariance, we first show (using Equations A.39 and A.63) that the expectation of $C_{between}$ is given by

$$E(C_{between}) = E(\bar{X}_{\bullet j}Y_{\bullet j}^*) = \rho_B\tau_X\tau_Y + (1/n)\rho_W\sigma_X\sigma_Y. \tag{A.66}$$

Then it follows that the estimator of the between-groups covariance is unbiased

$$E(\hat{\sigma}_{B,XY}) = E(C_{between}) - (1/n)E(C_{within}) = \rho_B\tau_X\tau_Y = \sigma_{B,XY}. \tag{A.67}$$

However, the estimator of the between-groups coefficient $\beta_{B,XY}$ is biased because the estimator of the between-groups variance τ_Y^2 is biased

$$\begin{aligned} Bias(\hat{\beta}_{B,XY}) &= \frac{\sigma_{B,XY} + Bias(\hat{\sigma}_{B,XY})}{\tau_Y^2 + Bias(\hat{\tau}_Y^2)} - \beta_{B,XY} \\ &= -\beta_{B,XY} \frac{Bias(\hat{\tau}_Y^2)}{\tau_Y^2 + Bias(\hat{\tau}_Y^2)}. \end{aligned} \tag{A.68}$$

Received February 12, 2015

Revision received May 23, 2016

Accepted May 26, 2016 ■