

# Master Research Report: Multilevel Multivariate Imputation by Chained Equations through Bayesian Additive Regression Trees

*Methodology and Statistics for the Behavioural, Biomedical and Social  
Sciences*

*Heleen Brügger*

**Word count:**

**Candidate Journal:**

**FETC Case Number:**

**Supervisors:**

MSc. T. Volker

Dr. G. Vink

MSc. H. Oberman

2814

Computational Statistics & Data Analysis

23-1778

---

Utrecht University

Utrecht University

Utrecht University

# 1 Introduction

Incomplete data is a common challenge in many fields of research. Frequently used ad hoc strategies to deal with missing data, as listwise or pairwise deletion, mean imputation or regression imputations often lead to erroneous inferences in realistic situations, due to biased estimates and inaccurate variance estimates [Austin et al., 2021, Enders, 2017, Kang, 2013, van Buuren, 2018]. Rubin defined three types of missing data mechanisms: Missing Completely At Random (MCAR) where the cause of the missing data is unrelated to the unobserved or observed data, Missing At Random (MAR) where the missing data is related to the observed data, and Missing Not At Random (MNAR) where the missing data is related to unobserved data [Rubin, 1976].

Multiple imputation (MI) is considered a valid method for dealing with incomplete data and allows us to separate the missing data problem from the analysis problem [Audigier et al., 2018, Austin et al., 2021, Burgette and Reiter, 2010, Enders, 2017, Grund et al., 2021, Hughes et al., 2014, Mistler and Enders, 2017, Van Buuren, 2007, van Buuren, 2018]. MI is used to impute each missing value in the dataset more than once given the observed data, considering necessary variation associated with the missingness problem. The multiply imputed datasets are analyzed, and the corresponding inferences are pooled according to Rubin’s rules [Austin et al., 2021, Carpenter and Kenward, 2013, Rubin, 1987, van Buuren, 2018]. However, specifying the imputation models, the models used to impute the missing data, can be challenging. The concept of congeniality dictates that the imputation models should be at least as general as the analysis model and preferably all-encompassing [Bartlett et al., 2015, Enders et al., 2018a, Grund et al., 2016, 2018b, Meng, 1994]. Otherwise, it will not capture every aspect of the data and the completed datasets cannot be properly analysed. So, when the complexity of data increases, specifying the imputation models becomes more difficult [Grund et al., 2018b, van Buuren, 2018].

Congeniality-issues become more pronounced when MI is used in a multilevel data context [Audigier et al., 2018, Dong and Mitani, 2023, Enders et al., 2020, 2018a,b, 2016, Grund et al., 2016, 2018a,b, 2021, Lüdtke et al., 2017, Mistler and Enders, 2017, Quartagno and Carpenter, 2022, Resche-Rigon and White, 2018, Taljaard et al., 2008, van Buuren, 2018]. Multilevel data is hierarchically structured, where, for example, students are nested within schools [Hox and Roberts, 2011, Hox et al., 2017]. When analysing multilevel data, the hierarchical structure should be taken into account, which can be done using multilevel models (MLM) [Hox and Roberts, 2011, Hox et al., 2017, Lüdtke et al., 2017]. They can contain both level-1, and level-2 variables, relating to the individual and class respectively, random intercepts, random slopes, and cross-level interactions [Hox and Roberts, 2011, Hox et al., 2017]. The complexity of the multilevel analysis model is built step-wise with non-linearities [Hox and Roberts, 2011, Hox et al., 2017]. Thus, when defining imputation models the hierarchical structure should be included, along with the complicated non-linearities from cross-level interactions. Since ignoring the multilevel structure will underestimate the intra-class correlation (ICC) [Hox and Roberts, 2011, Lüdtke et al., 2017, Taljaard et al., 2008, van Buuren, 2018], which can be interpreted as the expected correlation between two randomly sampled individuals from the same group or the proportion of the total variance at level-2 [Gulliford et al., 2005, Hox and Roberts, 2011, Shieh, 2012]. So, defining these models is quite challenging [Burgette and Reiter, 2010, Hox and Roberts, 2011, van Buuren, 2018].

One of the two MI frameworks, fully conditional specification (FCS) or otherwise known as chained equations, is believed to be more flexible than its counterpart, joint modeling (JM) [Audigier et al., 2018, Burgette and Reiter, 2010, Grund et al., 2018a, Van Buuren, 2007]. FCS iteratively imputes each incomplete variable conditional on complete and previously imputed variables [Enders et al., 2018a,b, 2016, Grund et al., 2018a, Hughes et al., 2014, Mistler and Enders, 2017, van Buuren, 2018]. In a multilevel context, FCS employs univariate linear mixed models to account for the hierarchical structure [Enders et al., 2018a, Mistler and Enders, 2017, Resche-Rigon and White, 2018]. Furthermore, FCS can be used to impute non-linearities, such as cross-level interactions, by using ‘passive imputation’ or defining a separate imputation model for the non-linearities [Grund et al., 2018b, van Buuren, 2018]. On the other hand, JM doesn’t account for cross-level interactions at all [Grund et al., 2018b, van Buuren, 2018]. Still, imputation models including cross-level interaction or non-linear terms in FCS is still very complicated [Grund et al., 2018b, 2021] and, thus, researchers’ focus has predominantly been on the inclusion of random intercepts and slopes, but not of cross-level interactions [Enders et al., 2020, 2018a,b, 2016, Grund et al., 2016, 2018a].

Using non-parametric tree-based models might solve this problem because they do not assume a specific data distribution and, thus, implicitly model non-linear relationships and interactions between the predictor variables, and handle continuous and categorical variables simultaneously [Breiman et al., 1984, Burgette and Reiter, 2010, Chipman et al., 2010, Hill et al., 2020, James et al., 2021, Lin and Luo,

2019, Salditt et al., 2023]. In a single-level imputation context, the use of tree-based, non-parametric models like regression trees, random forests, or Bayesian Additive Regression Trees (BART) simplified imputation models and performed better than parametric methods: the imputations showed better confidence interval coverage of the population parameters, lower variance and lower bias, especially in non-linear and interactive contexts [Burgette and Reiter, 2010, Silva and Gutman, 2022, Xu et al., 2016]. Waljee et al. [2013] also found lower imputation error when imputing with a random forest algorithm compared to MICE, KNN and mean imputation.

BART models have been implemented in a multilevel prediction context. However, multilevel-BART models (M-BART) have predominantly been implemented with only random intercepts [Chen, 2020, Tan et al., 2016, Wagner et al., 2020, Wundervald et al., 2022]. In a prediction context, Wagner et al. [2020] have found that this random intercept M-BART model provided better predictions with a lower mean squared error (MSE) compared to a parametric MLM, Tan et al. [2016] found higher area under the curve (AUC) values compared to a single level BART model and linear random intercept model, and Chen [2020] found better predictions and better coverage of the estimates compared to parametric models and a single-level BART model. Other researchers modeled the random intercept as an extra split on each terminal node and found a lower MSE compared to a standard BART model and parametric MLMs [Wundervald et al., 2022]. Dorie et al. [2022] developed a multilevel BART model that included random intercepts and random slopes by combining BART with the Stan algorithm. Where the random parts are modeled by Stan [Lee et al., 2017]. Their results showed that their algorithm ‘stan4bart’ showed better coverage of the population values and lower root mean squared error (RMSE) compared to BART models with varying intercept, BART models ignoring the multilevel structure, bayesian causal forests (BCF), and parametric MLMs.

In spite of these promising findings, M-BART models have yet to be implemented in a multilevel multiple imputation context. Thus, my research question will be: *Can multivariate imputation by chained equations through a multilevel bayesian additive regression trees model improve the bias, variance and coverage of the estimates in a multilevel context compared to current practices?* Given the success of non-parametric models in single-level MI, I anticipate that employing M-BART models in a multilevel missing data context will reduce bias, accurately model variance, and improve estimate coverage compared to classical multilevel imputation through *2l.pmm*, *2l.lmer*, *2l.pan*, *2l.jomo*, *rf* and single-level *pmm* and complete case analysis in the R-package MICE [Buuren and Groothuis-Oudshoorn, 2011].

The research report’s sections will cover theoretical background, methods for evaluating M-BART models, preliminary results, and discussion of next steps.

## 2 Method

### 2.1 Theoretical background

#### 2.1.1 Bayesian Additive Regression Trees (BART)

BART is a sum-of-trees model proposed by Chipman et al. [2010]. Regression trees are its building blocks [Chipman et al., 2010, Hill et al., 2020, James et al., 2021]. Regression trees divides the data into subgroups by recursively splitting the data into binary subgroups based on the predictors maximizing the homogeneity of the subgroups [Hastie, 2017, James et al., 2021, Salditt et al., 2023]. Recursive binary partitioning of the predictor space doesn’t assume a specific data form, making this a non-parametric model [Hastie, 2017, James et al., 2021, Salditt et al., 2023] and allows regression trees to model non-linearities well and automatically [Burgette and Reiter, 2010, Hill et al., 2020].

Chipman et al. [2010] define the BART model as:

$$f(\mathbf{x}) = \sum_{k=1}^m g(\mathbf{x}; T_k, M_k), \quad (1.1)$$

$$f(\mathbf{x}) = g(\mathbf{x}, T_1, M_1) + g(\mathbf{x}_i, T_2, M_2) + \dots + g(\mathbf{x}_i, T_m, M_m), \quad (1.2)$$

where  $f(\mathbf{x})$  is the overall fit of the model: the sum of many regression trees,  $\mathbf{x}$  are the predictor variables,  $T_k$  is the  $k^{\text{th}}$  tree and  $M_k$  is the collection of leaf parameters within the  $k^{\text{th}}$  tree [Chipman et al., 2010, Hill et al., 2020, James et al., 2021]. The data are assumed to arise from a model with additive normally distributed errors:  $Y = \sum_{k=1}^m g(\mathbf{x}; T_k, M_k) + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2)$ . Next to the sum-of-trees model, BART also includes a regularization prior that constrains the size and fit of each tree so that

each contributes only a small part to prevent overfitting [Chipman et al., 2010, Hill et al., 2020, James et al., 2021]. BARTs are estimated using the Bayesian back-fitting Markov Chain Monte Carlo (MCMC) algorithm. It updates individual trees, considering the remaining trees, their associated parameters, and the residual standard deviation ( $\sigma$ ). It fits a new tree to the partial residuals,  $r_i$ , treating them as the data, by perturbing the tree from the previous iteration. It can either *grow*, *prune*, or *change* a tree. *Growing* means adding additional splits, *pruning* removes splits, and *changing* changes decision rules. The algorithm stops after the specified amount of iterations. The partial residuals are defined as:

$$r_i = y_i - \sum_{k' < k} \hat{f}_{k'}^b(x_i) - \sum_{k' > k} \hat{f}_{k'}^{b-1}(x_i), \quad (2)$$

where  $\hat{f}_k^b(x_i)$  is the prediction of the  $k^{\text{th}}$  tree in the  $b^{\text{th}}$  iteration for person  $i$ .

### 2.1.2 Multilevel-BART (M-BART)

Chen [2020], Wagner et al. [2020] and Tan et al. [2016] define a multilevel-BART (M-BART) model including a random intercept building on the work of Lin and Luo [2019]. The M-BART algorithm breaks down the observed variable into fixed and random components. The fixed components are modeled by BART and the random components are modeled by a linear mixed effects model [Chen, 2020, Tan et al., 2016, Wagner et al., 2020]. The M-BART model including a random intercept can be identified as:

$$f(\mathbf{x}) = \sum_{k=1}^m g(\mathbf{x}; T_k, M_k) + \alpha_j, \quad (3)$$

where, now,  $f(\mathbf{x})$  is the overall fit of the model incorporating random intercept  $\alpha_j$  for cluster  $j$ .

## 2.2 Simulation study

For this research report, I will conduct a simulation study to examine the performance of three M-BART models in a multilevel prediction context compared to a single-level BART model.

### 2.2.1 Data generating mechanism

The population data-generating mechanism will be based on the following MLM:

$$y_{ij} = \beta_{0j} + \beta_{1j}X_{1ij} + \beta_{2j}X_{2ij} + \beta_{3j}X_{3ij} + \beta_{4j}X_{4ij} + \beta_{5j}X_{5ij} + \beta_{6j}X_{6ij} + \beta_{7j}X_{7ij} + \epsilon_{ij}, \quad (4)$$

$$\beta_{0j} = \gamma_{00} + \gamma_{01}Z_{1j} + v_{0j}, \quad (4.1)$$

$$\beta_{1j} = \gamma_{10} + \gamma_{11}Z_{1j} + v_{1j}, \quad (4.2)$$

$$\beta_{2j} = \gamma_{20} + \gamma_{21}Z_{1j} + v_{2j}, \quad (4.3)$$

$$\beta_{3j} = \gamma_{30} + \gamma_{32}Z_{2j} + v_{3j}, \quad (4.4)$$

$$\beta_{4j} = \gamma_{40} + v_{4j}, \quad (4.5)$$

$$\beta_{5j} = \gamma_{50} + v_{5j}, \quad (4.6)$$

$$\beta_{6j} = \gamma_{60} + v_{6j}, \quad (4.7)$$

$$\beta_{7j} = \gamma_{70}, \quad (4.8)$$

where  $y_{ij}$  is a continuous level 1 outcome variable for person  $i$  in group  $j$  and  $Z_{1j}$  and  $Z_{2j}$  are continuous level 2 variables. The random intercept  $\beta_{0j}$  is determined by the grand mean  $\gamma_{00}$ , the group effect  $\gamma_{01}Z_{1j}$  and the group-level random residuals  $v_{0j}$ . The regression coefficients  $\beta_{1j}$ ,  $\beta_{2j}$ , and  $\beta_{3j}$  for the continuous variables  $X_{1ij}$ ,  $X_{2ij}$ , and  $X_{3ij}$  depend on the intercepts  $\gamma_{10}$ ,  $\gamma_{20}$ , and  $\gamma_{30}$ , the cross-level interactions  $\gamma_{11}Z_{1j}$ ,  $\gamma_{21}Z_{1j}$ , and  $\gamma_{32}Z_{2j}$ , and the random slopes  $v_{1j}$ ,  $v_{2j}$ , and  $v_{3j}$ . The regression coefficients  $\beta_{4j}$ ,  $\beta_{5j}$  and  $\beta_{6j}$  are determined by the intercepts  $\gamma_{40}$ ,  $\gamma_{50}$  and  $\gamma_{60}$  and the random slopes  $v_{4j}$ ,  $v_{5j}$  and  $v_{6j}$ . The regression coefficient  $\beta_{7j}$  is determined by the intercept  $\gamma_{70}$ . The residuals and random slopes  $v_{0j}$ ,  $v_{1j}$ ,  $v_{2j}$ ,  $v_{3j}$ ,  $v_{4j}$ ,  $v_{5j}$ , and  $v_{6j}$  and  $\epsilon_{ij}$  follow a zero-mean normal distribution. The variance of  $v_{0j}$ , the group-level random residuals, were scaled such that the specified ICC value was obtained.  $v_{1j}$ ,  $v_{2j}$ ,  $v_{3j}$ ,  $v_{4j}$ ,  $v_{5j}$ , and  $v_{6j}$  all have a variance of 1.  $\epsilon_{ij}$  had a variance of 25.  $X_1$ ,  $X_2$ ,  $X_3$ ,  $X_4$ ,  $X_5$ ,  $X_6$  and  $X_7$  are multivariate normally distributed:  $\mathbf{X}_{ij} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , with  $\boldsymbol{\mu} = (0, 0, 0, 0, 0, 0, 0)$

and  $\Sigma = \text{diag}(6.25, 9, 4, 11.56, 4, 2.5, 19.36)$  with no co-variances. The level-2 variables  $Z_1$  and  $Z_2$  are also multivariate normally distributed:  $\mathbf{Z}_j \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ , with  $\boldsymbol{\mu} = (0, 0)$  and  $(\Sigma) = \text{diag}(1, 2.56)$  with no co-variances. The group-level effects ( $\gamma_{01}$  and  $\gamma_{02}$ ) were set to .5, the cross-level interactions ( $\gamma_{11}$ ,  $\gamma_{21}$ , and  $\gamma_{32}$ ) were set to .35, and the overall intercept ( $\gamma_{00}$ ) was set to 10. The within-group effect sizes ( $\gamma_{10}$ ,  $\gamma_{20}$ ,  $\gamma_{30}$ ,  $\gamma_{40}$ ,  $\gamma_{50}$ ,  $\gamma_{60}$ , and  $\gamma_{70}$ ) were varied in the simulations.

### 2.2.2 Simulation design

Table 1 shows the factors that will be varied in the simulation study. All these values are realistic in practice and/or previously proposed [Enders et al., 2020, 2018b, Grund et al., 2018b, Gulliford et al., 1999, Hox et al., 2017, Murray and Blitstein, 2003]. For each combination of varying parameters, 6 datasets will be simulated for each scenario to reduce computational time. 4 different models will be compared: a single level BART, single level BART with group-dummy's, a multilevel BART model incorporating a random intercept [Chen, 2020, Tan et al., 2016, Wagner et al., 2020, Wundervald et al., 2022], and a multilevel BART model combined with Stan to model the random parts of the models [Dorie et al., 2022]. The first three models will be performed with the package *dbarts* [Dorie, 2023a] and the last will performed with the package *stan4bart* [Dorie, 2023b] in R [R Core Team, 2023]. The default arguments from the function *rbart\_vi* will be used for all models, as well as the default priors.

**Table 1:** Simulation design

Parameter	Values
Number of clusters (j)	30, 50
Within-cluster sample size ( $n_j$ )	5, 15, 35, 50
Intraclass Correlation (ICC)	0, .05, .3, .5
Within-group effect size ( $\gamma$ )	.2, .5, .8

### 2.2.3 Evaluation

In this research report, the preliminary results will pertain to the evaluation of different BART models in terms of the relative bias and Mean Squared Error (MSE) of the estimates, which will be calculated as follows [Morris et al., 2019]:

$$Bias = \frac{1}{n_{sim}} \sum_{t=1}^{n_{sim}} (\hat{\theta}_t - \theta), \quad (5)$$

$$MSE = \frac{1}{n_{sim}} \sum_{t=1}^{n_{sim}} (\hat{\theta}_t - \theta)^2, \quad (6)$$

$$(5a)$$

where  $\hat{\theta}_t$  is the estimated parameter in simulation  $t$ ,  $\theta$  is the true value, and  $n_{sim}$  is the number of simulated datasets.

## 3 Results

The results of this research report are shown in figure 1 and 2. Figure 1 shows the average relative bias for all models, simulated datasets, every *ICC* value, and within-group effect size. On the x-axis we can see the different simulated datasets with their names specifying the total sample size, number of groups and group sizes. Figure 2 shows the average Mean Squared Error (MSE) for all models, datasets, *ICC* values, and within-group effect sizes ( $\gamma$ ).

In figure 1 we can see that when there is no multilevel structure in the dataset,  $ICC = 0$ , the models perform similarly in terms of relative bias: overall, the bias is around zero. We can see a slight increase in relative bias when the total sample size is small for all  $\gamma$ . When we increase the *ICC*,  $ICC > 0$ , we can see a divide in the performance of the models. Even when  $ICC = .05$ , the single level BART model (bart) increases in relative bias. bart shows higher levels of relative bias when the total sample size is small. However, when the total sample size is large, 2500, bart performs slightly better than the single level BART model including a group-dummy (gbart) for all  $\gamma$ . The same pattern is observed for the other models: when  $ICC = .05$ , the relative bias decreases when the total sample size increases. This pattern is further continued when increasing the *ICC*: the relative bias decreases when the total sample size increases for both  $ICC = .3$  and  $ICC = .5$ . However, when considering the stan4bart model, which models the random parts of the model in Stan and the fixed parts in BART, the relative bias stays

considerably constant when increasing the  $ICC$ : the relative bias is higher when the total sample size is small, but it does not seem to significantly increase with higher  $ICC$  or higher  $\gamma$ . This cannot be said for the bart and gbart models, which show a slight increase in relative bias when increasing the  $ICC$ , especially in the lower total sample sizes. The rbart model, which incorporates a random intercept, shows a similar pattern as the stan4bart model, but with an overall higher relative bias.

**Figure 1:** Bias of the estimates for all simulated datasets, all  $ICC$  values, and differing within group effect sizes using four models. bart = single level BART model, gbart = single level BART model with group-dummies, rbart = multilevel BART model incorporating a random intercept, stan4bart = multilevel BART model combined with Stan to model the random parts of the models. dataset names show the total sample size, number of groups, and group size. Bias is based on 6 replicates for every dataset.

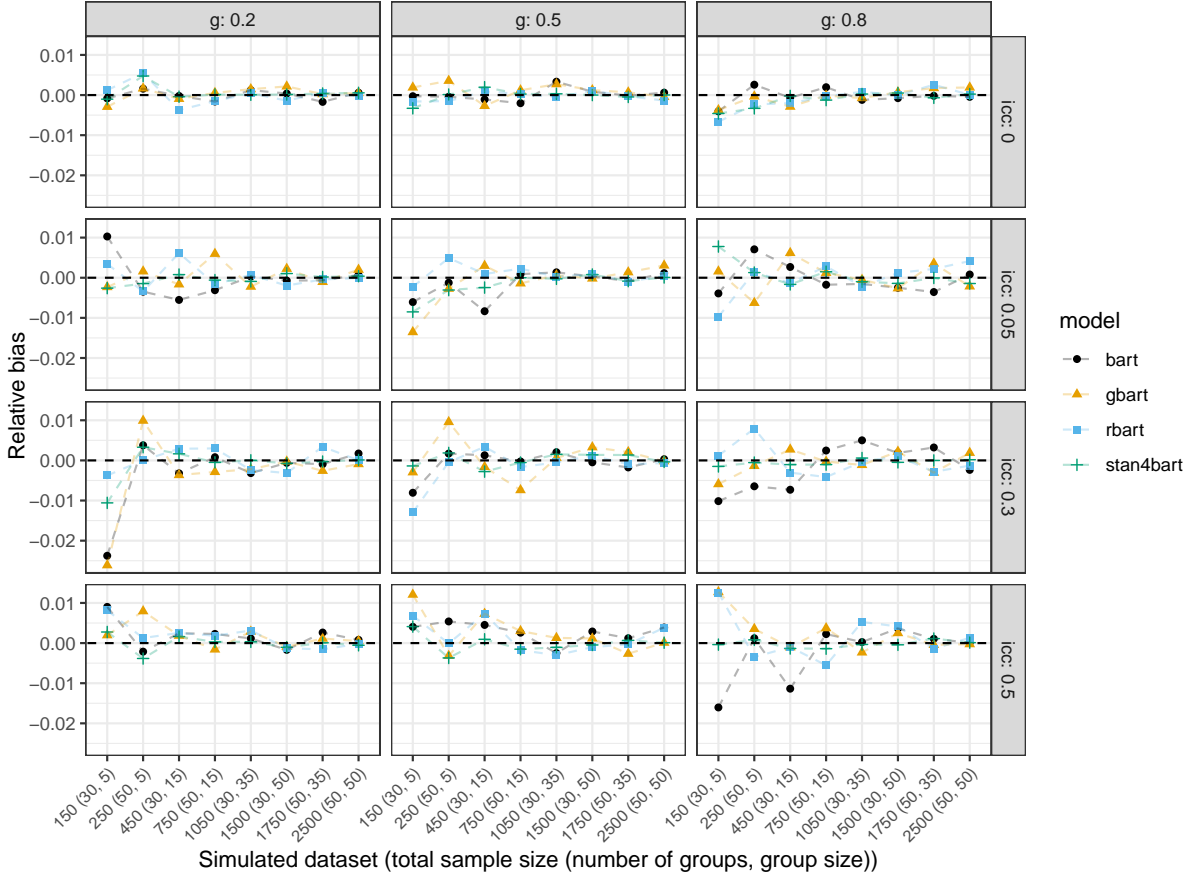
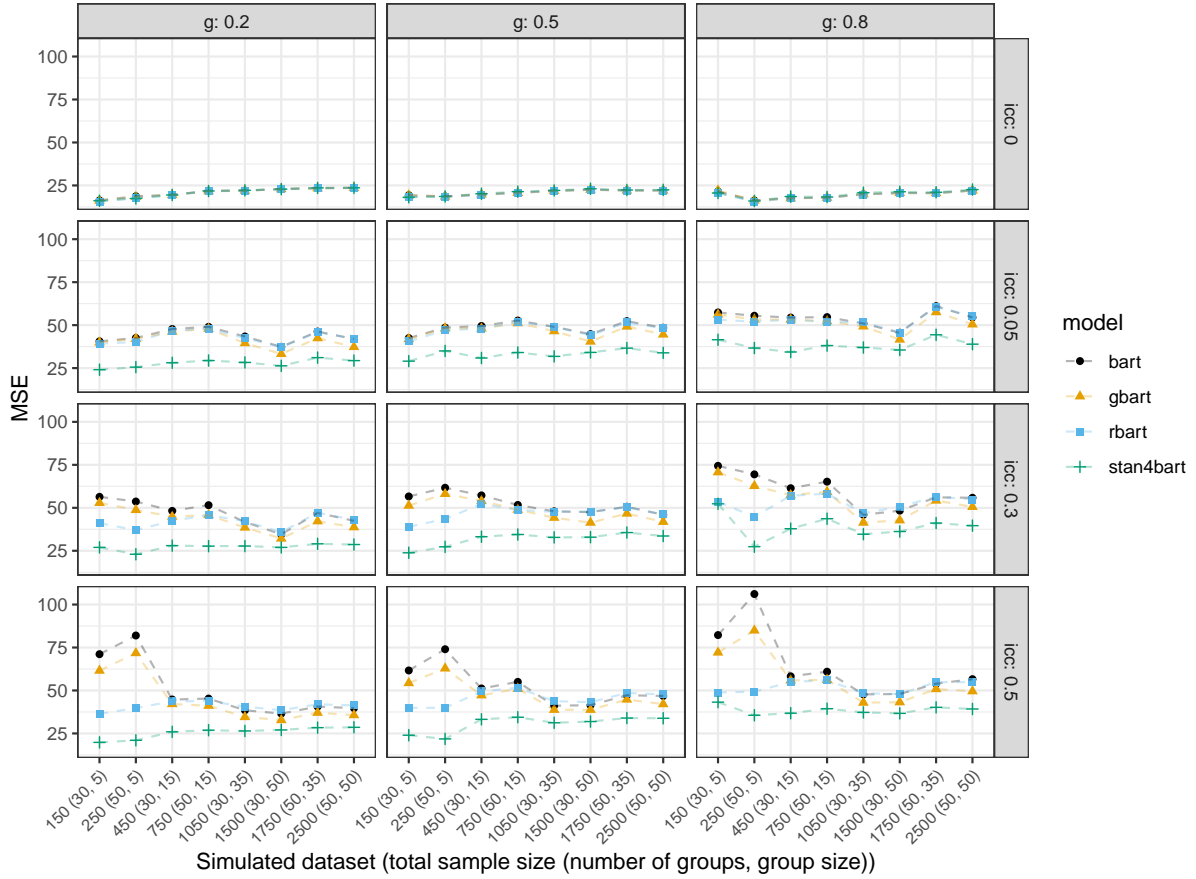


Figure 2 shows a similar pattern as figure 1. When there is no multilevel structure in datasets,  $ICC = 0$  the models perform well and almost exactly the same. When increasing the  $ICC$ , we can start to see a divide in the performance of the models. When  $ICC = .05$  the models bart, gbart and rbart perform similarly, with gbart outperforming the other two models when the group sizes are 35 or larger. However, stan4bart consistently outperforms the other three models in terms of MSE. Increasing the  $ICC$  to .3, divides the performance of the models further: in small datasets bart now has the highest MSE with gbart performing only a little better. rbart performs better than bart and gbart, but when the datasets increase in size, it performs similar to bart. Again, gbart outperforms bart and rbart when the groups sizes are 35 or larger. stan4bart is still has the overall lowest MSE. Increasing the  $ICC$  further to .5 exaggerates these patterns. Over all  $ICC$  values, increasing  $\gamma$  increases the MSE of the models.

**Figure 2:** Mean Squared Error (MSE) of the estimates for all simulated datasets, all ICC values, and differing within group effect sizes using four models. bart = single level BART model, gbart = single level BART model with group-dummies, rbart = multilevel BART model incorporating a random intercept, stan4bart = multilevel BART model combined with Stan to model the random parts of the models. dataset names show the total sample size, number of groups, and group size. MSE is based on 6 replicates for every dataset.



## 4 Discussion

In this research report I have investigated the performance of different BART models in terms of relative bias and MSE of the estimates. I considered four different models: a single-level BART model (*bart*), a single-level BART model including a group-dummy (*gbart*), a multilevel BART model including a random intercept (*rbart*), and a multilevel BART model combining Stan and BART (*stan4bart*). Chen [2020], Tan et al. [2016] and Wundervald et al. [2022] found that a BART model including a random intercept performed better than a single-level BART model, which is partly in agreement with my results: the MSE for the *rbart* showed more consistency, but when the total sample sizes were large, the model did not perform better than the *bart* and *gbart* models. The results indicate that the *stan4bart* model performs best out of the four models: it shows to lowest relative bias as well as the lowest MSE. These results are agreement with Dorie et al. [2022], who found that the *stan4bart* algorithm performed better in terms of coverage of the population values and RMSE compared to single-level BART models, BART models including a random intercept, Bayesian Causal Forests (BCF), and parametric MLMs.

Building on these results, I will implement the *stan4bart* as an imputation method within the package *MICE* [Buuren and Groothuis-Oudshoorn, 2011] in my thesis. I will compare the performance of the *stan4bart* model to other imputation methods: *2l.pmm*, *2l.lmer*, *2l.pan*, *2l.jomo*, *rf* and single-level *pmm* and complete case analysis in the R-package *MICE* [Buuren and Groothuis-Oudshoorn, 2011]. They will be evaluated in terms of relative bias, modeled variance, and the 95% confidence interval coverage of the estimates [Oberman and Vink, 2023]. The simulation design will be extended to include more parameters: the missing data mechanism and amount of missingness. The simulation design is shown in table 2. For each combination of paramteres, a 1000 replicated datasets will be generated.

**Table 2:** Simulation design thesis

Parameter	Values
Number of clusters ( $j$ )	30, 50
Within-cluster sample size ( $n_j$ )	5, 15, 35, 50
Intraclass Correlation (ICC)	0, .05, .3, .5
Missing data mechanism	MAR, MCAR
Amount of missingness	0%, 25%, 50%
Within-group effect size ( $\gamma$ )	.2, .5, .8



## References

- Audigier, V., White, I. R., Jolani, S., Debray, T. P. A., Quartagno, M., Carpenter, J., Van Buuren, S., and Resche-Rigon, M. (2018). Multiple Imputation for Multilevel Data with Continuous and Binary Variables. *Statistical Science*, 33(2).
- Austin, P. C., White, I. R., Lee, D. S., and Van Buuren, S. (2021). Missing Data in Clinical Research: A Tutorial on Multiple Imputation. *Canadian Journal of Cardiology*, 37(9):1322–1331.
- Bartlett, J. W., Seaman, S. R., White, I. R., Carpenter, J. R., and for the Alzheimer’s Disease Neuroimaging Initiative\* (2015). Multiple imputation of covariates by fully conditional specification: Accommodating the substantive model. *Statistical Methods in Medical Research*, 24(4):462–487.
- Breiman, L., Friedman, J. H., Olshen, R. A., and Stone, C. J. (1984). *Classification And Regression Trees*. Routledge, 1 edition.
- Burgette, L. F. and Reiter, J. P. (2010). Multiple Imputation for Missing Data via Sequential Regression Trees. *American Journal of Epidemiology*, 172(9):1070–1076.
- Buuren, S. V. and Groothuis-Oudshoorn, K. (2011). **Mice** : Multivariate Imputation by Chained Equations in R. *Journal of Statistical Software*, 45(3).
- Carpenter, J. R. and Kenward, M. G. (2013). *Multiple Imputation and Its Application*. Wiley, 1 edition.
- Chen, S. (2020). *A New Multilevel Bayesian Nonparametric Algorithm and Its Application in Causal Inference*. PhD thesis, Texas A&M University.
- Chipman, H. A., George, E. I., and McCulloch, R. E. (2010). BART: Bayesian additive regression trees. *The Annals of Applied Statistics*, 4(1).
- Dong, M. and Mitani, A. (2023). Multiple imputation methods for missing multilevel ordinal outcomes. *BMC Medical Research Methodology*, 23(1):112.
- Dorie, V. (2023a). Dbarts: Discrete bayesian additive regression trees sampler.
- Dorie, V. (2023b). *Stan4bart: Bayesian Additive Regression Trees with Stan-Sampled Parametric Extensions*.
- Dorie, V., Perrett, G., Hill, J. L., and Goodrich, B. (2022). Stan and BART for Causal Inference: Estimating Heterogeneous Treatment Effects Using the Power of Stan and the Flexibility of Machine Learning. *Entropy*, 24(12):1782.
- Enders, C. K. (2017). Multiple imputation as a flexible tool for missing data handling in clinical research. *Behaviour Research and Therapy*, 98:4–18.
- Enders, C. K., Du, H., and Keller, B. T. (2020). A model-based imputation procedure for multilevel regression models with random coefficients, interaction effects, and nonlinear terms. *Psychological Methods*, 25(1):88–112.
- Enders, C. K., Hayes, T., and Du, H. (2018a). A Comparison of Multilevel Imputation Schemes for Random Coefficient Models: Fully Conditional Specification and Joint Model Imputation with Random Covariance Matrices. *Multivariate Behavioral Research*, 53(5):695–713.
- Enders, C. K., Keller, B. T., and Levy, R. (2018b). A fully conditional specification approach to multilevel imputation of categorical and continuous variables. *Psychological Methods*, 23(2):298–317.
- Enders, C. K., Mistler, S. A., and Keller, B. T. (2016). Multilevel multiple imputation: A review and evaluation of joint modeling and chained equations imputation. *Psychological Methods*, 21(2):222–240.
- Grund, S., Lüdtke, O., and Robitzsch, A. (2016). Multiple imputation of missing covariate values in multilevel models with random slopes: A cautionary note. *Behavior Research Methods*, 48(2):640–649.
- Grund, S., Lüdtke, O., and Robitzsch, A. (2018a). Multiple Imputation of Missing Data at Level 2: A Comparison of Fully Conditional and Joint Modeling in Multilevel Designs. *Journal of Educational and Behavioral Statistics*, 43(3):316–353.

- Grund, S., Lüdtke, O., and Robitzsch, A. (2018b). Multiple Imputation of Missing Data for Multilevel Models: Simulations and Recommendations. *Organizational Research Methods*, 21(1):111–149.
- Grund, S., Lüdtke, O., and Robitzsch, A. (2021). Multiple imputation of missing data in multilevel models with the R package mdmb: A flexible sequential modeling approach. *Behavior Research Methods*, 53(6):2631–2649.
- Gulliford, M., Adams, G., Ukoumunne, O., Latinovic, R., Chinn, S., and Campbell, M. (2005). Intraclass correlation coefficient and outcome prevalence are associated in clustered binary data. *Journal of Clinical Epidemiology*, 58(3):246–251.
- Gulliford, M. C., Ukoumunne, O. C., and Chinn, S. (1999). Components of Variance and Intraclass Correlations for the Design of Community-based Surveys and Intervention Studies: Data from the Health Survey for England 1994. *American Journal of Epidemiology*, 149(9):876–883.
- Hastie, T. J., editor (2017). *Statistical Models in S*. Routledge, 1st edition.
- Hill, J., Linero, A., and Murray, J. (2020). Bayesian Additive Regression Trees: A Review and Look Forward. *Annual Review of Statistics and Its Application*, 7(1):251–278.
- Hox, J. and Roberts, J. K., editors (2011). *Handbook of Advanced Multilevel Analysis*. Routledge, 0 edition.
- Hox, J. J., Moerbeek, M., and Van De Schoot, R. (2017). *Multilevel Analysis: Techniques and Applications*. Routledge, Third edition. — New York, NY : Routledge, 2017. —, 3 edition.
- Hughes, R. A., White, I. R., Seaman, S. R., Carpenter, J. R., Tilling, K., and Sterne, J. A. (2014). Joint modelling rationale for chained equations. *BMC Medical Research Methodology*, 14(1):28.
- James, G., Witten, D., Hastie, T., and Tibshirani, R. (2021). *An Introduction to Statistical Learning: With Applications in R*. Springer Texts in Statistics. Springer US, New York, NY.
- Kang, H. (2013). The prevention and handling of the missing data. *Korean Journal of Anesthesiology*, 64(5):402.
- Lee, D., Carpenter, B., Li, P., Morris, M., Betancourt, M., Maverickg, Brubaker, M., Trangucci, R., Inacio, M., Kucukelbir, A., Buildbot, S., Bgoodri, Seantalts, Arnold, J., Tran, D., Hoffman, M., Margossian, C., Modrák, M., Adler, A., Sakrejda, K., Stukalov, A., Lawrence, M., Goedman, R. J., Van Horn, K. S., Vehtari, A., Gabry, J., Casallas, J. S., and Bales, B. (2017). Stan-dev/stan: V2.17.1. Zenodo.
- Lin, S. and Luo, W. (2019). A New Multilevel CART Algorithm for Multilevel Data with Binary Outcomes. *Multivariate Behavioral Research*, 54(4):578–592.
- Lüdtke, O., Robitzsch, A., and Grund, S. (2017). Multiple imputation of missing data in multilevel designs: A comparison of different strategies. *Psychological Methods*, 22(1):141–165.
- Meng, X.-L. (1994). Multiple-imputation inferences with uncongenial sources of input. *Statistical science*, pages 538–558.
- Mistler, S. A. and Enders, C. K. (2017). A Comparison of Joint Model and Fully Conditional Specification Imputation for Multilevel Missing Data. *Journal of Educational and Behavioral Statistics*, 42(4):432–466.
- Morris, T. P., White, I. R., and Crowther, M. J. (2019). Using simulation studies to evaluate statistical methods. *Statistics in Medicine*, 38(11):2074–2102.
- Murray, D. M. and Blitstein, J. L. (2003). Methods To Reduce The Impact Of Intraclass Correlation In Group-Randomized Trials. *Evaluation Review*, 27(1):79–103.
- Oberman, H. I. and Vink, G. (2023). Toward a standardized evaluation of imputation methodology. *Biometrical Journal*, page 2200107.
- Quartagno, M. and Carpenter, J. R. (2022). Substantive model compatible multilevel multiple imputation: A joint modeling approach. *Statistics in Medicine*, 41(25):5000–5015.
- R Core Team (2023). *R: A Language and Environment for Statistical Computing*. Vienna, Austria.

- Resche-Rigon, M. and White, I. R. (2018). Multiple imputation by chained equations for systematically and sporadically missing multilevel data. *Statistical Methods in Medical Research*, 27(6):1634–1649.
- Rubin, D. B. (1976). Inference and missing data. *Biometrika*, 63(3):581–592.
- Rubin, D. B. (1987). *Multiple Imputation for Nonresponse in Surveys*. Wiley, New York.
- Salditt, M., Humberg, S., and Nestler, S. (2023). Gradient Tree Boosting for Hierarchical Data. *Multivariate Behavioral Research*, pages 1–27.
- Shieh, G. (2012). A comparison of two indices for the intraclass correlation coefficient. *Behavior Research Methods*, 44(4):1212–1223.
- Silva, G. C. and Gutman, R. (2022). Multiple imputation procedures for estimating causal effects with multiple treatments with application to the comparison of healthcare providers. *Statistics in Medicine*, 41(1):208–226.
- Taljaard, M., Donner, A., and Klar, N. (2008). Imputation Strategies for Missing Continuous Outcomes in Cluster Randomized Trials. *Biometrical Journal*, 50(3):329–345.
- Tan, Y. V., Flannagan, C. A. C., and Elliott, M. R. (2016). Predicting human-driving behavior to help driverless vehicles drive: Random intercept Bayesian Additive Regression Trees.
- Van Buuren, S. (2007). Multiple imputation of discrete and continuous data by fully conditional specification. *Statistical Methods in Medical Research*, 16(3):219–242.
- van Buuren, S. (2018). *Flexible Imputation of Missing Data*. Chapman & Hall/CRC Interdisciplinary Statistics Series. CRC Press, Taylor & Francis Group, Boca Raton London New York, second edition edition.
- Wagner, J., West, B. T., Elliott, M. R., and Coffey, S. (2020). Comparing the Ability of Regression Modeling and Bayesian Additive Regression Trees to Predict Costs in a Responsive Survey Design Context. *Journal of Official Statistics*, 36(4):907–931.
- Waljee, A. K., Mukherjee, A., Singal, A. G., Zhang, Y., Warren, J., Balis, U., Marrero, J., Zhu, J., and Higgins, P. D. (2013). Comparison of imputation methods for missing laboratory data in medicine. *BMJ Open*, 3(8):e002847.
- Wundervald, B., Parnell, A., and Domijan, K. (2022). Hierarchical Embedded Bayesian Additive Regression Trees.
- Xu, D., Daniels, M. J., and Winterstein, A. G. (2016). Sequential BART for imputation of missing covariates. *Biostatistics*, 17(3):589–602.