# Adapting a Correlated Random Effects Method to Handle Ignorable Missingness in the Predictors and Non-ignorable Missingness in the Response

**Hanadi Alzahrani[1,4], Benn Macdonald[1], Caroline Haig[2], John Cleland[3]**
[1] School of Mathematics and Statistics, University of Glasgow
[2] Robertson Centre for Biostatistics, School of Health and Wellbeing, University of Glasgow
[3] School of Cardiovascular and Metabolic Health, University of Glasgow
Glasgow G12 8SQ, UK
Hanadi.Alzahrani@glasgow.ac.uk; Benn.Macdonald@glasgow.ac.uk; Caroline.Haig@glasgow.ac.uk;
John.Cleland@glasgow.ac.uk
[4] King Saud bin Abdulaziz University for Health Sciences
Jeddah/ Saudi Arabia

***Abstract -*** Missing data is a common problem that can be present in longitudinal studies and can seriously impact the statistical analysis estimates by producing biased estimates. A correlated random effect (CRE) method for longitudinal models using latent variables based on Gibbs sampling has previously been proposed to deal with missingness in the response and has demonstrated and performed well in scenarios that assume the missingness of the response variable is Missing Not at Random (MNAR). However, it is currently unable to accommodate incomplete data in the analysis model explanatory variables. The proposed Two-Step method addresses this problem of dealing with Missing at Random (MAR) explanatory variables by incorporating an additional step before utilizing the CRE method. The extra step uses the MICE algorithm, a common approach to handling MAR data and producing imputed datasets. The CRE method is then applied to the imputed MICE datasets. Using simulated longitudinal data, the Two-Step method is compared with the CRE method and some baseline models (scenarios where there is no missingness and scenarios where analysis is complete case), for different numbers of repeated measures and missing proportion factors. The Two-Step method performed similarly to when the CRE method is applied to scenarios where there was no missingness in the explanatory variables, outperforming the complete case scenario in terms of out-of-sample predictive performance and how closely the parameter estimations match the parameters that generates the data.

***Keywords***: Longitudinal data; Latent variable; Missing data; MICE algorithm; Correlated random effect.

## 1. Introduction

Missing data is a common problem that can be present in longitudinal studies, where data records are obtained repeatedly over time on the same subjects. For example, in medical sciences research, patients may miss scheduled appointments or drop out of the study for different reasons. Using complete case data analysis may lead to biased parameter estimates [1]. Three main missingness mechanisms are introduced by [2]; Missing Completely at Random (MCAR), where missing values are independent of observed and unobserved data, Missing at Random (MAR), where missing values depend on observed data, and Missing Not at Random (MNAR), where missing values depend on unobserved data. The MCAR and MAR mechanisms are defined as ignorable missingness because they assume that the observed data explain missingness, therefore, inference is carried out using the observed data only. On the other hand, MNAR is defined as non-ignorable missingness, meaning the missing data mechanism should be also explicitly considered [3]. This involves modelling the joint distribution of the data and the missingness model, simultaneously estimating the observed model and the nonresponse process. The presence of missing values can be in the response and in the explanatory variables in the analysis model, since the measurements of these variables are taken repeatedly over

time in longitudinal studies. The proposed method aims to deal with missing data in both the response and explanatory variables and estimate parameters of interest, where the response missingness is non-ignorable and the explanatory variables missingness is ignorable.

We extend the work that falls in the class of Correlated Random Effect (CRE) selection modelling to handle longitudinal data outlined by [4], where the limitation of the existing method is that it focuses on the missingness in the response and assumes the analysis model explanatory variables are fully observed. This model-based estimation and imputation procedure extends to accommodate incomplete predictors using the Multiple Imputation by Chained Equations (MICE) algorithm to impute missing values under the MAR assumption [5] based on its overall effectiveness and accessibility [6]. The proposed method aims to produce estimates of the analysis regression model where response and explanatory variables are incomplete, reflecting what we might encounter in real-life datasets.

## 2. Proposed Method

The proposed Two-Step method starts by using the MICE algorithm to produce multiple imputed datasets of the missing observations in the analysis model's explanatory variables. Next, the CRE method is applied to each imputed data set to handle missingness in the model response, and the analysis model is estimated simultaneously as a second step. Then, the overall parameter estimates are obtained by combining the posterior distributions. This is done by aggregating the posterior distributions from each dataset into a single distribution as outlined by [7]. In the CRE method we are following the approach proposed by [4]. However, in our work we are assuming our data follows a parametric Linear Mixed Effects model (LMM).

### 2.1. MICE Algorithm

The MICE algorithm begins by choosing a random sample of the incomplete variable's observed values and setting up the incomplete variable imputation model. Every iteration of the procedure involves sampling the incomplete variable model's parameters from its conditional distribution using the most recent completed data and the observed portion of the current variable. After that, given the other variables and parameters, missing values are imputed from the predicted distribution of the missing values. Lastly, it fills in the incomplete variable with the imputed values from the last iteration. It creates multiple imputed datasets by running them multiple times with varied initial values. The MICE algorithm works as a Bayesian simulation technique that uses a Gibbs sampler, which takes samples from the conditional distributions to obtain samples from the joint distribution. The conditional distributions in MICE represent the distributions of incomplete variables, given the observed data variables. The algorithm creates multiple copies $(K)$ of the data and replaces the missing values in each copy with predicted values from observed data. Then, a standard statistical method for each imputed dataset is applied. Finally, the pooled estimates are computed to get general results and to consider the uncertainty produced by the missing values [5]. The process of the MICE algorithm for a dataset consists of vectors of variables $X$ as outlined in Algorithm 1, based on [5]:

---

**Algorithm 1** MICE Algorithm

Define an imputation model $p(X_j^{mis}|X_j^{obs}, X_{-j}, X_{j`})$.
  Fill in initial values of $X_j^0$ by randomly drawing from $X_j^{obs}$.
     for $w = 1, ..., W$ iterations do
       for $j = 1, ..., J$ incomplete variables do
          Define $X_{-j}^w = (X_1^w, ..., X_{j-1}^w, X_{j+1}^{w-1}, ..., X_J^{w-1})$ as the latest imputed data without $X_j$.
          Draw $\delta_j^w \sim p(\delta_j^w|X_j^{obs}, X_{-j}^w, X_{j`})$.
          Draw imputations $X_j^w \sim p(X_j^{mis}|X_j^{obs}, X_{-j}^w, X_{j`}, \delta_j^w)$.
       end $j$
     end $w$

---

where $X_j$ is the $j$ th incomplete variable, $j = 1, ..., J$ , and $X_{j`}$ is the $j`$th complete variable, $j` = 1, ..., J`$ and $X_{-j}$ is all other incomplete variables except $X_j$. $X_j^{mis}$ and $X_j^{obs}$ represents the missing and observed observations in the $j^{th}$

variable, respectively. $\boldsymbol{\delta}_j$ represents the vector of the imputation model parameters for variable $j$. The MICE algorithm will run $K$ times in parallel and the number of iterations, $W$, recommended by [5] is between 5 and 10.

## 2.2. CRE Method

The CRE method proposed by [4] will be explained in this section, although we employ a parametric LMM, which is a standard framework for studying the relationship between longitudinal outcomes and predictor variables. For a continuous response measured over $m$ different time points from $n$ subjects and a set of predictors, some of which are partially observed, the response for the $i^{th}$ subject at the $t^{th}$ time point, which we denote by $Y_i(t)$, can be modelled as the following:

$$Y_i(t) = \mu + \sum_{j=1}^{J} \beta_j X_{ji}(t) + \sum_{j`=1}^{J`} \lambda_{j`} X_{j`i}(t) + u_i \tilde{Z}_i(t) + e_i(t), \tag{1}$$

where $J`$ and $J$ express the number of predictors of fixed effects that are fully observed and partially observed respectively. The fixed intercept represents the mean of the overall population, expressed as $\mu$, $\beta_j$ & $\lambda_{j`}$ denote the regression coefficients associated with the $j^{th}$ partially observed and $j`^{th}$ fully observed fixed effects, respectively, $X_{ji}(t)$ is the outcome of the $j^{th}$ partially observed fixed effect for subject $i$ at time $t$ and $X_{j`i}(t)$ is the outcome of the $j`^{th}$ fully observed fixed effect for subject $i$ at time $t$. Subject-specific random effects $u_i$ capture the longitudinal dependence and are assumed to be independent and identically distributed from $N(0, \sigma_B^2)$ and $\tilde{Z}_i(t)$ is the outcome of the random effect for subject $i$ at time $t$. The residuals, $e_i(t)$, are assumed to be independent and identically distributed from $N(0, \sigma_A^2)$. To consider the non-ignorable missing values in the response we will define a binary missing response indicator $U_i(t)$, where $U_i(t) = 0$ if $Y_i(t)$ is missing and $U_i(t) = 1$ if $Y_i(t)$ is observed. The latent response variable can be written as:

$$Y_i(t) = \begin{cases} Y_i^*(t), & if \quad U_i(t) = 1 \\ missing, & if \quad U_i(t) = 0 \end{cases} \tag{2}$$

We can then rewrite the regression model given in Equation (1) as follows:

$$Y_i^*(t) = \mu + \sum_{j=1}^{J} \beta_j X_{ji}(t) + \sum_{j`=1}^{J`} \lambda_{j`} X_{j`i}(t) + u_i \tilde{Z}_i(t) + e_i(t) \tag{3}$$

Now consider the following model for the missing response mechanism as:

$$U_i^*(t) = \alpha + \sum_{l=1}^{L} \theta_l X_{li}(t) + v_i \tilde{Z}_i(t) + \varepsilon_i(t), \tag{4}$$

where $L = J` + J$ is the total number of fixed effects in the response model, the fixed intercept represents the mean of the overall population, expressed as $\alpha$, $\theta_l$ denotes the regression coefficients of the $l^{th}$ fixed effects, which expresses the systematic influence of missingness due to the unobserved response variables. Subject-specific random effects $v_i$ capture the longitudinal dependence and are assumed to be independent and identically distributed from $N(0, \sigma_C^2)$ and the residuals $\varepsilon_i(t)$ are assumed to be independent and identically distributed from $N(0, 1)$. $U_i(t)$ is conditional on the proclivity $U_i^*(t)$ via a probit model, which divides the standard normal into two parts. If $U_i^*(t)$ is greater than zero, then $U_i(t) = 1$ and if $U_i^*(t)$ is less than zero then $U_i(t) = 0$ [8]. Following [4], we assume a correlation between the response variable $Y_i^*(t)$ and the response missing indicator variable $U_i^*(t)$ random effects, and therefore consider $u_i$ and

$v_i$ are correlated random vectors following a multivariate normal distribution with mean vector $\mathbf{0}$ and covariance matrix $\Sigma = \begin{pmatrix} \sigma_B^2 & \sigma_D^2 \\ \sigma_D^2 & \sigma_C^2 \end{pmatrix}$, where $\sigma_D^2$ represents the covariance between $u_i$ and $v_i$ random effects.

Bayesian inference is used to estimate the model parameters in Equation (3) and in Equation (4) using an iterative MCMC algorithm (Gibbs sampling) to simultaneously impute missing values in the response and produce analysis model estimates. For Gibbs sampling to be carried out, one needs to sample from the joint posterior of the model parameters and latent variables. Let $\mathbf{Y} = \big(Y_{11}(t), \dots, Y_{nm}(t)\big)$, $\mathbf{Y}^* = \big(Y^*_{11}(t), \dots, Y^*_{nm}(t)\big)$, $\mathbf{U} = \big(U_{11}(t), \dots, U_{nm}(t)\big)$, $\mathbf{U}^* = \big(U^*_{11}(t), \dots, U^*_{nm}(t)\big)$. The joint posterior density for the latent variables and the parameters associated with the proposed model is:

$$p\big(\Theta_{Y,U}, \mathbf{Y}^*, \mathbf{U}^* \mid \mathbf{Y}, \mathbf{U}\big) \propto p\big(\Theta_{Y,U}\big) \times$$
$$\prod_{i=1}^{n} \int_{u_i=-\infty}^{u_i=\infty} \int_{v_i=-\infty}^{v_i=\infty} \prod_{t=1}^{m} \times f(Y_i^*(t), U_i^*(t) \mid u_i, v_i) \times g(u_i, v_i) \times \tag{5}$$
$$\{I(U_i^*(t) > 0)I(U_i(t) = 1) + I(U_i^*(t) \le 0)I(U_i(t) = 0)\} du_i dv_i,$$

where $\Theta_{Y,U} = \{ \mu, \boldsymbol{\beta}, \boldsymbol{\lambda}, \alpha, \boldsymbol{\theta}, \sigma_A^2, \Sigma \}$ denote a set of analysis model parameters, $\boldsymbol{\beta}, \boldsymbol{\lambda}, \boldsymbol{\theta}$ each denotes a vector of the corresponding regression coefficients. $f(Y_i^*(t), U_i^*(t))$ is the joint distribution of $Y_i^*(t)$ and $U_i^*(t)$, $g(u_i, v_i)$ is the joint distribution of the random effects $u_i$ and $v_i$. I(A) is an indicator variable which takes the value 1 if A occurs and zero otherwise. The prior distribution is $p\big(\Theta_{Y,U}\big)$ for $\Theta_{Y,U}$. We consider non-informative conjugate priors due to lack of prior information and to produce closed form conditional distributions. Therefore, the prior distribution can be broken up and expressed as:

$$p(\mu, \boldsymbol{\beta}, \boldsymbol{\lambda}, \sigma_A^2) \propto \frac{1}{\sigma_A^2}; \qquad p(\alpha, \boldsymbol{\theta}) \propto 1; \qquad p(\Sigma) \propto IW(v, \Lambda) \tag{6}$$

Note, we use the Inverse Wishart (IW) distribution because in our application using a more vague prior for $\Sigma$ has shown non convergence. We set $v > p + 1$ and $\Lambda = (v - p - 1)I$, to make it weakly informative and a valid prior distribution, where $p$ is the dimension of the covariance matrix, $\Lambda$ is a $p \times p$ scale matrix and $I$ is a $p \times p$ identity matrix [9].

## 2.3. Estimation

The response and the observed predictor variables in the analytical model will be used in the Two-Step method to impute the incomplete predictor variables. Since the imputation and analysis steps are performed separately [6], this identifies and captures any relationships present in the data, potentially improving the imputation accuracy. Ten imputed MICE datasets will be created (i.e. K = 10 in Algorithm 1), which could increase the accuracy of the results but may be computationally costly. The LMM is applied in the MICE step to impute missing values for incomplete time-varying predictors given all other variables. The Two-Step method will produce posterior distributions for each parameter estimate. Therefore, these posteriors will be aggregated by mixing them into a single posterior distribution [7].

## 3. Simulated Data

The performance of the proposed method is examined using a simulated study; we generate data for longitudinal settings, using the LMM with a random intercept. We assume one incomplete continuous predictor $X_1(t)$ measured over time, and two fully observed predictors $X_2(t)$ and $X_3(t)$. The analysis model and missingness model of the response were generated based on:

$$Y_i^*(t) = \beta_0 + \beta_1 X_{1i}(t) + \beta_2 X_{2i}(t) + \beta_3 X_{3i}(t) + u_i \tilde{Z}_i(t) + e_i(t) \tag{7}$$

$$U_i^*(t) = \theta_0 + \theta_1 X_{1i}(t) + \theta_2 X_{2i}(t) + \theta_3 X_{3i}(t) + u_i \tilde{Z}_i(t) + \epsilon_i(t) \tag{8}$$

The values of the regression coefficients in the response missingness model in Equation 8 were derived to produce the desired missing data proportion, using a probit regression equation [10] to connect missingness probabilities of the response $\mathbf{Y}$ to values of $\mathbf{Y}$ through the latent missingness indicator regression model $\mathbf{U}^*$. The model predictor $X_2(t)$ is generated from $Unif(0,2)$, $X_3(t)$ from $Bern(0.6)$ and $X_1(t)$ from a LMM with a random intercept and the fully observed predictors $X_2(t)$ and $X_3(t)$ as explanatory variables. This allows us to reproduce the missing values of $X_1(t)$ using these variables ($X_2(t)$ and $X_3(t)$). We assume $n = 100$ number of subjects in the study and varied the values of the number of repeated measures per subject, assuming m = 2, 4 & 8. A number of combinations for the proportion of missing data was considered, namely 20%, 40% and 60% in the analysis model response with fixed 20% missingness in the incomplete predictor and 20%, 40% and 60% in the analysis model incomplete predictor with fixed 20% missingness in the analysis model response. The missingness of the incomplete predictor $X_1(t)$ is generated using the "deleteMARcensoring()" function in the "missMethod" package [11] in R. The performance of the proposed method is compared with the baseline models i.e. the model with fully observed variables (no missing data) and the model with available data (missing values remain in the dataset and no imputation is applied, in other words, complete case analysis is carried out). Also, we will compare the proposed method with the CRE method [4] assuming fully observed predictors, as well as when MICE is used to impute the model's response and incomplete predictor. Root Mean Square Error (RMSE) between data-generating parameters and estimated values will be used to assess the methods, as well as the RMSE for out-of-sample prediction. In order to perform out-of-sample prediction, a test dataset is generated using identical simulation settings as already outlined, but without missingness, for each number of repeated measures. We examine the distribution of these criteria across 100 replications.

## 4. Results

The MCMC simulations were performed for 50,000 iterations, with a thinning rate of 10 applied and half of the iterations designated as a burn-in phase. A single chain was produced due to the computational time and storage of the Two-Step method. To assess convergence, we examined the Geweke convergence statistic [12] for individual parameters and visually examine the trace plots for each parameter. The baseline methods were fitted using a Hierarchal Bayesian model using Hamiltonian Monte Carlo utilised by the "brm" function in R [13].

Fig.1 shows the results for 4 repeated measures. The results for 2 and 8 are omitted for brevity since the conclusion is consistent across the repeated measures. The degree of variation between the estimated and data-generating parameters is shown by the RMSE analysis, which offers insights into the accuracy of our estimates. Generally, the Two-Step method RMSE is comparable to the CRE method which has no missingness in the model predictor, except when there is a larger proportion of missingness (60%) in the incomplete predictor, in which case the Two-Step method tends to have a larger RMSE. The MICE imputation of the response and incomplete predictor has similar RMSE performance to the Two-Step method, except with higher proportion of missingness in the analysis model response and incomplete predictor (60%), where it tends to have larger RMSE values and uncertainty. Overall, the available data shows high RMSE values and uncertainty compared with the other methods.
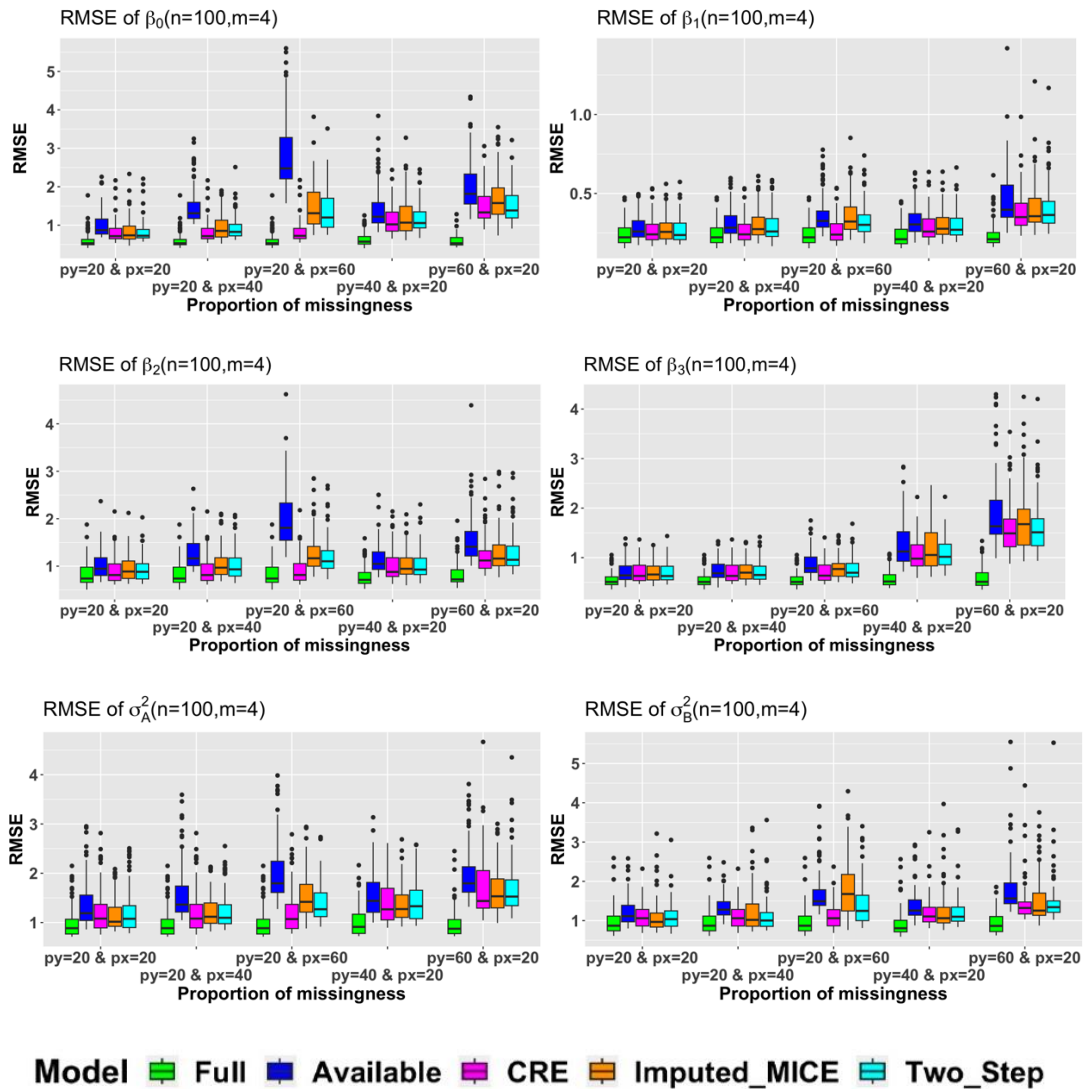
Fig. 1: Boxplots illustrating the RMSE of the analysis model parameters. Each boxplot represents one of the applied methods. The y-axis shows the RMSE values, and the x-axis represents one of the combinations of proportions of missingness in the model response and incomplete predictor. It's obvious that the available data method has larger RMSE values as compared to other methods.

In terms of out-of-sample performance, by looking at Fig. 2, we can see the method is robust enough to describe data it has not been built upon. Again, the results shown are for 4 repeated measures, with 2 and 8 repeated measures displaying the same conclusion (and are thus omitted for brevity). The results showed that the available data method has the lowest performance, while the MICE imputation approach for the model response and incomplete predictor has the second lowest out-of-sample prediction with higher RMSE values and lower RMSE density concentrated near smaller RMSE values. In contrast, the Two-Step and CRE methods have comparable and better overall performance and are similar to the full data, with 20% missingness in the model response. Not surprisingly, the full data method performs the best because it contains no missing values, providing a nice benchmark for the "best-case scenario". This indicates that the Two-Step method is performing well overall, as the performance is overall on par with the full data method.
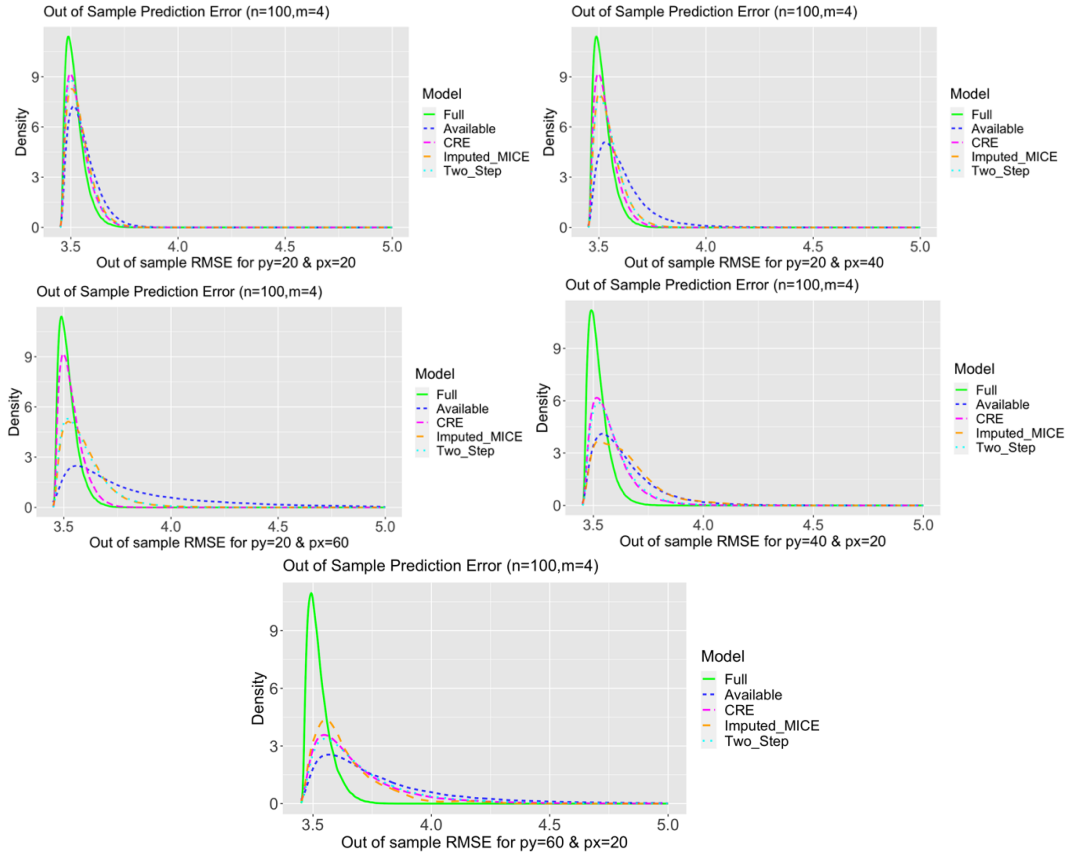
Fig. 2: The density plots of the out-of-sample RMSE for different methods across various missingness proportions in the response and the incomplete predictor in case of four repeated measures. Each density curve corresponds to one of the methods used, and each plot corresponds to a different combination of missingness proportion in the response and the incomplete predictor. The available data method has a lower density for smaller RMSE values and overall has a density shifted to the right of the plots (indicating higher error).

## 5. Conclusion

In longitudinal data analysis, it is common to have incomplete data for different reasons. Usually, ignorable missingness is assumed, however, it is possible to have non-ignorable missingness. Building on one of the recently proposed non-ignorable modelling frameworks [4], we proposed a method to accommodate ignorable missingness in the model predictor in addition to the missingness in the model response. The performance of the proposed Two-Step method was evaluated using RMSE for both the parameter estimation and out-of-sample prediction. The Two-Step method outperformed the available data, which uses only the observed values (complete case analysis) and the MICE imputation approach of both the response and incomplete predictor. It has a similar performance to the CRE method with fully observed explanatory variables, but with the advantage that it can handle missingness in the incomplete predictor.

Since the Two-Step method involves applying the CRE method to multiple imputed datasets, it requires $K$ times the computational time of the CRE method, which makes it more computationally costly. Depending on the size of $K$, the dataset sample size, and whether parallel computing is available, the computational time and storage space may be prohibitive. Scalability of the method can therefore be the focus of future research. Additionally, it is worth noting that for the current work, the MAR assumption for missingness in the explanatory variable was satisfied (due to the way missingness was simulated). Future work could assess how sensitive the approach is to the type of missingness in the explanatory variable, since the MICE algorithm may perform more poorly when the data is MNAR, thus potentially impacting the performance of the proposed Two-Step method.

## Acknowledgements

## References

[1]   K. J. M. Janssen, A. R. T. Donders, F. E. Harrell, Y. Vergouwe, Q. Chen, D. E. Grobbee, and K. G. M. Moons, "Missing covariate data in medical research: to impute is better than to ignore", *Journal of Clinical Epidemiology*, vol. 63, no. 7, pp. 721-727, 2010.

[2]   D. B. Rubin, "Inference and missing data", *Biometrika*, vol. 63, no. 3, pp. 581–592, 1976.

[3]   Z. Ma and G. Chen, "Bayesian methods for dealing with missing data problems", *Journal of the Korean Statistical Society*, vol. 47, no. 3, pp. 297–313, 2018.

[4]   P. Bhuyan, "Estimation of random-effects model for longitudinal data with nonignorable missingness using Gibbs sampling", *Computational Statistics*, vol. 34, no. 4, pp. 1693–1710, 2019.

[5]   S. van Buuren, *Flexible Imputation of Missing Data*. CRC press, 2018.

[6]   N. S. Erler, D. Rizopoulos, J. van Rosmalen, V. W. V. Jaddoe, O. H. Franco, and E. M. E. H. Lesaffre, "Dealing with missing covariates in epidemiologic studies: a comparison between multiple imputation and a full Bayesian approach", *Statistics in Medicine,* vol. 35, no. 17, pp. 2955–2974, 2016.

[7]   A. Gelman, J. B. Carlin, H. S. Stern, D. B. Dunson, A. Vehtari, and D. B. Rubin, *Bayesian data analysis,* CRC press, 2013.

[8]   V. E. Johnson, and J. H. Albert, *Ordinal data modelling*, Springer Science and Business Media, 2006.

[9]   N. K. Schuurman, R. P. Grasman, and E. L. Hamaker, "A Comparison of Inverse-Wishart Prior Specifications for Covariance Matrices in Multilevel Autoregressive Models", *Multivariate Behavioural Research*, vol. 51, no. 2-3, pp. 185-206, 2016.

[10] J. H. Albert and S. Chib, "Bayesian analysis of binary and polychotomous response data", *Journal of the American statistical Association.* vol. 88, no. 422, pp. 669-679, 1993.

[11] T. Rockel, "missMethods: Methods for missing data", *R Package Version 0.2.2,* 2020.

[12] J. Geweke, "Evaluating the accuracy of sampling-based approaches to the calculations of posterior moments", *Bayesian Statistics,* no. 4, pp. 641–649, 1992.

[13] *P.-C.* Bürkner, "brms: an R package for Bayesian multilevel models using Stan"*, Journal of Statistical Software.*, vol. 80, no. 1, 2017.