

Aus dem Institut für Medizinische Informationsverarbeitung, Biometrie
und Epidemiologie (IBE)
der Ludwig-Maximilians-Universität München



Dissertation
zum Erwerb des Doctor of Philosophy (Ph.D.)
an der Medizinischen Fakultät der
Ludwig-Maximilians-Universität zu München

**Handling of realistic missing data scenarios in clinical trials
using machine learning techniques**

vorgelegt von:

.....Halimuniyazi..Haliduola.....

aus:

.....Altay, China.....

Jahr:

.....2023.....

Mit Genehmigung der Medizinischen Fakultät der
Ludwig-Maximilians-Universität zu München

First evaluator (1. TAC member): Prof. Dr. Ulrich Mansmann

Second evaluator (2. TAC member): Prof. Dr. Anne-Laure Boulesteix

Third evaluator: Priv. Doz. Dr. Michaela Schunk

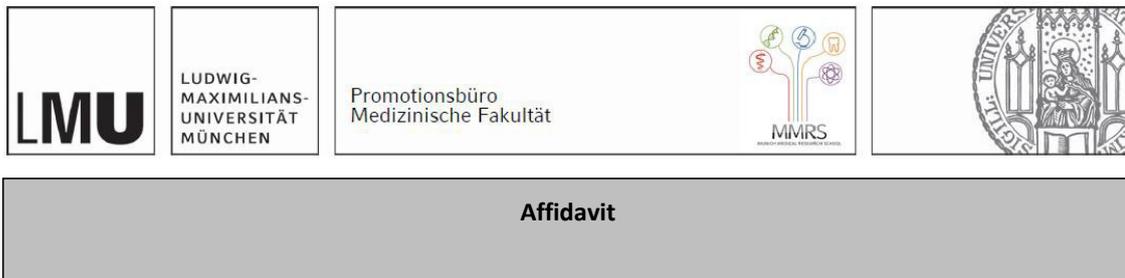
Fourth evaluator: Priv. Doz. Dr. Steffen Unkel

Dean: Prof. Dr. med. Thomas Gudermann

date of the defense:

03-Mar-2023

Affidavit



Haliduola, Halimuniyazi
Surname, first name

I hereby declare, that the submitted thesis entitled:

Handling of realistic missing data scenarios in clinical trials using machine learning techniques
.....

is my own work. I have only used the sources indicated and have not made unauthorised use of services of a third party. Where the work of others has been quoted or reproduced, the source is always given.

I further declare that the submitted thesis or parts thereof have not been presented as part of an examination degree to any other university.

Munich, 13-Mar-2023
place, date

Halimuniyazi Haliduola
Signature doctoral candidate

Confirmation of congruency



Confirmation of congruency between printed and electronic version of the doctoral thesis

Haliduola, Halimuniyazi
Surname, first name

I hereby declare, that the submitted thesis entitled:

Handling of realistic missing data scenarios in clinical trials using machine learning techniques
.....

is congruent with the printed version both in content and format.

Munich, 13-Mar-2023
place, date

Halimuniyazi Haliduola
Signature doctoral candidate

Table of content

Affidavit	3
Confirmation of congruency	4
Table of content.....	5
List of abbreviations	6
List of publications	7
Your contribution to the publications	8
1.1 Contribution to paper I	8
1.2 Contribution to paper II	9
2. Introductory summary	11
2.1 Motivation and main idea	11
2.1.1 Main idea for Paper I.....	12
2.1.2 Main idea for Paper II.....	13
2.2 Proposals	15
2.2.1 Proposal in Paper I	15
2.2.2 Proposal in Paper II	16
3. Paper I	18
4. Paper II	39
References	50
Acknowledgements.....	52

List of abbreviations

ANCOVA = analysis of covariance

bagging = bootstrap aggregating

CI = confidence interval

CV = cross validation

MAE = mean absolute error

MAR = missing at random

MCAR = missing completely at random

MI = multiple imputation

MMRM = mixed model for repeat measurement

MNAR = missing not at random

MSE = mean squared error

QRF = quantile regression forests

PMM = pattern mixture model

RF = random forests

RNN = recurrent neural network

SM = selection model

SMOTE = synthetic minority oversampling technique

SMOTER = synthetic minority oversampling technique for regression

SPM = shared parameter model

UBL = utility-based learning

UBR = utility-based regression

List of publications

Paper I

Haliduola, H. N., Bretz, F., Mansmann, U. (2022). Missing data imputation in clinical trials using recurrent neural network facilitated by clustering and oversampling. *Biometrical Journal*, 64 (5) (2022), pp. 863-882, <https://doi.org/10.1002/bimj.202000393>

Paper II

Haliduola, H. N., Bretz, F., Mansmann, U. (2022). Missing data imputation using utility-based regression and sampling approaches. *Computer Methods and Programs in Biomedicine*, Volume 226, 2022, 107172, ISSN 0169-2607, <https://doi.org/10.1016/j.cmpb.2022.107172>

Your contribution to the publications

1.1 Contribution to paper I

My contribution to Paper I begins with a literature review. Through reading the literatures, I realized that many classical approaches in dealing with missing data have been developed with underlying assumption of MAR (missing at random) or MNAR (missing not at random), respectively. For example, mixed model for repeated measures has been used under assumption of MAR; and pattern mixture model (Little, 1993), selection model (Rubin, 1976), and shared parameter model (Little, 1995) have been implemented as sensitivity analysis with underlying assumption of MNAR. However, missing data in realistic situations can be more complex, for example, MNAR and MAR mixed together in the same dataset. Many studies have shown that those classical models can reduce bias appropriately when (and only when) their underlying assumptions are satisfied. However, in some real-life situations (e.g., when MNAR and MAR mixed together), these methods perform poorly because they rely heavily on underlying assumptions (Enders, 2010). Unfortunately, it is difficult to test these underlying assumptions, so there is no easy solution to evaluate the performance of those models in analyzing real study data. I discussed these findings with my main supervisor Prof. Dr. Ulrich Mansmann and proposed to investigate machine-learning-based methods to impute the missing data in clinical trial settings. Prof. Mansmann was very supportive on this idea, and he pointed out that the new method should be able to handle MNAR problems first, and at the same time, it should also be able to handle MAR problems.

Following Prof. Mansmann's advice, a missing data prediction approach that based on machine learning techniques was developed to handle missing data problem in real-life situations (i.e., when both MNAR and MAR exist in the same dataset) in Paper I. The main idea is to handle MNAR by focusing on (giving more weights to) the missing part, meanwhile, and also to handle the MAR data by looking for precise individual (subject-level) information. The problem of MNAR is seen as an imbalanced machine learning exercise, i.e., to oversample the minority cases to compensate for the data that are MNAR in certain area. It should be noted that this framework is original, it was never published or discussed in anywhere else.

We evaluate our approach through comprehensive and objective simulation studies. Initially, I considered one scenario only in the simulation. After discussing with Prof. Mansmann and Prof. Anne-Laure Boulesteix (my second supervisor), to fully evaluate the performance of the proposed method, they suggested to extend the simulation studies to consider different scenarios (i.e., different sample size, different dropout proportions, and different missingness starting time-point). They pointed out that in light of the "evidence-based computational statistics" (Boulesteix et al. 2017a, b), the proposed methodology needs to be evaluated in an objective manner through simulation studies and needs to emphasize the plausibility of real-life simulation scenarios.

Another contribution from my side is the programming of the entire workflow. As shown in Figure 3 (in section 2.2.1), different software/tools have been used in different steps according to the availability of relevant package/library/standard toolbox. For example:

- SAS 9.4 was used for the simulation study data creation and statistical analysis for both simulation study and real data example. The following models were implemented using SAS: mixed model for repeat measurement (MMRM), analysis of covariance (ANCOVA), selection model (SM), pattern mixture model (PMM) and shared parameter model (SPM);
- R (version 3.6.0) is used in the longitudinal clustering (k-means trajectories was implemented using R package “KML” version 2011, Genolini and Falissard, 2011);
- Python (version 3.6) was used in the RNN implementation via library Keras (version 2.2.4, Falbel et al., 2015).

A detailed supporting information package was prepared by me, which includes all programming codes for the entire workflow (including SAS, R and Python codes), intermediate results, datasets and outputs. My main supervisor reviewed the package and validated some key steps and key results. The supporting information package was submitted to the journal for reproducibility check.

I wrote a draft of the manuscript, which my supervisors thoroughly reviewed and provided valuable constructive comments to improve it for submission to the journal.

1.2 Contribution to paper II

After working on the Paper I and reviewing literatures about the utility-based-learning, I realized that the standard predictive error measures like mean squared error (MSE) or mean absolute error (MAE) do not perform well in the regression tasks for the data with imbalanced distribution. For example, in some clinical studies, the extreme values of outcome variable are often missing due to reasons that are related to subjects' health status (subjects who dropout because the disease status is worsening), i.e., a MNAR issue. They have the disadvantage of being insensitive to the position of the outcome variable values (Ribeiro, 2011). In addition, the simple random oversampling with replacement will cause many duplicated cases in minority cluster (i.e., in that case the training data becomes very specific) and hence causes overfitting of the model. Therefore, after discussing with Prof. Mansmann, I proposed to incorporate the utility-based error measurement method (which is sensitive enough to the specific location of predictive error in the target variable scale) and sampling approach (i.e., the SMOTER = synthetic minority oversampling technique for regression, Torgo et al. 2013) to handle the issue of MNAR in clinical studies. This proposal was fully supported by my supervisors.

Like Paper I, we evaluate our proposed approach through comprehensive and objective simulation studies. We created random data for 600 subjects. The created target variable

and the covariates follow a normal distribution, missing data flags follow binomial distribution (i.e., we created mutually exclusive binary flags for MNAR, MAR and MCAR data). We simultaneously created correlated binary and normal data following a point-biserial method. In the utility-based regression process, the fixed coefficient of relevance function was used initially. After discussing with my main supervisor, he suggested to perform sensitivity analysis by defining a set of coefficients for relevance function to show how the chosen relevance function affects the imputation results.

I did the programming of the entire workflow using R. The R package “BinNor” (Amatya, 2020) was used for the data generation, and R package “Utility-Based Learning” (i.e., “UBL”, developed by Branco et al. in 2017) was used for the utility-based regression and SMOTER. The performance of the proposed method was compared with other standard methods like multiple imputation (MI, using R package “MICE”, van Buuren et al. 2011), random forests (RF, using R package “randomForest”, Liaw and Wiener, 2018), and quantile regression forests (QRF, using R package “quantregForest”, Meinshausen, 2017). All R codes, intermediate results, datasets and outputs were submitted to the journal for reproducibility check.

I wrote a draft of the manuscript, and my supervisors thoroughly reviewed it and provided valuable constructive comments to improve it for submission to the journal.

2. Introductory summary

2.1 Motivation and main idea

Missing data problem is a common challenge when designing and analyzing clinical trials, which are the data that are needed for the main analyses but are not collected. If the missing data are not properly imputed/handled, they may cause following issues: reduce the statistical power of the important analysis; they may bias/confound the treatment effect estimation; they may cause an underestimation of the variability in target variable. Three different types of missingness are defined in Rubin's 1976 paper. (1) MCAR (missing completely at random): when data are MCAR, "the probability of missingness does not depend on observed or unobserved measurements", for example, subjects who drop-out from the trial due to the reasons that are not related to their health status. (2) MAR (missing at random): when data are MAR, "the probability of missingness depends only on observed measurements conditional on the covariates in the model", for example, younger subjects (those who don't think it is necessary to measure their blood pressure as they consider themselves healthier) may more likely to have missing blood pressure. (3) MNAR (missing not at random): when data are MNAR, "the probability of missingness depends on unobserved measurements", for example, subjects leave the trial because of "lack of efficacy" (i.e., they are not convinced by effectiveness of the study drug and hence dropout from the trial).

Although all three types of missing data are well defined, it is very difficult to determine the association between missing data and unobserved outcomes in the real-world data; in other words, it is very difficult to justify the MAR assumption in any realistic situation. As EMA suggested in 2010, a combined strategy can be used, e.g., treat the discontinuations due to "lack of efficacy" as MNAR data, and treat the discontinuations due to "lost to follow-up" as MAR data.

Many statistical methods have been developed to handle missing data under the prerequisite assumption of either MNAR or MAR. However, in the real world, missing data are often mixed with different types of missing mechanisms. This violates the basic assumptions for missing data (i.e., either MNAR or MAR), which leads to a degradation in the processing performance of these methods (Enders, 2010). To handle the missing data problem in real-life situations (e.g., MNAR and MAR mixed together in the same dataset), we propose a missing data prediction framework that are based on machine learning techniques. As Breiman pointed out in his 2001 paper, in the statistical (machine) learning exercise, "the goal is not interpretability, but accurate information". Along this line of thought, our methods handle MNAR by focusing on (giving more sample weights to) the missing part, meanwhile, and also to handle the MAR data by looking for precise individual (subject-level) information. The problem of MNAR is seen as an imbalanced machine learning exercise, i.e., to oversample the minority cases to compensate for the data that are MNAR in certain area.

2.1.1 Main idea for Paper I

In Paper I, to handle the missing data problem in real-life situations (e.g., MNAR and MAR mixed together in the same dataset), we propose a missing data prediction framework that are based on machine learning techniques. We evaluate our proposed approach through objective simulation studies. As mentioned above, our methods handle MNAR by focusing on (giving more sample weights to) the missing part, meanwhile, and also to handle the MAR data by looking for precise individual (subject-level) information. The problem of MNAR is seen as an imbalanced machine learning exercise, i.e., to oversample the minority cases (Weiss, 2013) to compensate for the data that are MNAR in certain area. To be able to use the oversampling in longitudinal outcome variable (continuous scale), clustering through k-mean algorithm (Gower, 1971) needs to be performed the first. By using k-mean algorithm, subjects are clustered into “low responders”, “medium responders” and “good responders” according to their efficacy trajectories over-time. We consider the subjects discontinued from the study due to “lack of efficacy” as MNAR in our simulations. Therefore, the MNAR subjects are mostly low responders. Therefore, the amount of non-missing available training data in this category/cluster can be less than the other categories/clusters (therefore, such distribution of training data is imbalanced in nature). To compensate for the data that are MNAR in that area, and also to avoid the imputation results being driven by the non-missing data (from the subjects who completed the study) to a population average level, random oversampling (with replacement) for the minority cases is needed. See Figure 1 for the distribution of target variable in simulation data. The full data (including the non-missing data and the original values that are set to “missing” in the simulation) are presented in the left side, black dots are for non-missing data, red stars are for MNAR, and blue triangles are for MAR. The non-missing data (i.e., the original training data that to be used in learning process) are displayed in the middle again and clustered/categorized into three classes: green dots are low responders, orange dots are medium responders, and purple dots are good responders. It is shown that the available (non-missing) data cannot represent well the full data considering the existence of MNAR data. The green dots (i.e., the low responders) are relatively less in the original data (therefore, they are the minority cases). When we apply the random oversampling with replacement in this minority cluster, as displayed in the right side, the amount of data in the low responder area are increased to compensate for MNAR data in that area.

In this research, we use RNN (recurrent neural networks) model to fit the longitudinal subject trajectories. “RNN is a type of neural network that can learn from the past to predict the future outcomes” (Rumelhart et al., 1986; Schmidhuber, 1993). This makes it a useful tool to model the longitudinal clinical data. RNN allows us to model nonlinear data without any specific domain knowledge about the relationship between the variables, it automatically learns from the given data to optimize the model parameters and then provides predictions for the new test data. The optimal model hyperparameters (the specific RNN architectures in this case) are determined via the “bias-variance tradeoff

approach” (Claesen and De Moor, 2015). To consider the uncertainty of the single prediction and to optimize the prediction accuracy, bootstrap aggregating is used in this study. Considering the “evidence-based computational statistics” (Boulesteix et al. 2017a, b), the proposed method is evaluated through comprehensive simulation studies and implemented in real data from a clinical trial. The imputation results are evaluated at different levels, i.e., the individual subject level (e.g., prediction accuracy for the specific subjects) and the overall study population level (e.g., the treatment effect estimation).

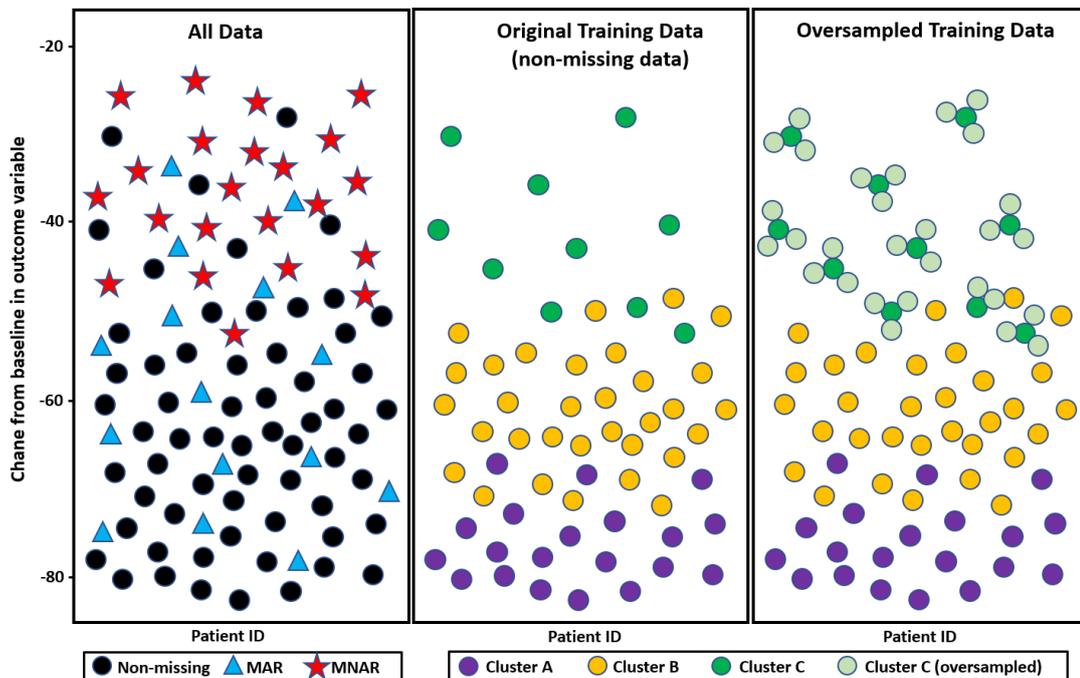


Figure 1. Illustration of clustering and oversampling for target variable – “Change from baseline in outcome variable”

2.1.2 Main idea for Paper II

In Paper I, we used a simple sampling approach (i.e., random oversampling with replacement) and a standard error measurement method (i.e., MSE, mean squared error). However, there are some drawbacks in those methods. The simple random oversampling with replacement will cause many duplicated cases in the minority cluster (i.e., the training data becomes very specific) and hence causes overfitting of the model (Chawla, et al. 2002). The standard error measurement method like MSE do not perform well in the regression tasks for the data with imbalanced distribution. For example, in some clinical studies, the extreme values of outcome variable are often missing due to reasons that are related to subjects’ health status (i.e., subjects who dropout because the disease status is worsening), i.e., a MNAR issue. They have the disadvantage of being insensitive to the position of the outcome variable values (Ribeiro, 2011). Take Figure 2 as an example, considering the red dots (i.e., the MNAR data), distribution of the black dots (i.e., the available non-missing data) are imbalanced over the target variable range (i.e., because of the MNAR in the high value area, there are less available data in that area). In that case, if the error measurement is not sensitive to the positions of the values of

target variable, during the training/learning process, the area tends to have MNAR will get less focus because of the smaller amount of data. This will cause bias in the aggregated statistical analysis as the impact of missing data are simply ignored. In that case, during the training/learning process, giving more focus on the area that tend to have MNAR data is really necessary and very important. This will compensate for the data that are MNAR in certain area and also will prevent the predictions being driven by the majority data (that located in other positions of the outcome variable scale). In summary, it is really necessary and very important to have error measurement methods that are sensitive to the location of the values, which can cope with the problem of imbalanced distribution of the target variable.

In Paper II, to overcome the above-mentioned drawbacks, we use the “synthetic minority oversampling technique for regression” (SMOTER) (Torgo et al. 2013) to oversample the relevant rare cases; and we use the imbalanced learning algorithm “utility-based regression” (UBR) (Torgo and Ribeiro, 2007) to consider both the importance/relevance of the locations (i.e., values of target variable) and the prediction errors simultaneously in the model parameters optimization process. We use the Quantile regression forests (QRF, Meinshausen, 2006) to estimate the conditional probability density (CPD) given covariates. The optimization process aims to maximize the integral for the product of the CPD multiplied by the utility function in each case. We evaluate our proposed approach through simulation studies with realistic missing data situations (i.e., when the MNAR, MAR and MCAR mixed together in the same dataset). The performance of proposed approach is evaluated objectively and compared with the commonly used missing data handling methods like random forests (RF) and conventional multiple imputation (MI). We also implemented our method in an antidepressant clinical trial data (publicly available datasets).

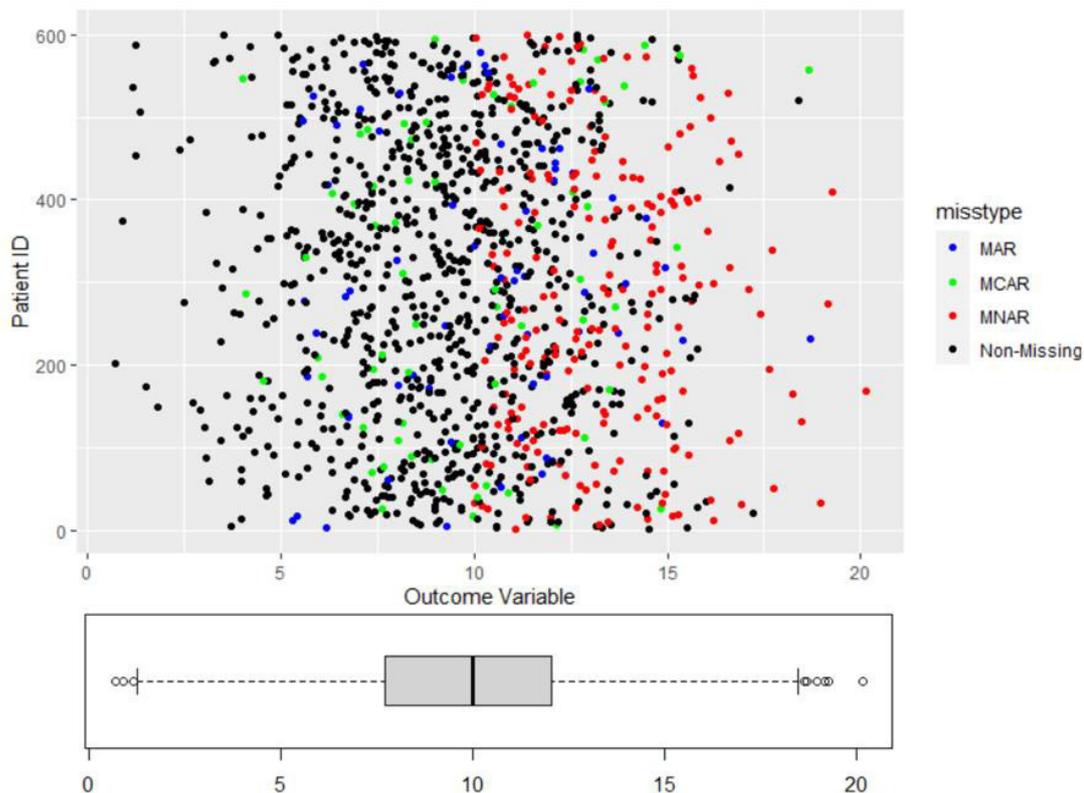


Figure 2. Simulation data – scatter plot and boxplot for the target outcome variable.

2.2 Proposals

In following sections, we elaborate our proposals in Paper I and Paper II respectively.

2.2.1 Proposal in Paper I

The proposed framework includes several necessary components, see Figure 3 for the proposed workflow. The details of each step and the methods please see Paper I.

- The step 1 is clustering, i.e., to cluster the subjects according to their longitudinal trajectories by treatment arm. This step is a key element in our approach for following reasons: 1. the clustering results (like the good, medium and low response) are used in the “stratified k-fold cross-validation (CV)” and the random oversampling process as the “categorization” of the outcome variable in continuous scale to balance the “majority” and “minority” classes; 2. the clustering results are used in the recurrent neural network model to indicate the longitudinal pattern of subjects’ trajectories (this is very important to borrow relevant information from the similar subjects).
- The step 2 is model selection, i.e., to optimize the RNN hyper-parameters (that have to be determined before the commencing of training process) via a “bias-variance tradeoff approach”.
- The step 3 is bootstrap aggregating, i.e., to generate 100 bagging samples with the minority clusters properly oversampled, and to impute the missing data

within each bagging. In this step, the RNN model (optimal model in terms of hyper-parameters) is executed for 100 times, in each execution, the internal parameters (also called “weights”) of the RNN model are updated and individual predictions for the dropouts are provided.

- The step 4 is to get final imputation for each subject at each visit by averaging the 100 predicted values from bagging.
- The step 5 is imputation performance evaluation, i.e., to evaluate the results in individual subject and study population level.

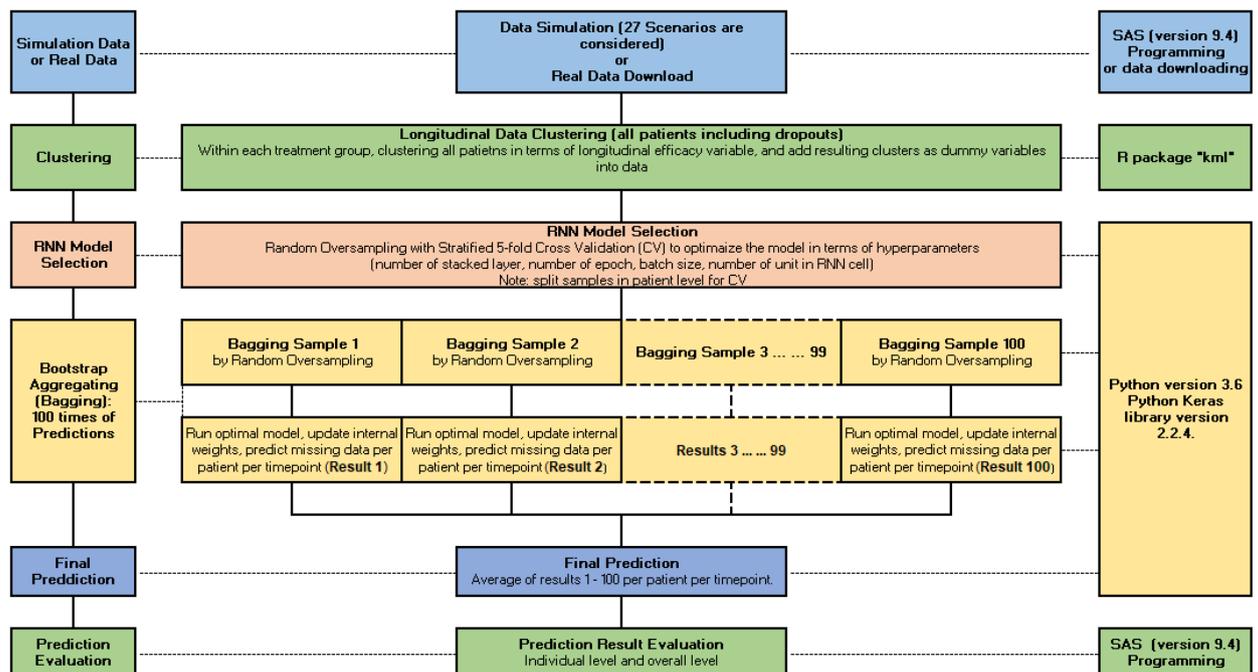


Figure 3. Diagram of Workflow– overall process for handling missing data in longitudinal continuous variable.

2.2.2 Proposal in Paper II

In Paper II, we aim to handle realistic situation when the MNAR, MAR and MCAR mixed together in the same dataset. We focus on a for a continuous target variable in clinical trial setting. The problem of MNAR is seen as an imbalanced machine learning exercise in our study. We propose a hybrid approach which consists of “synthetic minority oversampling technique for regression” (Torgo et al. 2013) and “utility-based regression” (Torgo and Ribeiro, 2007). In the first step, a relevance function is assigned to the outcome variable based on the distribution of the original training data (i.e., the original data with non-missing values) and a cut-off value for triggering the oversampling is defined. The second step is to pre-process the training data, in which the SMOTER (“synthetic minority oversampling technique for regression”) is used to oversample the relevant cases (i.e., the cases with relevance function values that are greater than the prespecified cut-off value in the first step). The third step is to apply the UBR (“utility-based regression”) on the processed training data (with relevant cases oversampled), then to

optimize the UBR parameters (i.e., model internal parameters) by simultaneously maximizing the relevance and minimizing the prediction error. In the final step, we use the optimized model (with final internal parameters) to predict the missing values for the target outcome variable. The details of each step and the methods please see Paper II.

3. Paper I

RESEARCH ARTICLE

Missing data imputation in clinical trials using recurrent neural network facilitated by clustering and oversampling

Halimu N. Haliduola^{1,2}  | Frank Bretz^{3,4}  | Ulrich Mansmann¹ 

¹Institute for Medical Information Processing, Biometry and Epidemiology (IBE), LMU Munich, Munich, Germany

²Alvotech Germany GmbH, Jülich, Germany

³Novartis Pharma AG, Basel, Switzerland

⁴Section for Medical Statistics, Medical University of Vienna, Vienna, Austria

Correspondence

Ulrich Mansmann, Institute for Medical Information Processing, Biometry and Epidemiology (IBE), LMU Munich, 81377 Munich, Germany.

Email:

mansmann@ibe.med.uni-muenchen.de



This article has earned an open data badge “**Reproducible Research**” for making publicly available the code necessary to reproduce the reported results. The results reported in this article could fully be reproduced.

Abstract

In clinical practice, the composition of missing data may be complex, for example, a mixture of missing at random (MAR) and missing not at random (MNAR) assumptions. Many methods under the assumption of MAR are available. Under the assumption of MNAR, likelihood-based methods require specification of the joint distribution of the data, and the missingness mechanism has been introduced as sensitivity analysis. These classic models heavily rely on the underlying assumption, and, in many realistic scenarios, they can produce unreliable estimates. In this paper, we develop a machine learning based missing data prediction framework with the aim of handling more realistic missing data scenarios. We use an imbalanced learning technique (i.e., oversampling of minority class) to handle the MNAR data. To implement oversampling in longitudinal continuous variable, we first perform clustering via k -mean trajectories. And use the recurrent neural network (RNN) to model the longitudinal data. Further, we apply bootstrap aggregating to improve the accuracy of prediction and also to consider the uncertainty of a single prediction. We evaluate the proposed method using simulated data. The prediction result is evaluated at the individual patient level and the overall population level. We demonstrate the powerful predictive capability of RNN for longitudinal data and its flexibility for nonlinear modeling. Overall, the proposed method provides an accurate individual prediction for both MAR and MNAR data and reduce the bias of missing data in treatment effect estimation when compared to standard methods and classic models. Finally, we implement the proposed method in a real dataset from an antidepressant clinical trial. In summary, this paper offers an opportunity to encourage the integration of machine learning strategies for handling of missing data in the analysis of randomized clinical trials.

KEYWORDS

clinical trial, k -mean clustering, missing data, oversampling, recurrent neural network

This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial](https://creativecommons.org/licenses/by-nc/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited and is not used for commercial purposes.

© 2022 The Authors. *Biometrical Journal* published by Wiley-VCH GmbH.

1 | INTRODUCTION

In clinical trials, missing data are the data that would be meaningful for the analysis but is not documented. Missing data, if not handled properly, will lead to lower statistical power as the sample size is reduced. In addition, dropouts from the trial may have poor outcomes or extreme values (e.g., treatment failure may lead to dropout). Therefore, the loss of these dropouts could lead to a bias in the estimated treatment effect (especially when the missing values are more likely in one treatment arm because it is not as effective as the other) and an underestimate of the variability. Missing data may also impact the external validity of the study outcome.

There are three types of missingness mechanisms (Rubin, 1976). (i) Missing completely at random (MCAR): if the probability of missingness does not depend on observed or unobserved measurements, for example, patients move to another city due to non-health related reasons. (ii) Missing at random (MAR): if the probability of missingness depends only on observed measurements conditional on the covariates in the model, for example, younger people may more likely to have blood pressure not measured. (iii) Missing not at random (MNAR): if the probability of missingness depends on unobserved measurements, for example, patients discontinue from the trial due to lack of efficacy, that is, patients do not come for the visit as the disease status worsening hence the data are missing.

The three types of missingness are clearly defined, but in practice it is typically not possible to be certain whether there is a relationship between missing values and the unobserved outcome variable. That is, it is not possible to ascertain whether the MAR/MCAR assumptions are appropriate in any practical situation. A mixed strategy may be considered. For example, assume that dropouts due to lack of efficacy are MNAR and loss to follow-up are MAR (European Medicines Agency, 2010). The consequence of MNAR is that the missing data cannot be simply predicted using observed data from that patient. In addition, the distribution of completers' data and MNAR data are different; thereby, it is not plausible to impute missing data using the completers' data. This is the central problem of missing data analysis in clinical trials (National Research Council of the National Academies, 2010).

Many established methods for handling missing data under the assumption of MAR are available. Under the alternative assumption of MNAR, the following classic models have been introduced as sensitivity analysis in the past decades: selection model (SM) (Heckman, 1976; Rubin, 1976), pattern mixture model (PMM) (Little, 1993, 1994, 1995), and shared parameter model (SPM) (Little, 1995). Let $(Y_{i,obs}, Y_{i,mis}, R_i)$ denote the data for i th patient, $Y_{i,obs}$ is for the observed component, $Y_{i,mis}$ is for the missing component, and R_i is the missingness indicator (1 = missing, 0 = observed). The full density function is described as $f(Y_{i,obs}, Y_{i,mis}, R_i | \Theta, \psi)$, where the parameters vectors Θ and ψ describe the response and missingness processes, respectively. The SM and PMM are developed by factorizing the full density function differently. The SM is based on the below factorization:

$$f(Y_{i,obs}, Y_{i,mis}, R_i | \Theta, \psi) = f(Y_{i,obs}, Y_{i,mis} | \Theta) f(R_i | Y_{i,obs}, Y_{i,mis}, \psi). \quad (1)$$

The first part is the marginal density of the response process and the second part is the density of the missingness process, conditional on the response. The PMM can be seen as a mixture of different populations, characterized by the observed pattern of missingness, it is based on the below factorization:

$$f(Y_{i,obs}, Y_{i,mis}, R_i | \Theta, \psi) = f(Y_{i,obs}, Y_{i,mis} | R_i, \Theta) f(R_i | \psi). \quad (2)$$

The SPM assumes that the response process Y_i and the missingness process R_i are conditionally independent of each other by sharing a random effect b_i . Therefore, the density function can be described as

$$f(Y_{i,obs}, Y_{i,mis}, R_i | \Theta, \psi) = \int f(Y_{i,obs}, Y_{i,mis} | b_i, \Theta) f(R_i | b_i, \psi) f(b_i) db_i. \quad (3)$$

As mentioned above, these models are introduced as sensitivity analysis assuming MNAR (i.e., assuming all missing data are MNAR alternatively). However, in reality, the composition of missing data may be more complex, for example, a mixture of MAR and MNAR. A relatively large number of empirical studies have examined the performance of the classic models. These studies suggest that those models can reduce or eliminate bias when their assumptions are met. However, in many realistic scenarios, the models can produce estimates that are even worse than those of MAR-based missing data-handling methods (Enders, 2010). Those models heavily rely on the underlying assumption; unfortunately, these assumptions are largely untestable, so there is no practical way to judge the model's performance in a real data analysis.

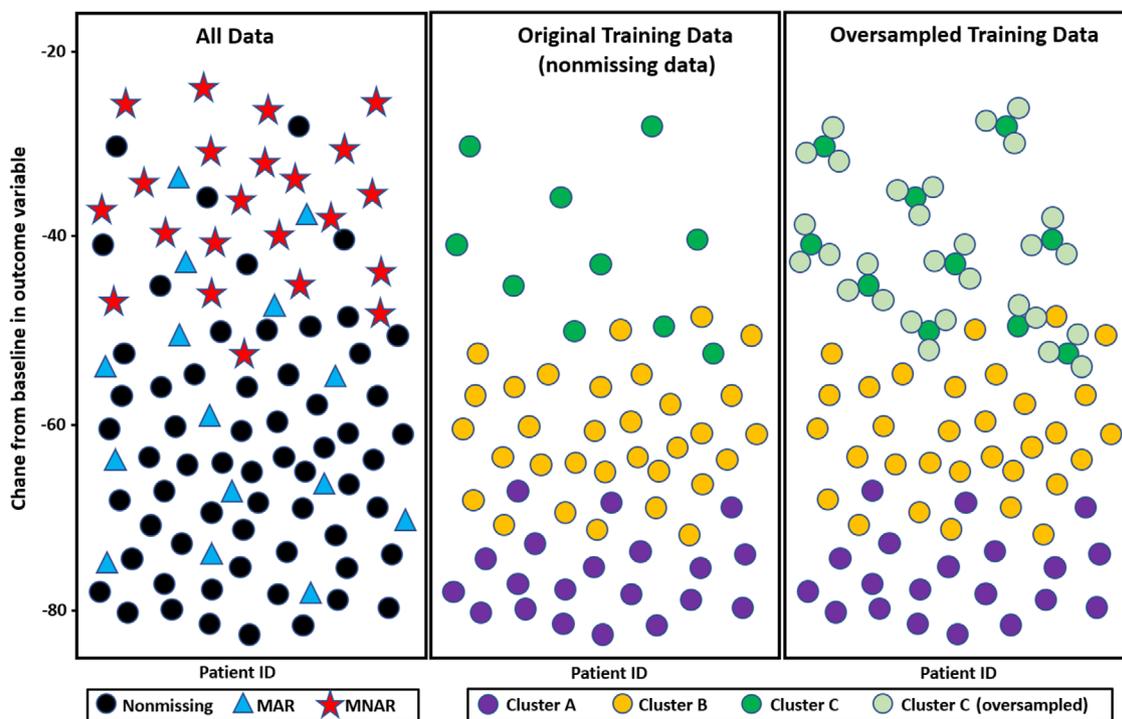


FIGURE 1 Demonstration of clustering and oversampling: cross-sectional of the longitudinal profiles – “change from baseline in outcome variable”

In this paper, with an aim of handling more realistic missing data scenarios (e.g., mixture of MAR and MNAR), a machine learning based missing data prediction framework is developed and evaluated using simulation data and real clinical trial data. According to Breiman (2001), in statistical learning “the goal is not interpretability, but accurate information.” In line with this thinking, the proposed framework handles MNAR by emphasizing what is missing, while it also covers MAR by seeking for accurate individual information. The MNAR problem is treated as an imbalanced learning task, that is, the minority class oversampling (Weiss, 2013) is used to compensate for the MNAR data. To implement oversampling in longitudinal continuous data, clustering via k -mean trajectories (Gower, 1971) is performed first. Patients are clustered into “good responder,” “medium responder,” and “low responder,” according to their longitudinal efficacy profiles. In our simulation study, to simplify the problem, we assume the dropouts due to lack of efficacy are MNAR. Therefore, those patients are mostly in the “low responders” class (of course it could be the other way round in reality, i.e., extreme good responders discontinue from the trial as they consider no need to continue the treatment). Depending on the proportion of MNAR data, the size of available data in the worst cluster can be smaller than the size of the other clusters (i.e., imbalanced distribution). In order to compensate for the MNAR in that cluster, and also to avoid the individual prediction being driven by the available data from the completers to an overall average level, random oversampling (with replacement) in the minority class is necessary. See Figure 1 for a display of the distribution of “change from baseline in an outcome variable” (cross-sectional of the longitudinal profiles) in a simulated example. The full data (including the nonmissing data and the values that are set to “missing” in the simulation) are presented in the left panel, black dot = nonmissing data, red star = MNAR, and blue triangle = MAR. The nonmissing data (i.e., the original training data) are repeated in the middle panel and clustered into three groups: green = “low responders,” orange = “medium responder,” and purple = “good responder.” It is clear that the nonmissing data is not a good representative of the full data considering the MNAR data. The “low responders” (green dots) are relatively rare in the original training data (i.e., the minority class). When random oversampling is applied in the minority class, as shown in the right panel (the light green dots are the oversampled cases), there are more data points in the “low responders” area to compensate for the MNAR, that is, the distribution of green and light-green dots approximates the distribution of green and red dots.

One may question the oversampling applied here, arguing that the distribution of data has been changed due to the oversampling, and this will impact the treatment effect estimation. First, the oversampled dataset will never be used to estimate the treatment effect. All of the efforts taken here are to minimize the individual error in statistical learning. Once

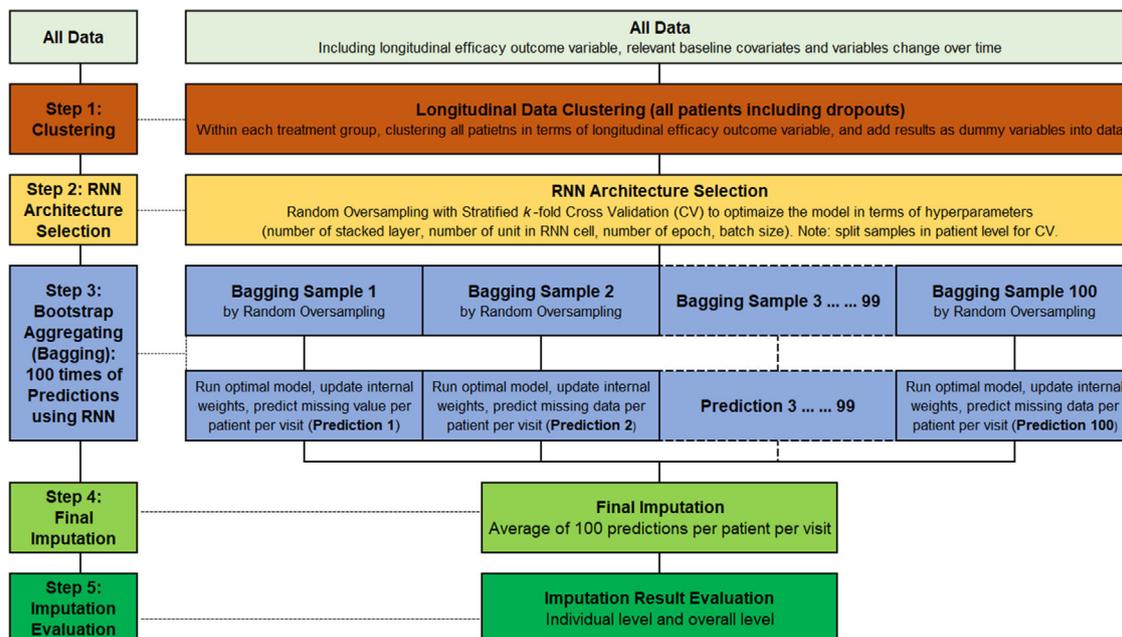


FIGURE 2 Overall workflow diagram: process to handle realistic missing data scenario in a longitudinal setting. Step 1: Structure the longitudinal efficacy response data by clustering; Step 2: select the optimal RNN architecture facilitated by oversampling and stratified K-fold CV; Step 3: perform multiple predictions for individual patients with missing data (facilitated by oversampling to avoid bias introduced by missing mechanisms); Step 4: create the final imputed dataset; Step 5: evaluate the imputation result

the prediction for each patient is optimized by minimizing the individual error, the treatment effect will be estimated based on the observed data plus the imputed data (i.e., no oversampling in the data analysis part).

We use recurrent neural networks (RNN) to model the longitudinal data. RNN is a type of neural network that can learn from the past to predict the future outcomes (Rumelhart et al., 1986; Schmidhuber, 1993). This allows us to exhibit temporal dynamic behavior for time sequence data, and thus it is a powerful tool for longitudinal clinical data modeling. RNN provides flexible nonlinear modeling without requiring any or much domain knowledge about the interrelationship between the variables, it learns automatically from the training data to estimate the weights and then predict the new data. Different RNN architectures are experimented to tune various hyperparameters, and the optimal model is selected via the bias-variance trade-off approach (Claesen & De Moor, 2015). To improve the accuracy of prediction and also to consider the uncertainty of a single prediction, bootstrap aggregating (bagging) is implemented. In light of the “evidence-based computational statistics” (Boulesteix et al., 2017, 2018), the proposed method is evaluated in the practically relevant simulation data and exemplified in a real clinical trial data. In the simulation data, the real-life plausibility of the simulation scenarios is emphasized. A mixed missing mechanism of MAR and MNAR is considered in the simulation. The real dataset is from an antidepressant clinical trial, which is one of the few publicly available datasets that can be used to demonstrate methods for handling missing data where a continuous outcome is measured repeatedly. The imputation results are evaluated at the individual patient level and the overall population level. Overall, the proposed methods provided plausible individual prediction for both of the MAR and MNAR data and reduced the bias of missing data in treatment effect estimation. Therefore, this paper offers an opportunity to encourage the integration of machine learning strategies for handling of missing data in the analysis of randomized clinical trials.

2 | METHODS

We propose a computational approach in this paper which comprises various individual components (see Figure 2 for the overall workflow). The proposed framework handles MNAR by emphasizing what is missing, while it also covers MAR by seeking for accurate individual information. The MNAR problem is treated as an imbalanced learning task, that is, the minority class oversampling is used to compensate for the MNAR data.

The first step is to cluster all patients (including the dropouts) according to their longitudinal efficacy profiles. Clustering structures the longitudinal data within each treatment group. It is a key concept in the proposed approach due to the following reasons: (i) the clusters (clustering results) are used in the stratified k -fold cross-validation (CV) and the random oversampling step as the “categorization” of the continuous target variable to balance the majority and minority clusters; (ii) the clusters are used (as dummy variable) in the RNN model to indicate the longitudinal pattern of patient efficacy profiles, which is important to borrow information from the similar patients (seeking for accurate individual information). Technical details about clustering are provided in Section 2.1.

The second step is to select the optimal RNN architecture. We use RNN to model the complex longitudinal data in a nonparametric manner (details about RNN are provided in Section 2.2). The RNN architectures feature a set of hyperparameters (e.g., number of units in each RNN cell, number of stacked layers, batch size, and number of epochs) that must be determined before training commences. In this step, the optimal RNN architecture (which can learn adequately from the training data and also performs equally well in the validation data) is selected via a bias-variance trade-off approach. Considering the data distribution with the presence of MNAR (as discussed in Section 1), stratified k -fold CV and oversampling of minority class are implemented in the RNN architecture selection process. Details are provided in Sections 2.3 and 2.4.

The third step is to generate, say, 100 bootstrap aggregating (bagging) samples with the minority classes oversampled, and to predict the missing data in each bagging sample. An ensemble method (i.e. bagging) is used to improve the prediction accuracy and also to consider the uncertainty of a single prediction. The optimal RNN model (in terms of hyperparameters) is executed 100 times, each time updating the internal weights of the RNN model and providing predictions for the missing data. In practice, the number of bagging can be even higher. We use 100 as a reasonable number in this paper as the proposed method is time consuming and computationally intensive. Within each bagging, minority classes are oversampled to compensate for the MNAR data.

The fourth step is to average all 100 predicted values for each patient at each visit. These are considered as final imputation.

The fifth step is to evaluate the imputation results at an individual patient level by visualizing the efficacy profiles. The treatment effect is estimated from the imputed data by commonly used statistical analysis methods.

The treatment effect estimated from the imputed data using the commonly used methods is compared with the treatment effect that was estimated using different methods including commonly used methods and the classic models (like SM, PMM, and SPM) that are applied without missing data imputation.

2.1 | Longitudinal data clustering

The k -mean clustering aims to partition n observations into k clusters in which each observation belongs to the cluster with the nearest mean. This is done by alternating an expectation phase and a maximization phase. In the expectation phase, the center of each cluster is determined, then in the maximization phase, each observation is assigned to its nearest cluster. This process is repeated until no changes in the clusters occur. For k -mean trajectories, different types of distance can be calculated. The R package `kml` is used in this paper (Genolini & Falissard, 2011). Considering the missing data, the classic Euclidian distance with Gower adjustment (Gower, 1971) is used. Hence the dropouts are also clustered based on their available data. Consider a set of n patients. The target variable is measured for each patient up to time t . Let Y_i denotes the patient i , and let Y_{ik} denotes the measurement for patient i at time k . The difference of the trajectories between two patients i and j can be calculated using the classic Euclidian distance with the Gower adjustment:

$$Dist_{GA}^E(Y_i, Y_j) = \sqrt{\frac{t}{\sum_{k=1}^t (\omega_{ijk})} \sum_{k=1}^t (Y_{ik} - Y_{jk})^2 \omega_{ijk}} \quad (4)$$

Here, ω_{ijk} equals 0 if Y_{ik} or Y_{jk} are missing, and 1 otherwise. Assuming the distribution of the target variable is different between the treatment group, we perform the clustering within each treatment group. The number of cluster should be decided on a case-by-case basis and should be prespecified. In this paper, we set the number of clusters to three within each treatment group with the idea of splitting the patients into “good response,” “medium response,” and “low response” categories according to their efficacy profile.

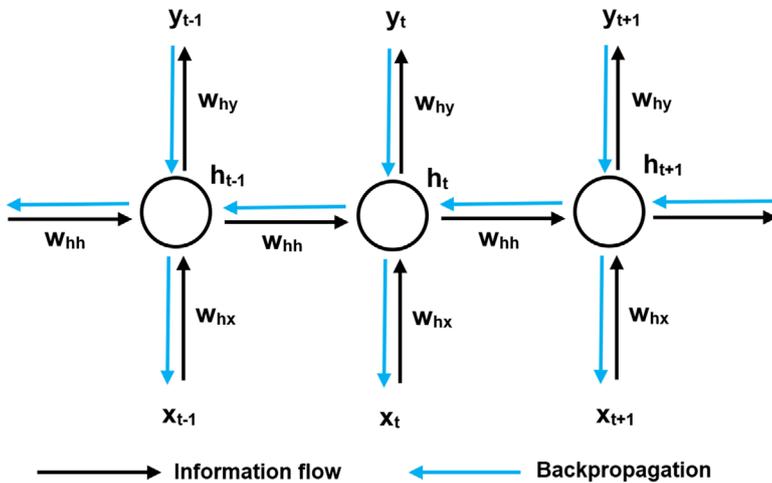


FIGURE 3 An RNN in time of the computation (modified based on LeCun et al., 2015)

2.2 | Recurrent neural network

An RNN is a type of neural network that can learn from the past to predict future outcomes. A basic RNN (Rumelhart et al., 1986; Schmidhuber, 1993) is shown in Figure 3. It uses hidden states (which temporarily store information about the past) to transfer information through time. At time t , the weighted sum of the input information x_t (e.g., the variables change over time) and the previous hidden state (h_{t-1}) is processed using an activation function (e.g., the weighted sum is squashed between -1 and 1 using a hyperbolic tangent (Tanh) function (see details of the Tanh function in Appendix Figure A.1). Then the processed information (hidden state h_t) is carried forward to the next time step. Meanwhile, the hidden state h_t can be output as prediction result y_t for time t using an appropriate activation function, for example, a Rectified Linear Unit (ReLU function; see Figure A.1) for a continuous outcome variable. In this way, an RNN can map an input sequence with elements x_t into an output sequence with elements y_t , with each y_t depending on all the previous $x_{t'}$ (for $t' \leq t$). The same internal weights (matrices w_{hx} , w_{hh} , w_{hy}) are used at each time step. For the internal weight optimization, a backpropagation algorithm can be applied to the computational graph of the unfolded network from the right to the left, that is, to compute the derivative of the error with respect to all the internal weights. The error, also called “loss” in machine learning, is calculated as the predicted value minus the observed value for a continuous outcome variable.

In addition to the basic RNN unit mentioned above (which contains a Tanh activation function), there are other types of RNN units. For example, the long short-term memory units (LSTM) (Gers et al., 1999; Hochreiter & Schmidhuber, 1997) and the gated recurrent unit (GRU) (Cho et al., 2014) are the most commonly used ones. Empirical evaluations show that LSTM and GRU perform superior over the basic RNN (Chung et al., 2014; Shewalkar et al., 2019). LSTM performs slightly better than GRU in terms of prediction accuracy (Shewalkar et al., 2019). We provide a basic introduction to LSTM in Appendix A.1. In this paper, LSTM is implemented using the Keras library version 2.2.4 (Falbel et al., 2015) in Python (version 3.6).

Neural networks use stochastic gradient descent which is an iterative method for optimizing an objective function with suitable smoothness properties. It can be regarded as a stochastic approximation of gradient descent optimization since it replaces the actual gradient (calculated from the entire dataset) with an estimate thereof (calculated from a randomly selected subset of the data, called “batch”). The machine learning algorithms consider the problem of minimizing an objective function that has the form of a sum:

$$Q(w) = \frac{1}{n} \sum_{i=1}^n Q_i(w), \quad (5)$$

where the internal weight w that minimizes $Q(w)$ is to be estimated. Each summand function Q_i is typically associated with the i th observation in the training dataset. When used to minimize the above function, a batch gradient descent method would perform the following iterations:

$$w^* = w - \eta \left(\frac{\partial \text{Loss}}{\partial w} \right), \quad (6)$$

where w^* is the new weight, w is the old weight, η is a step size (also called “learning rate”), the last part of this equation is the derivative of loss with respect to the weights. The `Keras` library implements the adaptive moment estimation (Adam) optimizer (Kingma & Ba, 2014). Empirical results demonstrate that Adam works well in practice and compares favorably to other stochastic optimization methods (Kingma & Ba, 2014). In Adam, the learning rate is initialized (e.g., default initial value = 0.001 in `Keras`) and then adapted automatically in training iterations.

There are three types of input variables in an RNN. (i) *Initial state*: the hidden state (h_0) at time step t_0 . In practice, the default approach is to set the initial state as zero. However, if the impact of the initial state is not negligible, it makes sense to train the initial state as a variable. Thereby the model can start to learn from a good default state. In a clinical study, for the continuous outcome variable, the baseline value of the outcome variable can be considered as the initial state in an RNN. (ii) *Series input*: the variables change over time (x_t). (iii) *Static input*: for example, relevant demographics and baseline characteristics. In the `Keras` library, static input can be implemented by passing external constants to the RNN, there are no internal model weights learnt for static inputs. For all types of input data, the continuous variables need to be standardized before feeding into RNN to make the calculation faster.

When training the RNN, the nonmissing data from the dropouts should also be used as this particular part of the data is quite important to learn a certain pattern of those dropouts. This may be more important for the dropouts due to lack of efficacy (MNAR) as the trend of their efficacy profiles can be very different from the patients who completed the study. If possible, variables to indicate the missing mechanism (whether MNAR or not) need to be included in the model. Including the dropouts in the model will lead to different lengths of time sequence in data. Using an RNN, a fixed length of time series input is expected in the current available deep learning packages. An effective way to handle this problem is to use sample weight, that is, to create a metric per patient per time to indicate which time points to use in the learning. For example, the weights are set to 1 for nonmissing time steps, and 0 for missing time steps. These metrics are then multiplied by the loss (e.g., the difference between the predicted value and the actual value) per patient and per time before training the RNN. For example, for patient j at time t , the final loss is

$$loss_{jt} = loss'_{jt} w_{jt}, \quad (7)$$

where $loss'_{jt}$ is the loss calculated before using the sample weight metrics, w_{jt} is the sample weight for patient j at time t . By having such sample weight, the missing time steps will be ignored in learning.

2.3 | Minority class oversampling

As mentioned in Sections 1 and 2.1, we cluster patients into “good responder,” “medium responder,” and “low responder” according to their longitudinal efficacy profiles within each treatment group. Depending on the proportion of MNAR data, the size of available data in the worst cluster can be smaller than in the other clusters (i.e., the worst cluster is the minority class; see Section 1 for more details). To compensate for the MNAR in that cluster, and also to avoid the individual prediction being driven by the available data from the completers to an overall average level, random oversampling (with replacement) in the minority class is necessary. This process involves randomly selecting examples from the minority class, with replacement, and adding them to the training dataset. The amount of oversampling is a hyperparameter of the system (Chawla et al., 2002), it should be decided on a case-by-case basis. We describe two slightly different oversampling approaches in Section 3.3.2 (for the simulation studies) and Section 4 (for the real data implementation). Oversampling is implemented in both Step 2 (RNN architecture selection) and Step 3 (bootstrap aggregating) of the proposed framework to compensate for the MNAR data in both model selection and individual prediction processes.

2.4 | Stratified k -fold CV

In machine learning, hyperparameters are the parameters whose values are used to control the learning process (this is different from the model internal parameter or weight whose values need to be optimized during the training process). In general, for RNN, hyperparameter may include learning rate (see details in Section 2.2), number of epoch (defined as the number of times that the learning algorithm work through the entire training dataset), and batch size (defined as the number of samples to work through before updating the internal weights of the model). In addition to those general hyperparameters, the number of units in each RNN cell and number of stacked RNN layers are also needed to be specified before

training commences (as a more complex model may lead to overfitting). The choice of hyperparameters can significantly affect the resulting model's performance; hence, a disciplined, theoretically sound search strategy is essential (Claesen & De Moor, 2015). The bias-variance trade-off is the most commonly used approach for hyperparameter tuning with a goal of selecting a model that can learn adequately from the training data and also performs equally well in validation data.

The k -fold CV is the standard tool for hyperparameter tuning to address the overfitting/underfitting problem. In k -fold CV, the original sample is randomly split into k approximately equal-sized subsets. Of the k subsets, a single subset is retained as the validation data for testing the model and the remaining $k - 1$ subsets are used as training data. The CV process is then repeated k times, with each of the k subsets used exactly once as the validation data. Ideally, all possible combinations of hyperparameters should be experimented using the CV approach. For each scenario, the loss over iteration history should be calculated and visualized (to facilitate the comparison) for both training data and validation data. The value of the hyperparameter that provides the least loss for both training data and validation data should be determined as the optimal value. Although the partition of the k fold is performed randomly, it does not guarantee to have a balanced distribution of the target outcome variable in each fold without supervision (especially with the presence of MNAR in the data). Stratified k -fold CV seeks to ensure that each fold is a representative subset of the whole data in terms of the variable of interest. This is very important for the MNAR as by stratification the dropouts due to lack of efficacy will be included in each fold equally, which means in each time when repeating the training process, the certain patterns of target variable in those dropouts (their nonmissing part) will be learnt adequately. Whenever the joint application of CV and oversampling concerns, the "overoptimism" issue should be emphasized and distinguished from overfitting (Santos et al., 2018). If the entire original data is oversampled first and then CV is performed later on, the same samples may appear in both of the training and validation partitions, thereby the model performs "very well" in both partitions. This is known as the "overoptimism" issue. Therefore, a better approach would be that the dataset is first divided into k stratified partitions and the oversampling happens in the training data part only. The validation data are never oversampled or seen by the model in the training stage, thereby allowing a proper evaluation of the model's performance for the generalization purpose.

3 | SIMULATION STUDY TO EVALUATE PERFORMANCE OF METHODS

We evaluate the proposed method by means of an extensive simulation study. In designing the simulation study, we used realistic missing data scenarios. We consider a longitudinal continuous clinical score as the efficacy variable. Further, we consider a mixed missing mechanism of MAR and MNAR in the simulation. One of the advantages of the simulation study in this context is that the "missing values" are known (as the complete data are generated first and some values are set to "missing"), and this can be used as a benchmark to evaluate the performance of the imputation methods (both at the individual level and the overall level).

3.1 | Design of simulation study

The patient baseline characteristics and longitudinal efficacy data are simulated assuming a two-arm parallel designed clinical trial. Each patient is designed to be treated and assessed biweekly from Week 0 (baseline) up to Week 16 ("primary endpoint"). Different sample sizes (300, 400, or 500 patients), overall dropout rate (20%, 30%, or 40%) and monotone missingness starting time-point (Week 8, Week 10, or Week 12) are considered in the simulation. In total, $3 \times 3 \times 3 = 27$ scenarios are simulated. In each scenario, the patient is randomly assigned to the treatment groups (Test : Control = 1 : 1). A longitudinal clinical score is considered as efficacy variable which decreases over time in general. To take account of the inpatient correlation, the change from previous visit values in the score (per patient per visit) is modeled considering several fixed factors and a random effect. Similar idea as for the PMM (i.e., data patterns are different for MNAR and completers), a mixed missing mechanism of MAR (discontinue due to lost to follow-up), and MNAR (discontinue due to lack of efficacy) is considered in the simulation. The control group is more impacted by the missing data as the proportions of MNAR and MAR are much higher in the control group than in the test group. The MAR is influenced by the treatment group but not by any other covariates (i.e., within each treatment group, it is actually an MCAR scenario). The complete simulation data before setting the missing values are kept for all patients at all visits for the purpose of imputation result evaluation. The details about the data generation process are described in Appendix A.2.

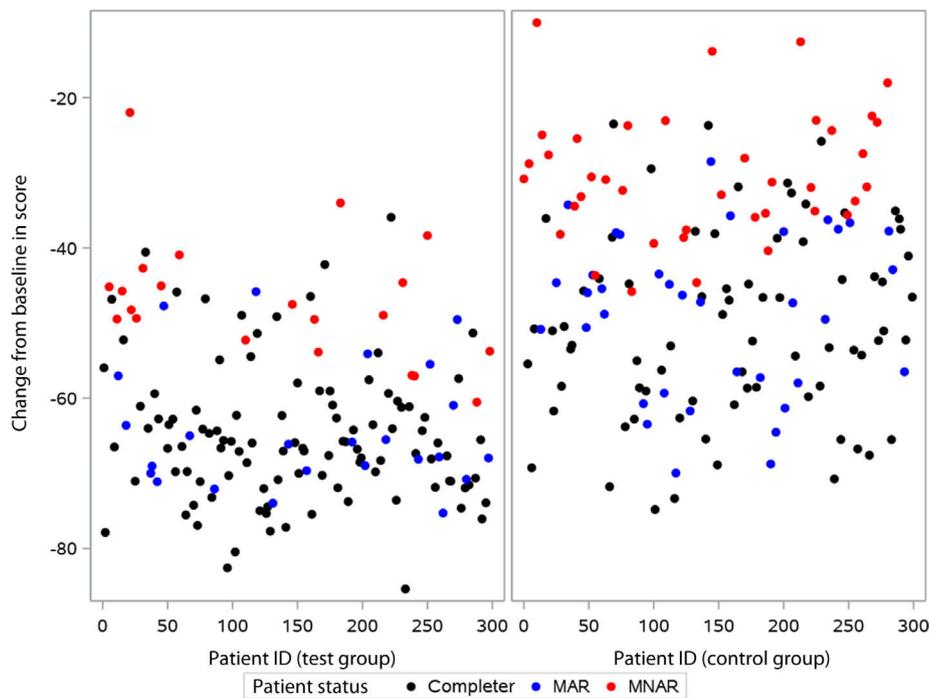


FIGURE 4 Simulation Scenario 300-40-10: scatter plot for “change from baseline in score at Week 16.” The black dots are the completers, the blue dots are MAR, and the red dots are MNAR

In this paper, we report the simulation scenario with total sample size = 300, overall dropout rate = 40%, missingness starting at Week 10 (called “Scenario 300-40-10” for short) as an example, as it is one of the scenarios that is most impacted by MNAR. In the test group, 21(14.0%) and 24(16.0%) patients discontinued due to “lack of efficacy” and “lost to follow-up,” respectively. In the control group, 39(26.0%) and 36(24.0%) patients discontinued due to “lack of efficacy” and “lost to follow-up,” respectively. As a cross section of the longitudinal profiles, the distribution of “the primary endpoint: change from baseline in score at Week 16” is shown in Figure 4. It is clear that the MAR (blue dots) are randomly distributed over the whole data space, but the MNAR data (red dots) are mostly presented in the “low response” area. Therefore, realistic missing data scenarios are successfully “mimicked” in the simulation study. A proper missing data handling method should compensate for the MNAR data and also provide accurate prediction to the MAR data.

In addition, we also simulated the scenarios with only MNAR data, using data from the 27 scenarios mentioned above, but with the MAR values replaced by the known true values (i.e., only MNAR are present in the data). An expected result of these simulations is that the classic models (i.e. PMM, SM, and SPM) perform the best and the proposed approach also performs fairly well.

3.2 | Measuring performance of the proposed methods

To measure the performance of the methods at an overall population level, the treatment effect is estimated using different methods (as described below), and the results are compared in a forest plot. The treatment effect is defined as the difference of change from baseline in score at Week 16 between the treatment group.

Analysis methods

- (i) A mixed model for repeat measurement (MMRM) was used for longitudinal data from Week 2 to Week 16. MMRM included treatment, visit, and treatment×visit as fixed effects, baseline score value as covariate, and subject as the random effect.
- (ii) Analysis of covariance (ANCOVA) model used for the data at Week 16 only, ANCOVA included treatment as factor and baseline score value as covariate.
- (iii) Classic models including PMM, SM, and SPM were used for longitudinal data from Week 2 to Week 16. Details about PMM, SM, and SPM are provided in the [Supporting Information](#) (with the SAS code and instructions).

Analysis dataset and the corresponding analysis method:

- (i) The “true” treatment effect: the simulated complete efficacy data before setting any missing value are analyzed by MMRM and ANCOVA. The “true” treatment effect is used as a benchmark to evaluate the performance of the proposed method and other methods.
- (ii) The imputed data (i.e., nonmissing data + data imputation by RNN prediction facilitated with clustering and oversampling) is analyzed by MMRM and ANCOVA models.
- (iii) To illustrate the role of clustering and oversampling in handling of MNAR data, the imputed data (i.e., nonmissing data + data imputation by a straightforward RNN prediction without facilitation with clustering and oversampling) is also analyzed by MMRM and ANCOVA models.
- (iv) The nonimputed data are analyzed by following commonly used methods: MMRM, ANCOVA, PMM, SM, and SPM.

In addition to the overall level treatment effect comparison, we are also interested in the prediction performance at the individual level. The patient profiles (the observed values and the 100 predicted values) are visualized for all dropouts. The mean value (final imputation) and variability of prediction (measured as the first and third quartiles) at each study week are provided in the plot. The “true” values of missing data (i.e., the simulated complete efficacy data before setting any missing value) are also provided for each patient to visually check the accuracy of the imputation.

3.3 | Simulation results

3.3.1 | Results of longitudinal clustering

As mentioned in Section 2, clustering is the first and very important step. Within each treatment group, patients are clustered into three categories: “good responder,” “medium responder,” and “low responder” according to their longitudinal efficacy profiles. Figure A.3 in the Appendix shows the individual patient profiles by cluster for Scenario 300-40-10. It is clear that the k -mean clustering captures the longitudinal data structure very well for all patients (including the dropouts). It must be noted that the interpretation of clusters has to be taken with caution as k -mean clustering is unsupervised learning, that is, the clusters/categories are not labeled in the input data. However, the clustering is helpful to structure the longitudinal profile patterns and to seek for similar patients. In addition, as mentioned in Section 1, to learn the pattern of minority clusters adequately in such setting (i.e., the worst clusters in each group with less completers compared to other clusters) and also to avoid the prediction that has been driven by the majority available data from the completers, it is necessary to oversample the small clusters with less nonmissing data.

3.3.2 | Tuning the RNN hyperparameters

Different RNN hyperparameters are tuned using stratified k -fold CV and oversampling (as described in Section 2.4). As mentioned above, it is necessary to oversample the small clusters with less available data. For the hyperparameter tuning purpose, to have a consistent oversampling approach that can be generalized in all simulation scenarios, the following rules are used: the “good” and “medium responder” clusters, and the completers from the “low responder” cluster are oversampled to the size of the largest cluster in that treatment group; in addition, to utilize the available data (nonmissing part) from the dropouts, the dropouts from the worst cluster are also 1:1 random sampled (with replacement). During the CV process, the loss (measured as mean squared error (MSE)) changing over iteration history from the fivefold CV and the MSE at the last training and validation iteration is compared and the appropriate values for the hyperparameters are selected. Based on a substantial number of experiments, the optimal model for the simulation data is determined as LSTM with one single layer and nine units in each cell, iteration epoch of 3000, and batch size of 50.

3.3.3 | Imputation result evaluation at the individual level

The optimal RNN model is performed 100 times. At each time, the data (including the available data from the dropouts) are randomly sampled (minority clusters are oversampled following the same approach as mentioned in Section 3.2.2), the

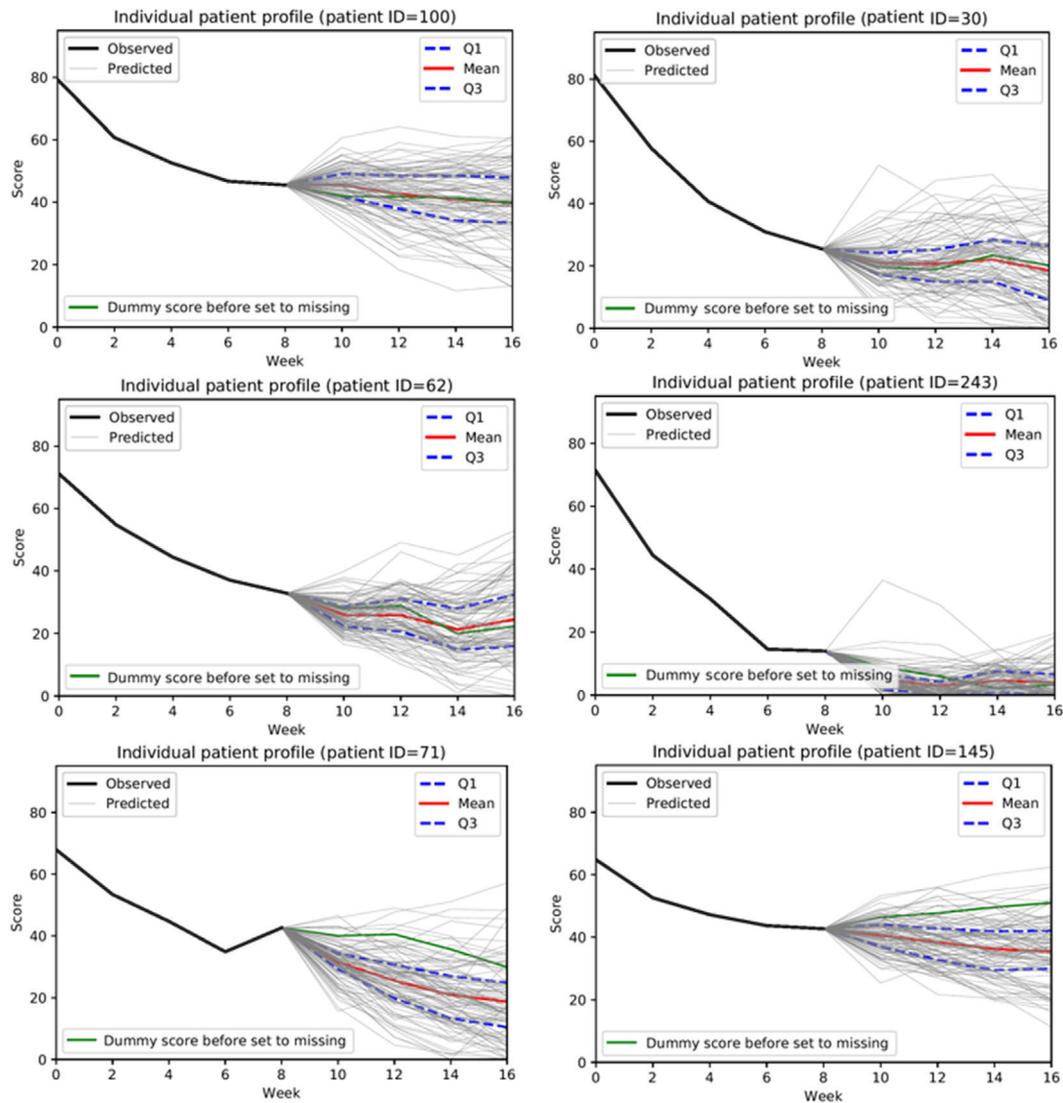


FIGURE 5 Simulation data: examples of individual prediction compared with actual values. The solid black lines are the observed data, each dashed gray line is the prediction from each bagging, and the mean (red lines) and quartiles (blue lines) from the 100 predictions are also provided. The solid green lines are the complete data before setting the missing values (i.e., the true known values)

internal weights are updated and predictions for the missing data are provided within each bagging. For each patient, at each time point, the average of all 100 predicted values is considered as the final imputation. Some examples of the individual patient profile are provided in Figure 5. For the majority of the dropouts, the predictions are close to the actual values (e.g., the first four patients in the figure). The proposed methods provided fairly good predictions for “good,” “medium,” and “low responders.” For some dropouts (< 15% of dropouts), when the inpatient variability is large or the scores are extremely high at the end of the trial, the imputations are not good as expected (e.g., patient numbers 71 and 145 in Figure 5). Another reason for such bad prediction could be a lack of relevant predictors in the data. Due to the difficulties in the longitudinal data simulation, only a few relevant variables are generated in the simulation data (see details of data generation process in Appendix A.2).

3.3.4 | Imputation result evaluation at the overall population level

The results from different methods (as described in Section 3.2) are presented in a forest plot (Figure 6) for Scenario 300-40-10. The proposed imputation method (i.e., RNN imputation facilitated by clustering and oversampling) + standard analysis method (MMRM or ANCOVA) provided the best estimation of the true treatment effect (i.e., the estimates are the closest

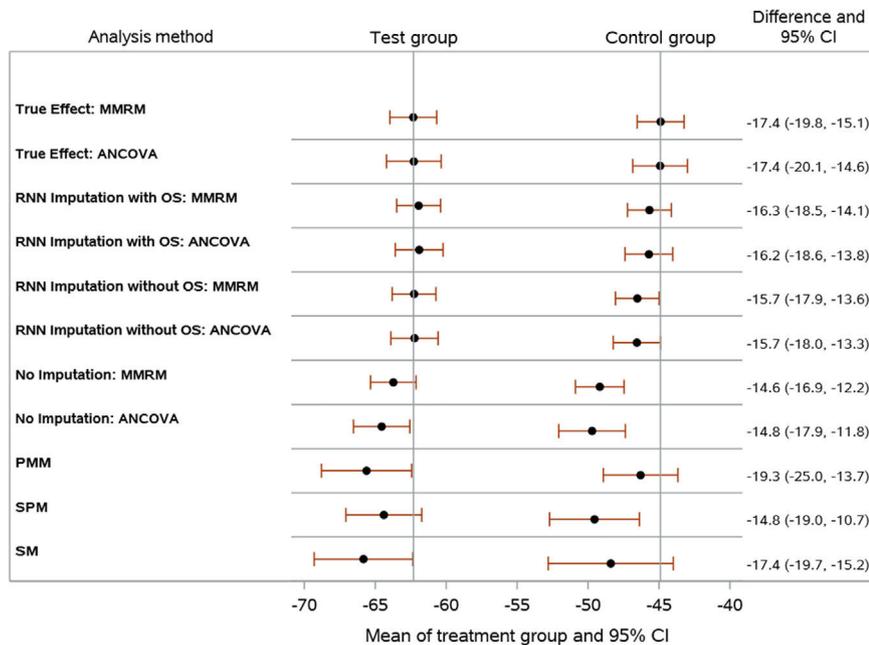


FIGURE 6 Forest plot for analysis results of “change from baseline in score at Week 16” using different methods. True effect: estimation from the complete efficacy data before setting any missing value. RNN imputation: nonmissing data + imputed data by the proposed method. OS = oversampling. No imputation: nonmissing data only. PMM, pattern mixture model; SPM, shared parameter model; SM, selection model

to the true effect in both treatment groups). RNN imputation without clustering and oversampling also provided good estimation (only with a slight bias in the control group), hinted that if the impact of MNAR is considered as ignorable, a simply RNN imputation (without clustering and oversampling) can also provide accurate imputation. Since MMRM and ANCOVA assume MAR, therefore, it is not surprising to have a considerable bias when the MMRM and ANCOVA are applied without imputing the missing value (given the presence of both MNAR and MAR in the data). The bias is larger in the control group where the impact of the MNAR data is much heavier. There is a systematic bias in the results from the classic models in both treatment groups (i.e., the models including PMM, SPM, and SM that are applied without imputing the missing value). Similar to MMRM and ANCOVA that are applied without data imputation, these models tend to overestimate the treatment effect when both MAR and MNAR are present in the data. Since those models heavily rely on the underlying assumption (e.g., a pure MNAR), when a mixture of MAR and MNAR is present in the data, the impact of missing data is somehow not properly handled by these models. The proposed method performed equally well in all 27 simulated scenarios. Similar patterns (as discussed above) are observed in all other simulation scenarios. The analysis results for other 26 scenarios are provided in the [Supporting Information](#). In general, based on the simulation study, the impact of missing data on the treatment effect estimation in a realistic scenario (e.g., mixture of MAR and MNAR) is properly handled by the proposed method by providing accurate individual predictions for both MAR and MNAR data.

In addition, scenarios with only MNAR data are also simulated, see the example of Scenario 300-40-10 in Figure 7 (in this case, the actual proportion of missing data is 20% due to the absence of MAR). In general, as expected, the classical models (i.e., PMM, SM, and SPM) perform better than MMRM and ANCOVA when only MNAR is present in the data, and the proposed approach also performs fairly well.

4 | REAL DATA EXAMPLE

The proposed method is implemented in a real dataset from an antidepressant clinical trial. Original data are from an antidepressant clinical trial with four treatments; two doses of an experimental medication, a positive control, and a placebo (Goldstein et al., 2004). Hamilton 17-item rating scale for depression (HAMDI7) is observed at baseline and weeks 1, 2, 4, 6, and 8. To mask the real data, Week 8 observations are removed. Two arms are created: the control group (original placebo arm, $N = 88$) and a test group created by randomly selecting patients from the three nonplacebo arms ($N = 84$). There are 21(25.0%) and 23(26.1%) dropouts in test group and control group, respectively. Within each treatment group, patients are clustered into three subgroups in terms of their HAMDI7 score profile. The individual patient profiles by cluster are provided in Figure 8. It is clear that the “good responder” and “low responders” are relatively small clusters. Small clusters are randomly oversampled to the size of the largest cluster (“medium responder”) within each treatment group. The oversampling process in the real datasetting is much simpler than the ones used in the simulation data given

FIGURE 7 Forest plot for analysis results of “change from baseline in score at Week 16” using different methods in the scenario with MNAR only in the data. True effect: estimation from the complete efficacy data before setting any missing value. RNN imputation: nonmissing data + imputed data by the proposed method. OS, oversampling; No imputation: nonmissing data only; PMM, pattern mixture model; SPM, shared parameter model; SM, selection model

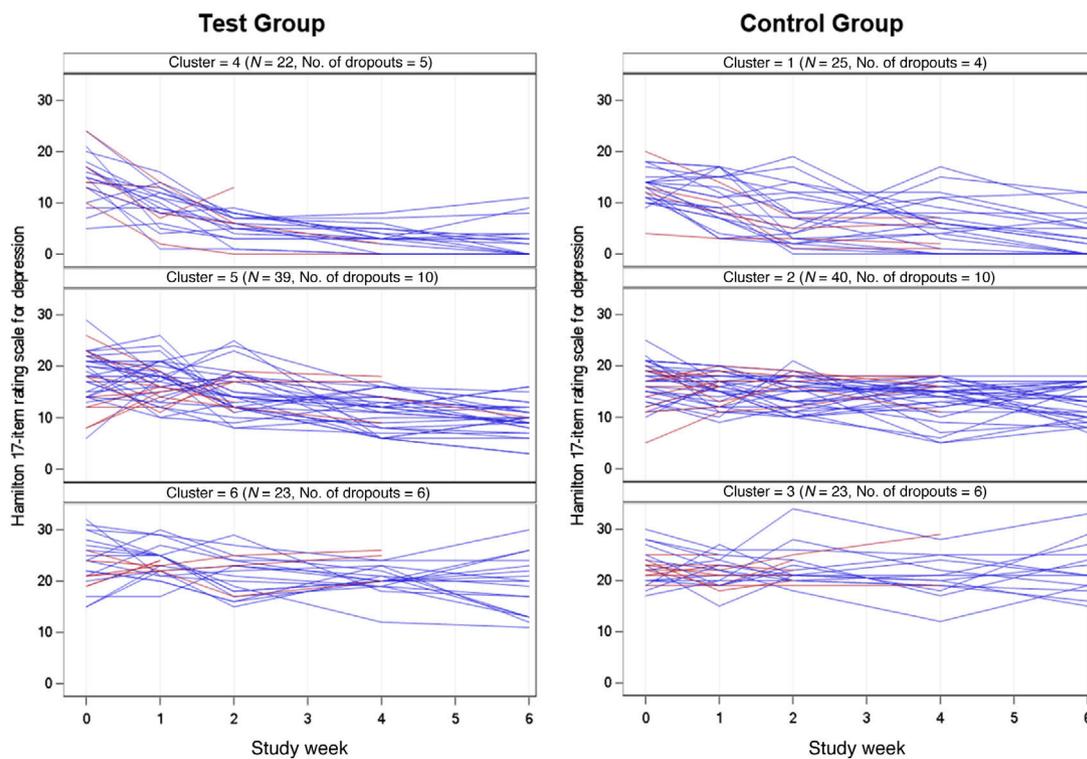
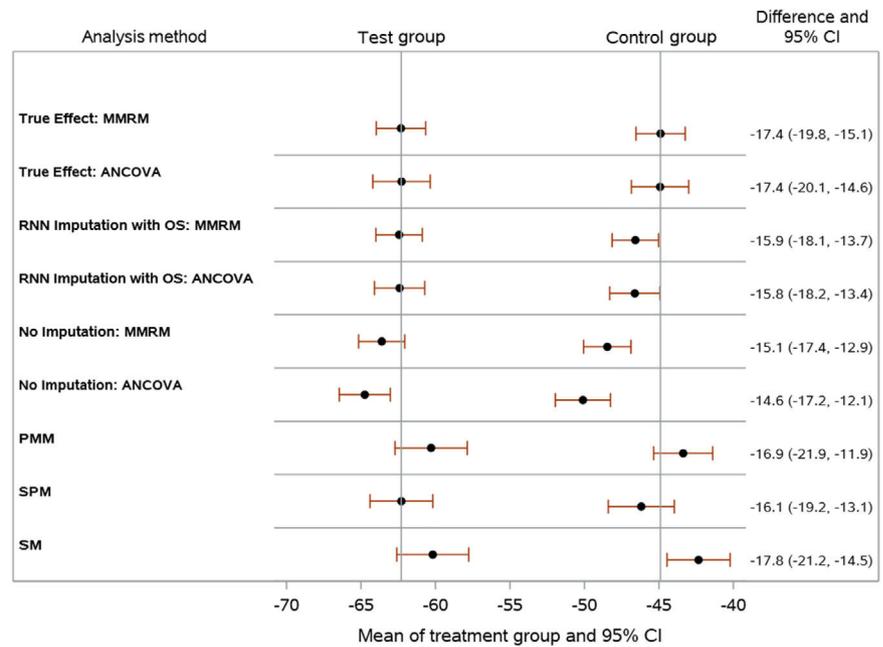


FIGURE 8 Real data: patient profile by the treatment group and cluster. Blue lines are for completers, and red lines are for dropouts

the nature of the data. The dropout reasons are not available in the published real dataset; this makes it difficult to make an assumption about the missing mechanism. Based on the clustering results, the dropouts in “low” and “good” responder clusters could be considered as MNAR, as they may discontinue from the trial due to lack of efficacy or they responded so well before completing the trial and considered that no need to continue with the treatment. Especially the dropouts in the “low responder” clusters (6/23 in the test group vs. 9/23 in the control group), without proper handling of the missing data, their impact on the treatment effect estimation can be nonnegligible.

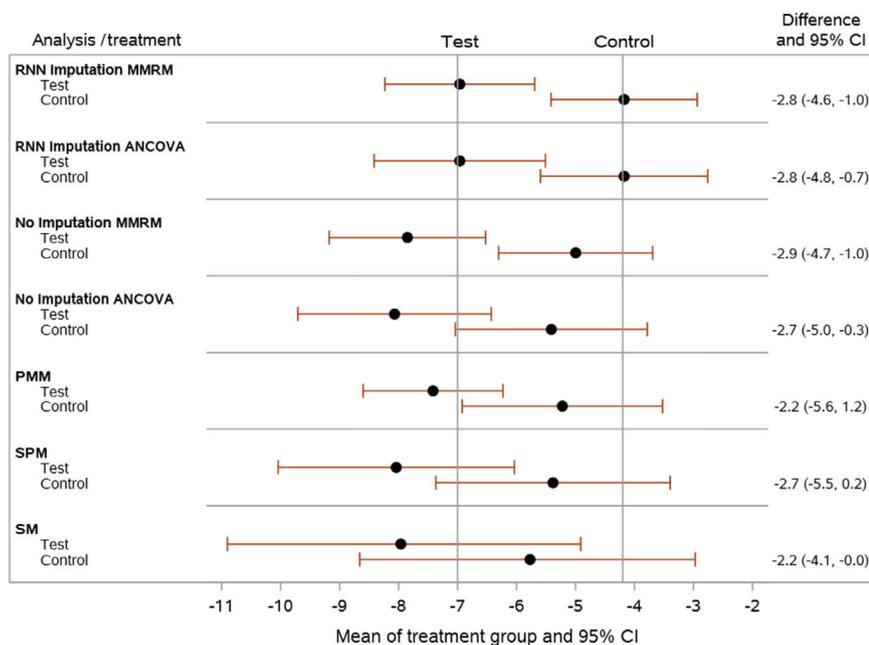


FIGURE 9 Real data: forest plot for analysis results of change from baseline in HAMD17 score at Week 6 using different methods. RNN imputation with OS: nonmissing data + imputed data by the proposed method. No imputation: nonmissing data only. PMM, pattern mixture model; SPM, shared parameter model; SM, selection model

All available variables are used in the RNN model, HAMD17 score as the outcome variable, input variables including gender, treatment, baseline HAMD17 value, HAMD Total score, Patient Global Impression of Improvement (PGI-I), and cluster result (as dummy variables). Based on many experiments, the optimal RNN model for the real data is determined as LSTM with one single layer and seven units in each cell, iteration epoch of 3000, and batch size of 60. The optimal model has been run 100 times, within each bagging, data (including the available data from dropouts) is randomly sampled (smaller clusters are oversampled as mentioned above), the internal weights are updated, and predictions for the missing data are provided within each bagging. The average of 100 predictions is considered as the final prediction. The detailed outputs for clustering and individual prediction are available in the [Supporting Information](#).

The change from baseline in the HAMD17 score is analyzed using the different methods as described in Section 3.4, and the results are presented in the forest plot (Figure 9). The proposed imputation method (i.e., RNN imputation facilitated by clustering and oversampling) + standard analysis method (MMRM and ANCOVA) provided the most conservative estimation for the treatment effect in both treatment groups. Considering the dropouts in the “low responder” clusters (as mentioned above), such conservative estimates may make sense to take into account for the potential impact of MNAR data. Similar to what was observed in the simulation data, there is a systematic bias in the results from the other methods (i.e., MMRM, ANCOVA, PMM, SPM, and SM that are applied without imputing the missing value). Although there are no considerable discrepancies in the point estimates for the difference between treatment groups (maybe due to the dropout rates are similar between treatment groups), the estimates for each treatment group are quite different from the estimates using the proposed method. In general, compared with the proposed method, other methods tend to be optimistic, which may lead to aggressive estimation and hence introduce bias in the study conclusion (especially in the cases when the dropout rate or the efficacy pattern of dropouts are not comparable between the treatment group).

5 | DISCUSSION

As mentioned in the Introduction, it is not possible to ascertain whether the MAR assumptions are appropriate in any practical situation. Therefore, at least a sensitivity analysis to evaluate the impact of MNAR should be warranted if the assumption of MAR cannot be fully justified. In this paper, a machine learning based missing data prediction framework has been developed for longitudinal clinical data with an aim of handling more realistic missing data scenarios. Overall, based on the simulation study, the proposed method provided accurate prediction for both MAR and MNAR data and reduces the bias of missing data in treatment effect estimation. RNN demonstrates the powerful predictive capability for longitudinal data and unrestricted flexibility for nonlinear modeling. Even without being facilitated by any other manner, a straightforward implementation of RNN can provide a fairly good prediction for longitudinal data if there are no severe

MNAR issues in the data. The k -mean trajectory clustering is a crucial step in the proposed method, not only because it facilitates the oversampling and stratified k -fold CV in longitudinal continuous data, but also it is important to borrow information from similar patients by including the cluster information in the RNN model to indicate the longitudinal pattern of efficacy profile. The classic Euclidian distance with Gower adjustment is used in this paper, more insights are needed in the future for other feasible distance metrics and their impact on the imputation results. As the fundamental principle for the imbalanced learning, balancing the classes is key to improve the prediction accuracy for the MNAR data, especially in clinical trials in which the low responders are relatively less and part of them leave the trial due to lack of efficacy. Oversampling the minority class will ensure the efficacy pattern of the low responders been adequately learnt by the model and will also avoid the individual prediction driven by the majority of patients who completed the trial to the overall average level. A simple random oversampling approach (e.g., equalizing the clusters) is taken in this paper, and more insights are needed in the future for other imbalanced learning techniques like different types of oversampling, undersampling, and the combination of both. In addition, it should be noted that the variability of prediction at each study week (measured by the quartiles) is increasing over time (as shown in Section 3.3.4). This hints that the prediction may be less confident for the distant time steps than the near ones. Therefore, when using this method, one should be cautious for the too early dropouts (i.e., the patients with only limited profile available).

In contrast, the commonly used analysis methods (like MMRM and ANCOVA that are applied without imputing the missing value) and classic models (like PMM, SM, and SPM) did not perform as well as the proposed method and showed systematic bias in the treatment effect estimation. These methods tend to overestimate the treatment effect when MNAR is present in the data. This finding is supported in real trial data, that is, a similar pattern of the systematic bias is observed in the real data from an antidepressant clinical trial with a dropout rate of 25%. The performance of classic models in missing data context might need more insights from the practical point of view. Those models heavily rely on the underlying assumption, for example, assuming only MNAR in the data. This kind of assumption can be violated in reality (e.g., missing data can be a mixture of MAR and MNAR), hence leading to suboptimal performance of the model. Additionally, unlike the implementation of joint modeling in complete data (where the binary variable is useful to define the conditional distribution of continuous variable), in the context of missing data, the binary variable only provides the information about missing yes or no. This information actually can also be gained from the continuous variable itself (if its value missing or not). Without providing further information about any feature of the missing data (like potential patterns of efficacy profile or similar patient), the value of such the second process is weakened.

The computational approach comprises necessary components to handle the problem: Step 1: Clustering structures the longitudinal data within each treatment group; Step 2: the RNN models the complex longitudinal data, and the optimal RNN architectures are selected via stratified k -fold CV; Step 3: individual prediction is based on bagging, and the minority class oversampling provides the necessary database for honest predictions; Step 4: average of the bagging predictions is considered as final imputation; Step 5: the imputation results are evaluated at the different levels. Steps 1–3 are reflecting the MNAR problem for longitudinal data with monotones missing patterns. It fits the definition of MNAR, where the missingness depends on the unobserved profile. It is obvious that the proposed methods also incorporate the MAR mechanism by seeking for accurate individual information. The limitation of this paper consists in its special setting studied: monotones missing patterns, continuous longitudinal outcome, three cluster approach with a quite standard metric, and a fixed percentage of MAR and MNAR observations. It is not clear how the proposed strategy will behave in settings that deviate from our assumptions. Therefore, the paper offers an opportunity to encourage the integration of machine learning strategies for handling of missing data in the analysis of randomized clinical trials.

ACKNOWLEDGMENTS

The authors thank Prof. Anne-Laure Boulesteix and Lara Donik for their valuable contribution to this work. The authors also thank the two anonymous reviewers, the associate editor and the editor for their generous and constructive detailed comments that helped us to improve the paper.

CONFLICT OF INTEREST

The authors have declared no conflict of interest.

DATA AVAILABILITY STATEMENT

The data that support the findings of this study are openly available in Supporting Information at <https://doi.org/10.1002/bimj.202000393>.

OPEN RESEARCH BADGES

 This article has earned an Open Data badge for making publicly available the digitally-shareable data necessary to reproduce the reported results. The data is available in the [Supporting Information](#) section.

This article has earned an open data badge **“Reproducible Research”** for making publicly available the code necessary to reproduce the reported results. The results reported in this article could fully be reproduced.

ORCID

Halimu N. Haliduola  <https://orcid.org/0000-0002-9157-8187>

Frank Bretz  <https://orcid.org/0000-0002-2008-8340>

Ulrich Mansmann  <https://orcid.org/0000-0002-9955-8906>

REFERENCES

- Boulesteix, A. L., Wilson, R., & Hapfelmeier, A. (2017). Towards evidence-based computational statistics: Lessons from clinical research on the role and design of real-data benchmark studies. *BMC Medical Research Methodology*, *17*(1), 138. <https://doi.org/10.1186/s12874-017-0417-2>
- Boulesteix, A. L., Binder, H., Abrahamowicz, M., & Sauerbrei, W. (2018). On the necessity and design of studies comparing statistical methods. *Biometrical Journal*, *60*, 216–218. <https://doi.org/10.1002/bimj.201700129>.
- Breiman, L. (2001). Statistical modeling: The two cultures. *Statistical Science*, *16*(3), 199–215. <https://doi.org/10.1214/ss/1009213726>.
- Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). Smote: Synthetic minority over-sampling technique. *Journal of Artificial Intelligence Research*, *16*(2002), 321–357. <https://doi.org/10.1613/jair.953>
- Cho, K., van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H. & Bengio, Y. (2014). Learning phrase representations using RNN encoder-decoder for statistical machine translation. arXiv:1406.1078 [cs.CL].
- Chung, J., Gulcehre, C., Cho, K. & Bengio, Y. (2014). Empirical evaluation of gated recurrent neural networks on sequence modeling. In NIPS 2014 Workshop on Deep Learning, December 2014. Neural Information Processing System Foundation.
- Claesen, M. & De Moor, B. (2015). Hyperparameter search in machine learning. arXiv:1502.02127 [cs.LG].
- Enders, C. K. (2010). *Applied missing data analysis* (pp. 295–301). Guilford Press.
- European Medicines Agency. (2011). *Guideline on missing data in confirmatory clinical trials*. (11–12). EMA.
- Falbel, D., Allaire, J., Chollet, F., RStudio, Google, Tang, Y., Van Der Bijl, W., Studer, M., & Keydana, S. (2015). *Keras*. GitHub. <https://github.com/fchollet/keras>
- Genolini, C., & Falissard, B. (2011). KmL: A package to cluster longitudinal data. *Computer Methods and Programs in Biomedicine*, *104*(3), e112–21. <https://doi.org/10.1016/j.cmpb.2011.05.008>.
- Gers, F. A., Schmidhuber, J., Cummins, F. (1999). Learning to forget: Continual prediction with LSTM. In *9th International Conference on Artificial Neural Networks: ICANN '99* (pp. 850–855). <https://doi.org/10.1049/cp:19991218>.
- Goldstein, D. J., Lu, Y., Detke M. J., Wiltse, C., Mallinckrodt, C., & Demitrack, M. A. (2004). Duloxetine in the treatment of depression: A double-blind placebo-controlled comparison with paroxetine. *Journal of Clinical Psychopharmacology*, *24*, 389–399.
- Gower, J. C. (1971). A general coefficient of similarity and some of its properties. *Biometrics*, *27*(4), 857–871. <https://doi.org/10.2307/2528823>
- Heckman, J. J. (1976). The common structure of statistical models of truncation, sample selection and limited dependent variables and a simple estimator for such models. *Annals of Economic and Social Measurement*, *5*(4), 475–492. <http://www.nber.org/chapters/c10491>
- Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, *9*(8), 1735–1780.
- Kingma, D. P. & Ba, J. (2014). Adam: A method for stochastic optimization. Paper presented at The Third International Conference for Learning Representations, San Diego, 2015. arxiv:1412.6980 [cs.LG].
- LeCun, Y., Bengio, Y., & Hinton, G. E. (2015). Deep learning. *Nature*, *521*(7553), 436–444. <https://doi.org/10.1038/nature14539>
- Little, R. J. A. (1993). Pattern-mixture models for multivariate incomplete data. *Journal of the American Statistical Association*, *88*, 125–134.
- Little, R. J. A. (1994). A class of pattern-mixture models for normal incomplete data. *Journal Biometrika*, *81*(3), 471–483. <https://doi.org/10.2307/2337120>.
- Little, R. J. A. (1995). Modeling the drop-out mechanism in repeated-measures studies. *Journal of the American Statistical Association*, *90*, 1113–1121.
- National Research Council of the National Academies. (2010). *The prevention and treatment of missing data in clinical trials*. (53–54). National Academies Press.
- Olah, C. (2015). Understanding LSTM networks. <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>
- Rubin, D. B. (1976). Inference and missing data. *Biometrika*, *63*(3), 581–592. <https://doi.org/10.1093/biomet/63.3.581>
- Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, *323*(6088), 533–536.
- Santos, M. S., Soares, J. P., Abreu, P. H., & Araujo, H. J. (2018). Cross-validation for imbalanced datasets: Avoiding overoptimistic and overfitting approaches. *IEEE Computational Intelligence Magazine*, *13*(4), 59–76. <https://doi.org/10.1109/mci.2018.2866730>
- Schmidhuber, J. (1993). *Netzwerkarchitekturen, Zielfunktionen und Kettenregel* [Network architectures, objective functions, and chain rule] [Habilitation (postdoctoral thesis)]. Institut für Informatik, Technische Universität München.

Shewalkar, A., Nyavanandi, D., & Ludwig, S.A. (2019). Performance Evaluation of deep neural networks applied to speech recognition: RNN, LSTM and GRU. *JAISCR: Journal of Artificial Intelligence and Soft Computing Research*, 9, 235–245.
 Weiss, G. M. (2013). *Imbalanced learning: Foundations, algorithms and applications*. Wiley-IEEE Press.

SUPPORTING INFORMATION

Additional supporting information may be found in the online version of the article at the publisher’s website.

How to cite this article: Haliduola, H. N., Bretz, F., & Mansmann, U. (2022). Missing data imputation in clinical trials using recurrent neural network facilitated by clustering and oversampling. *Biometrical Journal*, 64, 863–882. <https://doi.org/10.1002/bimj.202000393>

APPENDIX A

A.1 | A short introduction to LSTM

The LSTM unit is more complex than the basic RNN unit. It contains the following elements: two states (cell state and hidden state), input, output, and three gates (forget gate, input gate, and output gate; see Figure A.1). The cell state is the key to LSTM, the horizontal line running through the top of the diagram. The cell state runs straight down the entire chain, with only some minor linear interactions. This allows the LSTM to have the ability to remove or add information to the cell state, carefully regulated by the gates which in a way optionally let information through. The gates are composed of a sigmoid function (see details in Figure A.1), and an element-wise multiplication operation. The sigmoid function squashes information between 0 (“let nothing through”) and 1 (“let everything through”). In practice, the learning capability of LSTM can be improved by including more than one unit in each cell (i.e., one LSTM cell can contain several concatenated LSTM units).

As shown in Figure A.1, at each time step, first, for the information that comes from current input vectors (x_t) and the hidden state at the previous time step (h_{t-1}), the forget gate decides what information to throw away from the cell state. The forget gate is expressed as

$$f_t = \sigma_{sig}(w_{fx}x_t + w_{fh}h_{t-1} + b_f), \tag{A.1}$$

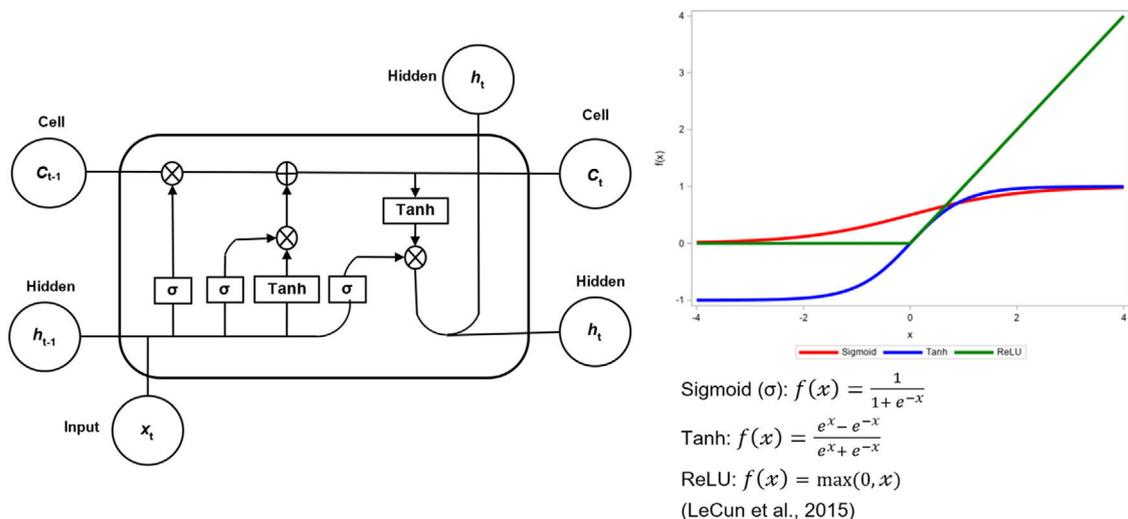


FIGURE A.1 One LSTM unit at time t (reproduced based on Olah, 2015) and the activation functions used in this paper. Note: c_{t-1} = cell state at the previous time step, h_{t-1} = hidden state at the previous time step, x_t = current input vectors, f_t = forget gate, i_t = input gate, \tanh = Tanh activation function, o_t = output gate, h_t = hidden state at current time step, c_t = cell state at current time step, \otimes = element-wise multiplication, \oplus = vector addition

where σ_{sig} is the sigmoid function; w (weights) and b (bias or intercept) is the parameter matrices to be learned. The next step is to decide what new information to store in the cell state. This has two parts. The first part is an input gate i_t , which decides what values to update:

$$i_t = \sigma_{sig}(w_{ix}x_t + w_{ih}h_{t-1} + b_i). \quad (\text{A.2})$$

The second part is a Tanh function (where information is squashed between -1 and 1), which creates a vector of new candidate values which could be added to the cell state, it is expressed as

$$\text{Tanh}_t = \sigma_{Tanh}(w_{Tanhx}x_t + w_{Tanhh}h_{t-1} + b_{Tanh}). \quad (\text{A.3})$$

Then the input gate and the outcome of the Tanh function are combined (element-wise multiplication) as an input section to create an update to the state, the input section is expressed as

$$\text{Tanh}_t \otimes i_t. \quad (\text{A.4})$$

In the next step, the cell state at the previous time step c_{t-1} is updated to the current cell state c_t . This is done by adding the element-wise product of the c_{t-1} and the forget gate f_t (i.e., forgetting the things are decided to forget earlier) and the input section (i.e., adding the new candidate values which are scaled by how much that decided to update the cell state). The current cell state is expressed as

$$c_t = c_{t-1} \otimes f_t + \text{Tanh}_t \otimes i_t. \quad (\text{A.5})$$

The final step is to decide what to store in the current hidden state. This has two parts. First, for the information that comes from current input vectors and previous hidden states, the output gate decides what parts to output. The output gate is expressed as

$$o_t = \sigma_{sig}(w_{ox}x_t + w_{oh}h_{t-1} + b_o). \quad (\text{A.6})$$

The second part is the current cell state going through a Tanh function. Then the element-wise product of these two parts is stored as hidden state for the current time step, expressed as

$$h_t = o_t \otimes \sigma_{tanh}(c_t). \quad (\text{A.7})$$

The current hidden state will be carried forward for the next time step, and it can also be gained as prediction result for the current time step by using an appropriate activation function, for example, a ReLU function for the continuous outcome variable, in that case, the output vector can be expressed as

$$y_t = \text{ReLU}(w_{yh}h_t + b_y). \quad (\text{A.8})$$

A.2 | Details on the data generation process used in the simulation study

The patient baseline characteristics and longitudinal efficacy data are simulated assuming a parallel designed clinical trial. Each patient is designed to be treated and assessed biweekly from Week 0 (baseline) up to Week 16. In total, $3 \times 3 \times 3 = 27$ scenarios are simulated. In each scenario, the patient is randomly assigned to the treatment group (test group : control group = 1 : 1). The longitudinal clinical score decreases over time in general. To take account of the inpatient correlation, and to reflect the influence of the relevant covariates (including the baseline variables and the variables change over time) on the longitudinal profile, the change from previous visit values (per patient per visit) is modeled using several fixed factors and a random effect (more details are provided in Figure A.2). Similar ideas as for the pattern mixture model (i.e., data patterns are different for MNAR and completers), a mixed missing mechanism of MAR (lost to follow-up) and MNAR (dropout due to the lack of efficacy) is considered in the simulation. For the completers and MAR, a completely random effect is considered. For the MNAR, the change from previous visit values shrinks over time (by using the absolute value of the random normal function) so that their scores decrease less or even increase over time. The change from baseline value is calculated by summing up the change from previous visit values up to certain visits within each patient.

Steps	Data Structure	Data Content
Step a: Patient level information	One record per patient	Baseline characteristics and treatment variables Treatment: test=1 (50%), control=0 (50%); Gender: 0 = female (50%), 1 = male (50%); Body weight (kg) ~ N(75, 10**2); Baseline clinical score ~ N(70, 5**2); Discontinuation flags: randomly assign certain number of patients as dropouts (MCAR or MNAR), rules are described as below: (a) the proportion of any missing data in Control group is larger than the Test group (i.e., the difference of proportion between treatment group is 1/2 * dropout rate); (b) in Test group, the proportion of MNAR is slightly lower than MCAR; (c) in Control group, the proportion of MNAR is higher than MCAR.
Step b: Longitudinal struction frame	One record per visit per patient	Visit: Week 0, 2, 4, 6, 8, 10, 12, 14, and 16; Missing record flag: for the dropouts, the monotone missingness started from Week 8, Week 10, or Week 12 based on simulation scenario. Concomitant Medication flags: randomly flag 20% records as taken ConMed, and 80% as not taken ConMed.
Step c: Complete longitudinal score data	One record per visit per patient	Change from previous visit (CFPV) values are modeled as described in below text; Change from baseline (CFB) = cumulative CFPV; Absolute score = baseline + CFB
Step d: Create missing value in longitudinal score data	One record per visit per patient	Set absolute score and CFB to missing value if visit is flagged as missing visit at Step b

FIGURE A.2 Data generation process in simulation study

The control group is more impacted by the missing data as the proportions of MNAR and MAR are much higher in the control group than in the test group. The MAR is influenced by the treatment group but not by any other covariates (i.e., within each treatment group, it is actually an MCAR scenario). The complete simulation data before setting the missing values are kept for all patients at all visits (for the purpose of prediction evaluation). The data generation process is shown in Figure A.2.

Longitudinal clinical score is generated as follows:

- (i) the score with a range of 0–100, the higher the score the worse the disease status. At baseline, the score from all patients follow a normal distribution of mean = 70, SD = 5; body weight follow normal distribution of mean = 75 kg, SD = 10 kg; treatment: test group = 1, control group = 0, that is, the average treatment effect difference = $\frac{4}{3} - 1$ (test - control); gender: 0 = female, 1 = male.
- (ii) at postbaseline visits (from Week 2 to Week 16), the score decreases over time for most cases, the changes from the previous visit (per visit per patient) are modeled as follows:

Change from the previous visit in Score =

$$\frac{75}{\text{BodyWeight}} \frac{\text{BaselineScore} (3 + \text{Treatment}) (\text{Gender} + b)(1 + \text{ConMed})}{70} \frac{30}{\sqrt{1 + \text{Week}^2}} + \beta f(\epsilon) \tag{A.9}$$

where ConMed (concomitant medication) changes over time randomly: 1 = Yes (20%), 0 = No (80%), *b* is the coefficient to determine the importance of gender and ConMed, here *b* = 5 which is chosen empirically; Week = 2, 4, 6, 8, 10, 12, 14, 16, the coefficient of 30 for week is chosen empirically to determine the magnitude of the change from previous visit data; the error function *f*(ϵ) follows a standard normal distribution, β is the coefficient of error term which determines the data pattern, that is, for MNAR (dropout due to lack of efficacy) $\beta = -3.5$ (and the absolute value of *f*(ϵ) is used); for the completers or MAR (dropout due to lost to follow-up), $\beta = 5$ (without an absolute function), the values for coefficients are chosen empirically.

- (iii) The change from baseline at each visit is calculated as accumulation of all changes from previous visit; the absolute score at postbaseline visits is calculated as baseline + change from baseline. For the dropouts, the monotone missingness started from Week 8, Week 10, or Week 12 according to the simulation scenario.

SAS version 9.4 is used for the data generation.

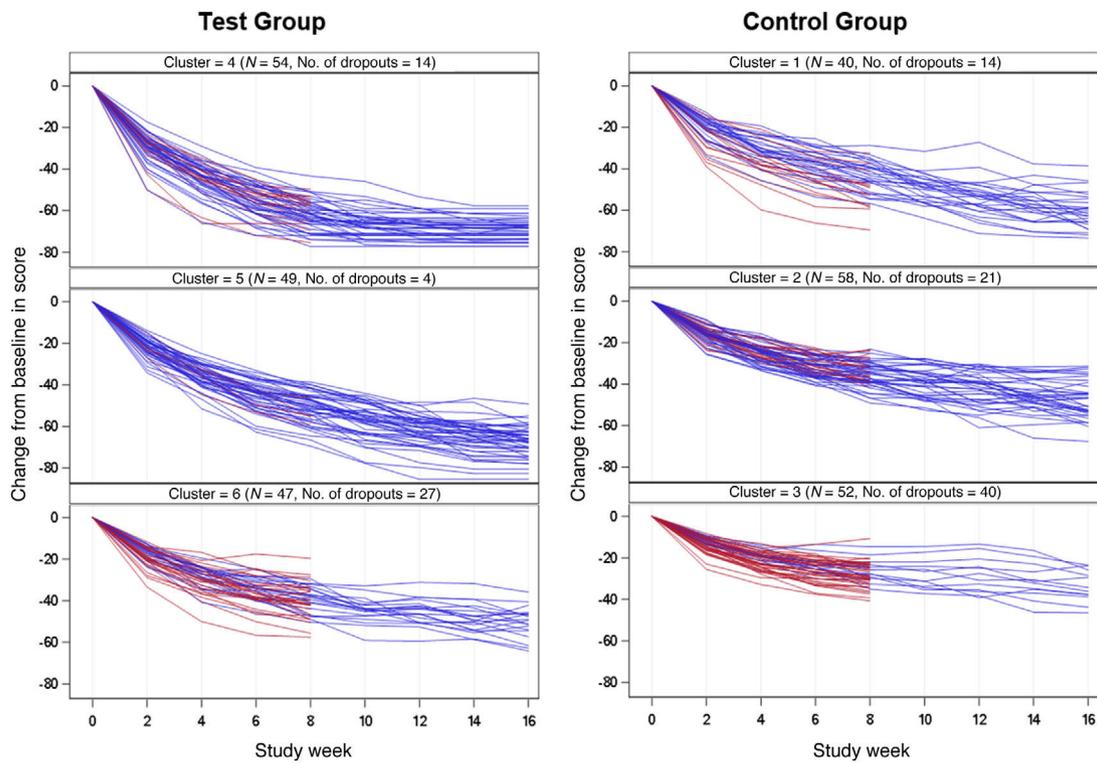


FIGURE A.3 Simulation data: patient profile by the treatment group and cluster (Scenario 300-40-10). Blue lines are for completers, and red lines are for dropouts

A.3 | Longitudinal clustering results in the simulation study

Within each treatment group, patients are clustered into three categories: “good responder,” “medium responder,” and “low responder” according to their longitudinal efficacy profiles. The individual patient profiles by cluster are provided in Figure A.3 for simulation Scenario 300-40-10.

4. Paper II



Contents lists available at ScienceDirect

Computer Methods and Programs in Biomedicine

journal homepage: www.elsevier.com/locate/cmpb

Missing data imputation using utility-based regression and sampling approaches

Halimu N. Haliduola^a, Frank Bretz^{b,c}, Ulrich Mansmann^{a,*}^aInstitute for Medical Information Processing, Biometry and Epidemiology – IBE, LMU Munich, Munich, Germany^bNovartis Pharma AG, Basel, Switzerland^cSection for Medical Statistics, Medical University of Vienna, Vienna, Austria

ARTICLE INFO

Article history:

Received 1 October 2021

Revised 4 February 2022

Accepted 2 October 2022

Keywords:

Missing data

Utility-based regression

SMOTER

Machine learning

ABSTRACT

Data are often missing not at random (MNAR) in scientific experiments. We treat the MNAR problem as an imbalanced learning task. Standard predictive error measures of regression (e.g., mean squared error) are not suitable for imbalanced learning problems, such as in clinical trials where extreme values tend to be MNAR. We investigate hybrid imbalanced learning approaches that combine utility-based regression (UBR) with synthetic minority oversampling technique for regression (SMOTER) in cross-sectional trial settings. UBR optimizes the product of the conditional probability density (estimated by quantile regression forests) and a utility function which takes the relevance of the target variable value and the prediction error into account. SMOTER oversamples the relevant rare cases. Simulations show that the proposed method provides plausible predictions and reduces the bias for realistic missing data scenarios when compared with standard approaches like random forests and multiple imputation (systematic bias is observed in those methods, i.e., a tendency to underestimate the mean and standard deviation given the presence of MNAR in the area of high values of the target variable). The proposed method is implemented in a real dataset from an antidepressant clinical trial, and similar pattern of the systematic bias from commonly used methods is observed in the real data compare to the proposed method. Therefore, we encourage the integration of utility-based learning strategies for handling of missing data in the analysis of clinical trials.

© 2022 The Authors. Published by Elsevier B.V.

This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>)

1. Introduction

Missing data are the data that would be meaningful for the analysis but is not documented. The missing data, if not handled properly, will lead to lower statistical power for the analysis, and may lead to a bias in the estimated treatment effect and an underestimate of the variability. There are three types of missing mechanism [1]. (i) Missing Completely at Random (MCAR): if the probability of missingness does not depend on observed or unobserved measurements, e.g., patients move to another city due to non-health related reasons. (ii) Missing at Random (MAR): if the probability of missingness depends only on observed measurements conditional on the covariates in the model, e.g., younger people may more likely to have blood pressure not measured. (iii) Missing Not at Random (MNAR): if the probability of missingness

depends on unobserved measurements, e.g., patients discontinue from the study due to lack of efficacy.

For the handling of missing data, many methods have been developed under the assumption of MAR or MNAR, respectively. In reality, however, missing data are often a mixture of different types. This makes the assumptions on the missing mechanism violated, which leads to poor performance of the handling methods [2]. To handle realistic missing data scenarios, Haliduola et al. [3] proposed a machine learning based missing data imputation framework where the MNAR problem is treated as an imbalanced learning task (since the MNAR cases are mostly distributed in one tail of the target variable). Take Fig. 1 as an example, depending on the proportion of MNAR data, regions that tend to have MNAR may have a smaller amount of available data than other regions (i.e., an imbalanced distribution). Imbalanced learning is necessary to compensate for the MNAR in that region and to avoid individual predictions being driven by the available non-missing data to an overall average level. They proposed oversampling of minority classes (i.e., the classes with extreme value of the target variable

* Corresponding author.

E-mail address: mansmann@ibe.med.uni-muenchen.de (U. Mansmann).

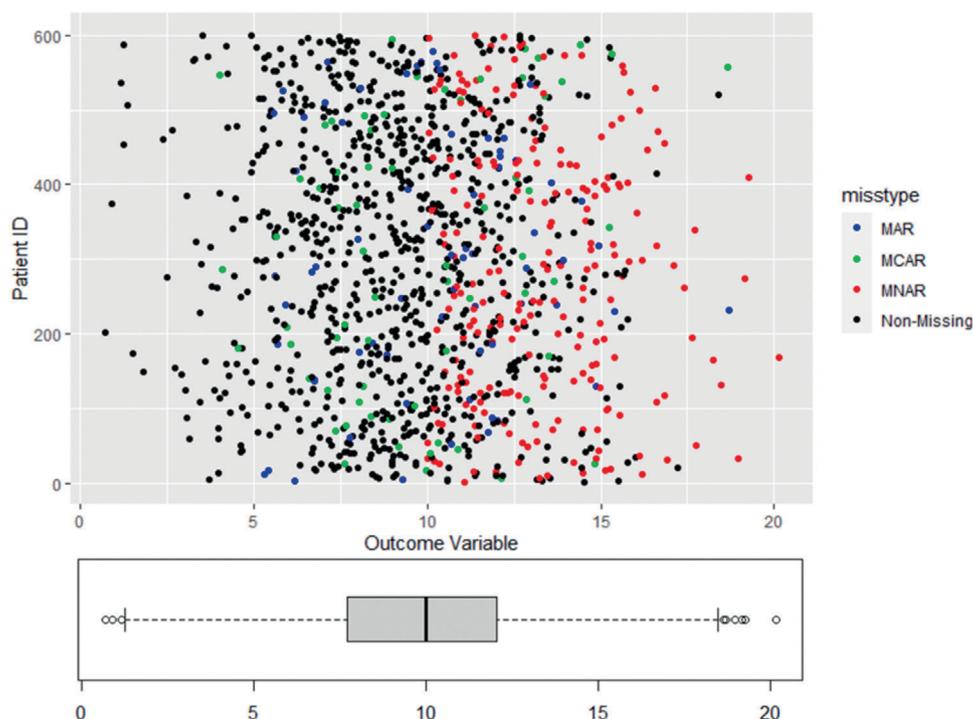


Fig. 1. Simulation data: scatter plot and boxplot for the target variable. The back dots are the non-missing data, blue dots are the MAR, green dots are MCAR, Red dots are MNAR. The details of the data generation process are described in [Section 3.1](#).

that tend to be MNAR), followed by recurrent neural networks to model the data. This framework was shown to be effective for the handling of the missing data based on simulation studies and a real clinical trial data.

Haliduola et al. [3] used a simple random oversampling with replacement and a standard error measure (i.e., mean squared error, MSE). However, these methods come with drawbacks. First, in a simple random oversampling with replacement, random sets of copies of minority class cases are added to the data, which may lead to many duplicates in the minority class. During the learning process, the decision region for the minority class may become very specific and the model will give more focus in that region. For example, in a tree-based learning process, this may lead to new splits in the decision trees, which will result in more terminal nodes (leaves) as the learning algorithm tries to learn more and more specific regions of the minority class, and eventually this will cause overfitting of the model [5]. Secondly, the standard predictive error measure like MSE is not suitable for a regression problem with imbalanced distribution of target variable values in the training data (like in MNAR problem where the extreme values tend to be missing). Their weakness is that they are not sensitive to the location of target variable values [4]. See [Fig. 1](#) as an example, considering the MNAR data (red dots), distribution of available non-missing data (black dots) are imbalanced across the range of target variable (i.e., less available data in the area of high values due to MNAR). If the error measure is not sensitive to the location of target variable values, the area of high values will get less focus in the training process due to the smaller amount of data in that area, and thus the impact of missing data on the aggregated estimation will be ignored. In such cases, it is important to give more focus on the area of high values in the training process to compensate for the MNAR and to avoid the prediction being driven by the frequent cases in the other locations of the target variable. Therefore, it is necessary to have an error metric that is sensitive to the location of the errors, which copes with imbalanced distribution of target variable values.

In this paper, to avoid model overfitting caused by the simple random oversampling, we use the synthetic minority oversampling technique for regression (SMOTER) [7] to oversample the relevant rare cases; and, to overcome the drawbacks of standard error measure, we use the imbalanced learning technique utility-based regression (UBR) [6], which takes both relevance (or importance) of the target variable values and the prediction errors into account in the optimization process. For simplicity, we consider cross-sectional data only. Quantile regression forests [9] are used to estimate the conditional probability density. The optimization process involves determining the maximum integral of the product of the conditional probability density function and the utility function for each case. In light of the “evidence-based computational statistics” [11,12], we evaluate the proposed method in an extensive simulation study using realistic missing data scenarios (i.e., mixture of MCAR, MAR, and MNAR data). The performance of proposed method is evaluated comprehensively in terms of the central tendency and variability of imputed data, prediction accuracy, and a performance comparison with commonly used methods like random forests and multiple imputation. Finally, we illustrate the proposed method with a real dataset from an antidepressant clinical trial, which is one of the few publicly available datasets that can be used to demonstrate methods for handling missing data where a continuous outcome is measured.

2. Methods

In this paper, we aim to handle realistic missing data scenario (i.e., mixture of MCAR, MAR, and MNAR data) in a continuous outcome variable. We treat the MNAR problem in clinical trials as an imbalanced learning task. We investigate a hybrid imbalanced learning approach that combines utility-based regression (UBR) [6] with synthetic minority oversampling technique for regression (SMOTER) [7] in the missing data imputation. First, we assign a relevance to the target variable values based on their distribution in the training data (i.e., available non-missing data) and

define a threshold for oversampling. The second step is data pre-processing, where we use the SMOTER method to oversample cases with relevance greater than the threshold. In the third step, we apply utility-based regression on the oversampled training data, and the model parameters are optimized by maximizing the relevance and minimizing the error simultaneously. The final step is to use the optimal model to predict the missing target variable values.

2.1. Utility-based regression

Let Y be a target variable and X predictor vector. The most commonly used error measures in regression are the mean squared error

$$MSE = \frac{1}{n} \sum_{i=1}^n (\bar{y}_i - y_i)^2,$$

and the mean absolute error

$$MAE = \frac{1}{n} \sum_{i=1}^n |\bar{y}_i - y_i|.$$

The standard predictive error measures are not suitable for a regression problem with imbalanced distribution of target variable values in the training data (like in MNAR problem where the extreme values tend to be missing). Their weakness is that they are not sensitive to the location of target variable values [4]. Take Fig. 1 as an example, if the error measure is not sensitive to the location of target variable values, the area of high values will get less focus in the training process due to the smaller amount of data in that area (due to MNAR), and thus the impact of missing data on the aggregated estimation will be ignored. In such cases, it is important to give more focus on the area of high values in the training process to compensate for the MNAR and to avoid the prediction being driven by the frequent cases in the other locations of the target variable. Therefore, it is necessary to have an error metric that is sensitive to the location of the errors, which copes with imbalanced distribution of target variable values. It should be noted that the example in Fig. 1 is used to demonstrate the idea, it could be the other way round in reality, i.e., the lower values tend to be MNAR.

Utility is a function of both the error of the prediction and the relevance (or importance) of both true and predicted values. Together, the relevance and loss information give a utility function, which provides more reliable evaluation of a regression model. The ultimate goal of utility-based regression is to maximize the utility, which is achieved by maximizing the relevance and minimizing the error simultaneously. In following sections, we use the notations for utility-based regression defined by Torgo and Ribeiro [6] and Ribeiro [4].

The relevance is the crucial property that distinguishes non-uniform cost/benefit regression problems from those standard regression problems. The relevance function $\vartheta(Y): y \rightarrow [0, 1]$ is a continuous function that expresses the domain-specific importance concerning the target variable domain y by mapping it into a $[0, 1]$ scale of relevance, where 0 represents the minimum and 1 represents the maximum (see an example in Fig. 2). To take both predicted value (\bar{y}) and true value (y) into account, the joint relevance function is defined as weighted average:

$$\vartheta(\bar{y}, y) = (1 - p) \vartheta(\bar{y}) + p \vartheta(y)$$

where $p \rightarrow [0, 1]$ is the weight, e.g., $p = 0.5$. The actual form the relevance function is domain specific and defined by the user based on the problem in hand.

For the missing data problem in a continuous target variable (like the example mentioned above), a relevance function can be defined to assign more relevance/importance to the extreme values

in one tail or both tails according to the distribution of available data. For example, we use boxplot whiskers or summary statistics like the first quartile (Q1) and the third quartile (Q3) to identify the extreme values. In the example in Fig. 2, the extreme values are identified using the boxplot whiskers, and then the maximum relevance of 1 is assigned to the extreme cases, and minimum relevance of 0 is assigned to the median value. A monotone cubic spline interpolation line over a set of maximum and minimum relevance points is the actual shape of the relevance function [8]. Using the boxplot to identify the extreme values, a coefficient needs to be specified to determine how far the whiskers extend to the extreme data points in the boxplot (e.g., a coefficient of 1.5 as in the standard boxplot). The choice of the coefficient should be based on the specific problem in hand and it should be pre-specified. For example, a coefficient smaller than 1.5 can be considered to assign high relevance to more data points. A range of the coefficients can also be considered to perform the sensitivity analysis. In our example, considering the presence of MNAR in the area of high values and the MCAR/MAR spread out across the whole range of target variable, it makes sense to assign more relevance for both high and low extreme values. Assigning high relevance in both tails may also avoid disproportionately heavy in one tail over the other in the prediction.

The cost of a prediction is defined as product of the relevance and the loss (or error) function,

$$c(\bar{y}, y) = \vartheta(\bar{y}, y) C_{max} L(\bar{y}, y)$$

where $\vartheta(\bar{y}, y)$ is the joint relevance function, C_{max} is the maximum cost that is only assigned when the relevance is maximum (i.e., $\vartheta(\bar{y}, y) = 1$). The term $\vartheta(\bar{y}, y) C_{max}$ can be seen as a case-specific maximum cost value, i.e., the maximum penalty we get if \bar{y} is the “worst possible” prediction for the particular case under consideration. $L(\bar{y}, y)$ is the loss function. It is important to scale the loss function to $[0, 1]$. Torgo and Ribeiro [6] defined a percentage-type loss function as the difference between the maximum and minimum relevance in the interval between the true and predicted values.

$$L(\bar{y}, y) = \frac{\max_{i \in \bar{y}..y} \vartheta(i) - \min_{i \in \bar{y}..y} \vartheta(i)}{1}$$

The total cost can be calculated by summing up all individual cost values. It is important to notice that when asserting the cost of a prediction, it is necessary to take both the true and the predicted values into account. Predicting an irrelevant value for a case that has an actual extreme value is not the only cost that can occur. It may be equally serious to predict an extreme value for a frequent case, as it causes false alarm that could lead to serious cost. Therefore, the joint relevance function is used in the above cost function. In addition, it makes sense to use weight $p = 0.5$ in the joint relevance function to give equal importance to both types of error.

The benefit of a prediction is defined as product of the relevance of true value and the complementary of the loss,

$$b(\bar{y}, y) = \vartheta(y) B_{max} (1 - L(\bar{y}, y))$$

where $\vartheta(y)$ is the relevance function of true value, B_{max} is the maximum reward that is only assigned when the relevance is maximum. In the benefit function, only the relevance of the true value is considered as the purpose is to assert how well a model predicts the test cases that are relevant (i.e., rewards the accurate prediction for the relevant values). The total benefit can be calculated by summing up all individual benefit values.

The utility of a prediction is the net balance between its benefits and costs, defined as,

$$U(\bar{y}, y) = b(\bar{y}, y) - c(\bar{y}, y)$$

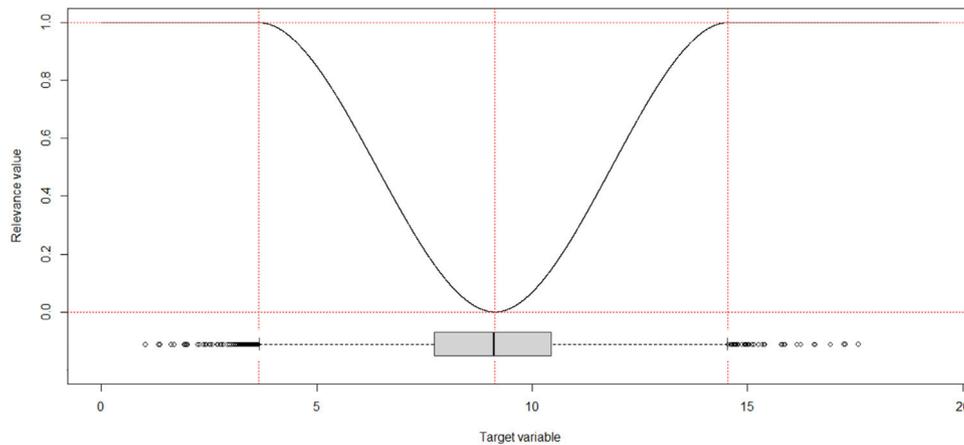


Fig. 2. An example of relevance function to assign more importance to the extreme values according to the distribution of available data.

The total utility can be calculated by summing up all individual utility values. The mean utility can also be calculated as utility-based model performance metrics.

2.2. Quantile regression forests

As mentioned in Section 2.1, the ultimate goal of utility-based regression is to optimize the utility, which is achieved by maximizing the relevance and minimizing the error simultaneously. In this paper, we use the optimization process proposed by Rau et al. [14]. This method uses quantile regression forests (QRF, [9]) to estimate the conditional probability density which is a crucial element in the optimization process. To elaborate the main idea of QRF, we start with the random forests (RF, [16]) and quantile regression [18].

The random forests build k trees in parallel using n independent observations (y_i, x_i) , $i = 1, \dots, n$. Each tree is based on the bootstrapped data (random sampling with replacement, e.g., use 2/3 as the original data size) and random subset of variables (e.g., use 1/3 of all feature variables). This kind of variety is what makes random forests more effective than individual decision tree. Let θ denote the random parameter vector that determines how a tree is grown (e.g., which variables are considered for split points at each node), the corresponding tree is denoted by $T(\theta)$, let L_f denote the leaves of the tree ($L = 1, \dots, m$). For every $x \in X$, there is only one leaf L_f can be obtained when dropping x down the tree. Denote this leaf by $L_f(x, \theta)$ for tree $T(\theta)$. For a single tree, the weight vector $w_i(x, \theta)$ is a positive constant if observation x_i is part of leaf $L_f(x, \theta)$ and 0 if not, and the weights $w_i(x, \theta)$ sum to 1. The prediction of a single tree k , given the feature $X = x$, is then weighted average of the original observations y_i ,

$$\bar{u}_t(x) = \sum_{i=1}^n \omega_i(x, \theta) y_i,$$

where t is the t th single tree, $t = 1, \dots, k$. The conditional mean $E(Y|X = x)$ is approximated by the averaged prediction of k single trees, each constructed with an independent and identically distributed vector θ_t . Let $\omega_i(x)$ be the average of $\omega_i(\theta)$ over the trees, defined as,

$$\omega_i(x) = k^{-1} \sum_{t=1}^k \omega_i(x, \theta_t).$$

The predictions of random forests are then the weighted conditional mean,

$$\bar{u}(x) = \sum_{i=1}^n \omega_i(x) y_i.$$

The weighted conditional mean is estimated by minimizing the MSE:

$$E(Y|X = x) = \arg \min_{\bar{y}} E\{(\bar{y} - y)^2 | X = x\}.$$

The conditional mean describes only one aspect of the conditional distribution of a target variable Y , while the quantile regression aims to provide more information about the conditional distribution, e.g., the conditional quantiles [18]. For $X = x$, the conditional distribution function $F(y|X = x)$ is given by the probability of Y is smaller than $y \in \mathbb{R}$ (\mathbb{R} is the space for the target variable),

$$F(y|X = x) = P(Y \leq y | X = x).$$

For a continuous distribution function, given $X = x$, the α -quantile $Q_\alpha(x)$ is then defined such that the probability of Y being smaller than $Q_\alpha(x)$ is exactly equal to α ($0 < \alpha < 1$). The quantiles $Q_\alpha(x)$ give more information about the conditional distribution of Y , which is defined as,

$$Q_\alpha(x) = \inf\{y : F(y|X = x) \geq \alpha\}.$$

The loss function L_α is defined as the weighted absolute deviations,

$$L_\alpha(y, q) = \begin{cases} \alpha |y - q| & y > q \\ (1 - \alpha) |y - q| & y \leq q \end{cases}$$

The conditional quantiles are estimated by minimizing the expected loss $E(L_\alpha)$,

$$Q_\alpha(x) = \arg \min_q E\{L_\alpha(Y, q) | X = x\}.$$

For quantile regression forests, trees are grown as in the standard random forests algorithm [9]. The conditional distribution is then estimated by the weighted distribution of observed target variables, where the weights ($\omega_i(x)$) attached to observations are identical to the original random forests algorithm. The key difference from the standard random forests is that, for each node in each tree, QRF keeps the value of all observations in this node (not just their mean as in the standard random forests), and assesses the conditional distribution of those observations.

For $X = x$, the conditional distribution function of Y is given by,

$$F(y|X = x) = P(Y \leq y | X = x) = E(I_{\{Y \leq y\}} | X = x)$$

where $I_{\{Y \leq y\}}$ is the indicator function, which equals to 1 if $Y \leq y$ otherwise 0. Just as $E(Y|X = x)$ is approximated by a weighted mean of Y , define an approximation to $E(1_{\{Y \leq y\}}|X = x)$ by the weighted mean over the observations of $1_{\{Y \leq y\}}$ as the prediction of QRF,

$$\bar{F}(y|X = x) = \sum_{i=1}^n \omega_i(x) 1_{\{Y \leq y\}}.$$

The optimization process uses a method proposed by Rau et al. [14], which use QRF to estimate the conditional probability density. In regression, for each case, this process involves determining the maximum integral of the product of the conditional probability density function and the utility function. The optimal prediction for $X = x$ is given by,

$$\bar{y}(X = x) = \arg \max[\bar{y}] \int pdf(y|X = x) U(\bar{y}, y) dy$$

where $pdf(y|X = x)$ is the conditional probability density estimation for $X = x$, and $U(\bar{y}, y)$ is the utility evaluated on the true value y and predicted value \bar{y} . Final predictions are the conditional means take target variable utility into account. We use the R package “UBL” (stands for “Utility-Based Learning”, [13,15]) in this paper.

2.3. SMOTER

Synthetic Minority Oversampling Technique (SMOTE) was introduced by Chawla et al. [5] for the classification task. This algorithm operates in the feature space rather than target variable space (as all rare cases have the same target minority class). The minority class is oversampled by taking each minority sample and introducing synthetic examples along the line segments joining any/all of the k minority class nearest neighbors (e.g., $k = 5$). For example, if the amount of oversampling needed is 200%, only two neighbors from the k nearest neighbors are chosen and one sample is generated in the direction of each. Synthetic samples are generated in the following way: take the difference between the feature vector under consideration and its nearest neighbor. Multiply this difference by a random number between 0 and 1, and then add it to the feature vector under consideration.

Torgo et al. extended the SMOTE for regression task (i.e., the SMOTER) in 2013. Three key components were addressed in the extension: the relevance function (i.e., the $\emptyset(Y)$ as discussed in Section 2.1) and the user-specified threshold for the relevance values were used to define the relevant (rare) cases and the frequent cases (e.g., relevance threshold = 0.8); the same approach as in the original algorithm was used to generate the synthetic feature samples; the weighted average of the target variable values of the two seed examples (i.e., the case and the selected neighbor) was used as the synthetic value for the target variable (the weights are calculated as an inverse function of the distance of the generated new case to each of the two seed examples). We use the R package “UBL” [15] for the implementation of SMOTER.

In practice, it is common to implement the SMOTER together with the undersampling of frequent cases. However, in this paper, we do not consider the undersampling for following reasons: 1) In realistic missing data scenarios, the MNAR data are located in certain area of the target variable, but the MAR/MCAR data may spread out in the whole range of the target variable. In the training process, it may not be a conservative approach to give less focus on the locations where MAR/MCAR data may appear; 2) The undersampling reduces the size of the training data, this may not be a favorable approach in clinical trials in which the total amount of data is normally not massive.

In SMOTE, the amount of oversampling is a hyperparameter of the system [5]. We fine-tuned the appropriate amount of oversampling using the cross-validation (CV) approach. It is important to

	MCAR	MAR	MNAR	OUT-COM	COV1	COV2	COV3	COV4	COV5	COV6	COV7
MCAR	1	0	0	0	0	0	0	0	0	0	0
MAR	0	1	0	0	0.4	0	0	0	0	0	0
MNAR	0	0	1	0.5	0	0.2	0.2	0.2	0.2	0.2	0.2
OUTCOME	0	0	0.5	1	0	0.5	0.5	0.5	0.5	0.5	0.5
COV1	0	0.4	0	0	1	0	0	0	0	0	0
COV2	0	0	0.2	0.5	0	1	0.2	0.2	0.2	0.2	0.2
COV3	0	0	0.2	0.5	0	0.2	1	0.2	0.2	0.2	0.2
COV4	0	0	0.2	0.5	0	0.2	0.2	1	0.2	0.2	0.2
COV5	0	0	0.2	0.5	0	0.2	0.2	0.2	1	0.2	0.2
COV6	0	0	0.2	0.5	0	0.2	0.2	0.2	0.2	1	0.2
COV7	0	0	0.2	0.5	0	0.2	0.2	0.2	0.2	0.2	1

Fig. 3. Example of correlation matrix used in the simulation data generation.

note that only the training data should be oversampled during the CV process, the validation data should never be oversampled to avoid the “overoptimism” issue [19].

3. Simulation study to evaluate performance of methods

3.1. Design of simulation study

To demonstrate the idea of the utility-based regression and sampling approaches, we consider the cross-sectional data only in this paper. In the simulation study, random data is generated for 600 subjects. The outcome variable and covariates (predictors) are normally distributed. Missing data indicators are binary variables (i.e., separate indicator variables for MCAR, MAR and MNAR). Correlated normal and binary data are generated simultaneously using the point-biserial correlation approach of Demirtas and Dogana [20]. Suppose that X and Y follow a bivariate normal distribution with a correlation of ρ_{XY} . If X is dichotomized to produce X_D , then the resulting correlation between X_D and Y can be given as point-biserial correlation,

$$\delta_{X_D Y} = \rho_{XY} \left(\frac{h}{\sqrt{p(1-p)}} \right)$$

where p is the proportion of the observations above the point of dichotomization, and h is the ordinate (probability density function) of the normal curve at the same point.

In the simulation study, we simultaneously generate one outcome variable and seven covariates (each normally distributed with mean 10 and variance 10) and 3 missingness indicators using a given correlation matrix (see Fig. 3 for an example). The MCAR flag (with missing data proportion of 5%) is independent from any other variables. The MAR flag (with missing data proportion of 5%) is correlated with the first covariate only (correlation coefficient = 0.4) and independent from the outcome variable and the other covariates. To evaluate the performance of imputation method properly, we consider higher proportion of MNAR data (i.e., 25%). The MNAR flag is positively correlated with outcome variable (i.e., the higher values tend to be missing, correlation coefficient = 0.5) and the second to seventh covariates (correlation coefficient = 0.2). The outcome variable is correlated with the MNAR flag and the second to seventh covariates (correlation coefficient = 0.5). The first covariate is correlated with MAR flag only. The second to seventh covariate are correlated with the outcome, therefore they are also correlated with each other (correlation coefficient = 0.2). See Fig. 1 as an example for the distribution of the outcome variable. We use the R package “BinNor” [21] in the data generation. Since the higher values of outcome variable tend to be missing (MNAR), the mean of the available non-missing data is an underestimation of the true value. A proper missing imputation method should compensate for the MNAR and reduce the bias

in the aggregated estimation. In this paper, we perform the simulation with 100 replications.

We impute the missing data using proposed method, i.e., UBR facilitated by SMOTER (`ubr.smt`). In the SMOTER process, we identify the relevant extreme values based on the summary statistics of available training data, i.e., the data points \leq the first quartile (Q1) or \geq the third quartile (Q3) are oversampled. The amount of oversampling is determined as 3 times as the available data in both tails based on the cross-validation. In the UBR process, we assign relevance function to target variable using the boxplot with a coefficient of 0.75 (i.e., half of the standard coefficient). Based on the summary statistics of available data, a coefficient of 0.75 is considered as appropriate to assign relevance to the high target variable values where tend to have MNAR and also the low extreme values. A range of coefficients (0.5, 0.6, 0.7, 0.8 and 0.9) are also experimented to illustrate the impact of relevance function on the imputation performance. As mentioned above, the relevance function is defined according to the distribution of the available data, and there is a shift in the central tendency of the available data due to MNAR in the area of high values. This shift is also reflected in the relevance function, which leads to more relevance given in the area of high values (this is considered as a conservative approach given the presence of MNAR in that area only in this case).

3.2. Measuring performance of the proposed methods

To compare the performance of proposed method (i.e., UBR facilitated by SMOTER), we impute the missing data using other methods including:

- `ubr.org` = UBR without facilitating by SMOTER.
- `qrf.smt` = QRF facilitated by SMOTER, details of QRF are described in Section 2.2. We use the R package “`quantregForest`” [10] in the implementation.
- `qrf.org` = QRF without facilitating by SMOTER.
- `rf.smt` = random forests facilitated by SMOTER, details of RF are described in Section 2.2. We use the R package “`randomForest`” [17] in the implementation.
- `rf.org` = random forests without facilitating by SMOTER.
- `mi` = traditional multiple imputation under the assumption of MAR. In addition to those machine learning-based methods, comparisons with the most commonly used traditional statistical methods (i.e., multiple imputation) are also considered meaningful. We use the R package “`MICE`” (van Buuren et al. [24]) with 200 multiple imputations. MICE stands for Multivariate Imputations by Chained Equations, which generates multiple imputations for incomplete multivariate data by Gibbs sampling. The algorithm imputes an incomplete target column by generating “plausible” synthetic values given other columns (covariates) in the data. The imputation method for the missing continuous outcome variable is predictive mean matching ([22] and [23]).

We perform the following measures to compare the performance of difference methods:

- Calculate the mean and standard deviation (SD) of the imputed outcome variable by different imputation methods as mentioned above, and compare with the mean and SD of true value (i.e., the complete outcome variable before set the missing values). If the estimations are close to the mean and SD of true value then the imputation method is appropriate. To show the bias that caused by missing data, the mean and SD of available non-missing data are also provided.
- Perform one sample *t*-test on the imputed data with a null-hypothesis of mean = 10, the larger *p*-values indicate better imputation performance.

- Perform a simple linear regression of imputed value versus the true value, and compare the intercepts (close to 0 is better) and the slopes (close to 1 is better).

3.3. Simulation results

We visualize the performance measures from 100 studies using the boxplot. In Fig. 4, the boxplots for the mean values from 100 studies per scenario are presented. The true means follow normal distribution around 10 (the blue box). The bias caused by the missing data is substantial, the means estimated from non-missing available data are significantly lower than the true means (i.e., `noimp`, the brown box on the right in below figure). The means estimated based on imputed data by the proposed method (i.e., UBR + SMOTER) are the closest to the true means (the green box) when comparing with other methods. The means from the UBR without SMOTER (the light green box) are the second closest estimation of the true means. The QRF and RF perform very similarly (the boxes labeled as `qrf.org` and `rf.org`), which is expected as the goal is to provide the conditional mean as prediction. When facilitating by SMOTER, QRF and RF perform better than without SMOTER but still are not as good as the proposed method (the boxes labeled as `qrf.smt` and `rf.smt`). The traditional multiple imputation is not as good as the proposed method (the purple box). In general, all other methods tend to underestimate the mean given the presence of MNAR in the area of high values of the target variable.

It is also important to evaluate the performance of imputation method in terms of the variability of imputed data. As shown in Fig. 5, similar as for the central tendency measure (i.e., the mean), the proposed method provides the closest estimation for the SD, followed by the UBR without facilitating by SMOTER. All other methods tend to underestimate the SD given the presence of MNAR in the area of high values of the target variable.

We perform sensitivity analysis in terms of the coefficient of the relevance function. A range of coefficients (i.e., 0.5, 0.6, 0.7, 0.8 and 0.9) are experimented and results are shown in Fig. 6 (A for distribution of mean and B for distribution of SD). It is clear that the relevance function impacts the performance of UBR considerably. The coefficient here is a parameter to determine how far the whiskers extend to the extreme data points in the boxplot when defining the relevance function. The higher coefficients result in high relevance been assigned to the more extreme cases (e.g., for less data points), this may increase the variability of the predicted values. As mentioned in Section 3.1, there is a shift in the relevance function due to MNAR in the area of high values, this leads to even less lower extreme values been assigned with high relevance. Therefore, higher coefficients result in higher estimated mean and SD in this case. As mentioned above, all commonly used methods tend to underestimate the mean and SD given the presence of MNAR in the area of high values of the target variable. It would be equally worse to overestimate the mean and SD (e.g., in the case of coefficient = 0.9). Therefore, it is important to pre-specify a proper relevance function according to the distribution of available data and make a plausible assumption on the missing data (i.e., the possible locations of target variable scale where the missing data tend to occur). It is also important to perform sensitivity analysis with different relevance functions (and associated parameters) to check the appropriateness and robustness of the primary analysis.

We perform one sample *t*-test on the imputed data (imputed by different methods) with a null-hypothesis of mean = 10 and present the distribution of *p*-values in Fig. 7. For the true data (where no missing data), the *p*-values are mostly greater than 0.05 as expected. For the proposed method (i.e., UBR + SMOTER), the majority of the *p*-values are greater than 0.05. While for other

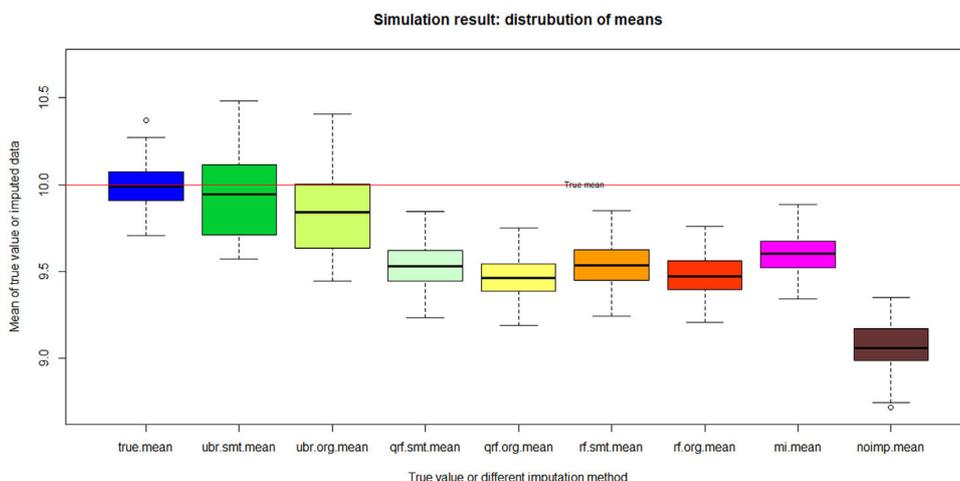


Fig. 4. Simulation result – distribution of means of imputed data by different methods. ubr = utility-based regression (coefficient of relevance function = 0.75), smt = SMOTER data, org = original data, qrf = quantile random forests, rf = random forests, mi = multiple imputation, noimp = no imputation.

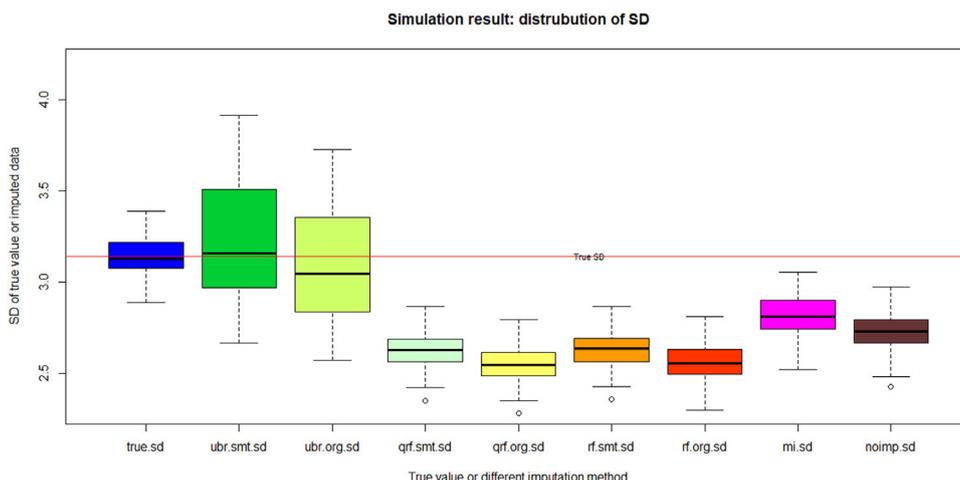


Fig. 5. Simulation result – distribution of SDs of imputed data by different methods. ubr = utility-based regression (coefficient of relevance function = 0.75), smt = SMOTER data, org = original data, qrf = quantile random forests, rf = random forests, mi = multiple imputation, noimp = no imputation.

methods, the p-values are quite small (mostly < 0.05). Although the p-value is sample size dependent, but the trend is clear to show that the proposed method is better than other methods in terms of the ability to reduce the bias of missing data in the aggregated estimation.

We perform simple linear regression for the true value versus the imputed value (by different method). The intercept and the slope from the linear regression are visualized using the boxplots in Fig. 8. The proposed method (i.e., UBR + SMOTER) gives the least intercept and the greatest slope (i.e., closest to 1), suggesting the best performance among all the methods.

4. Real data example

We implement the proposed method in a real dataset from an antidepressant clinical trial, which is available on the website of London School of Hygiene and Tropical Medicine [25]. Original data are from an antidepressant clinical trial with four treatments; two doses of an experimental medication, a positive control, and placebo [26]. There are 26.1% and 25.0% patients with missing Hamilton 17-item rating scale for depression (HAMD17) at Week 6 in Control group (i.e., placebo, $N = 88$) and Test group (i.e., created by randomly selecting patients from the three non-placebo arms, $N = 84$), respectively.

We use the HAMD17 at Week 6 as the target variable (cross-sectional data), use the treatment group and the available baseline variables as predictors (including the gender, baseline HAMD17 value, HAMD Total score and Patient Global Impression of Improvement (PGI-I)). The reasons for discontinuation are not available in the published dataset, this makes it difficult to make assumption about the missing mechanism. We define the relevance function according to the summary statistics of the available data. In the pre-processing, the data points $\leq Q1$ or $\geq Q3$ are oversampled using SMOTER method. The amount of oversampling is determined as twice as the original available data in both tails based on the cross-validation. In the UBR process, maximum relevance of 1 is assigned to the data points $\leq Q1$ or $\geq Q3$, and minimum relevance of 0 is assigned to median value (note: it is not the boxplot method in this case and therefore no coefficient to be determined). A monotone cubic spline interpolation line over a set of maximum and minimum relevance points is the actual shape of the relevance function.

To compare the imputation performance, we impute the missing data using the methods as described in Section 3.2. The imputed outcome variable (i.e., change from baseline in HAMD17 score at Week 6) is analyzed using the analysis of covariance (ANCOVA) model with treatment as factor and baseline value as covariate. To show the bias that caused by the missing data, we also

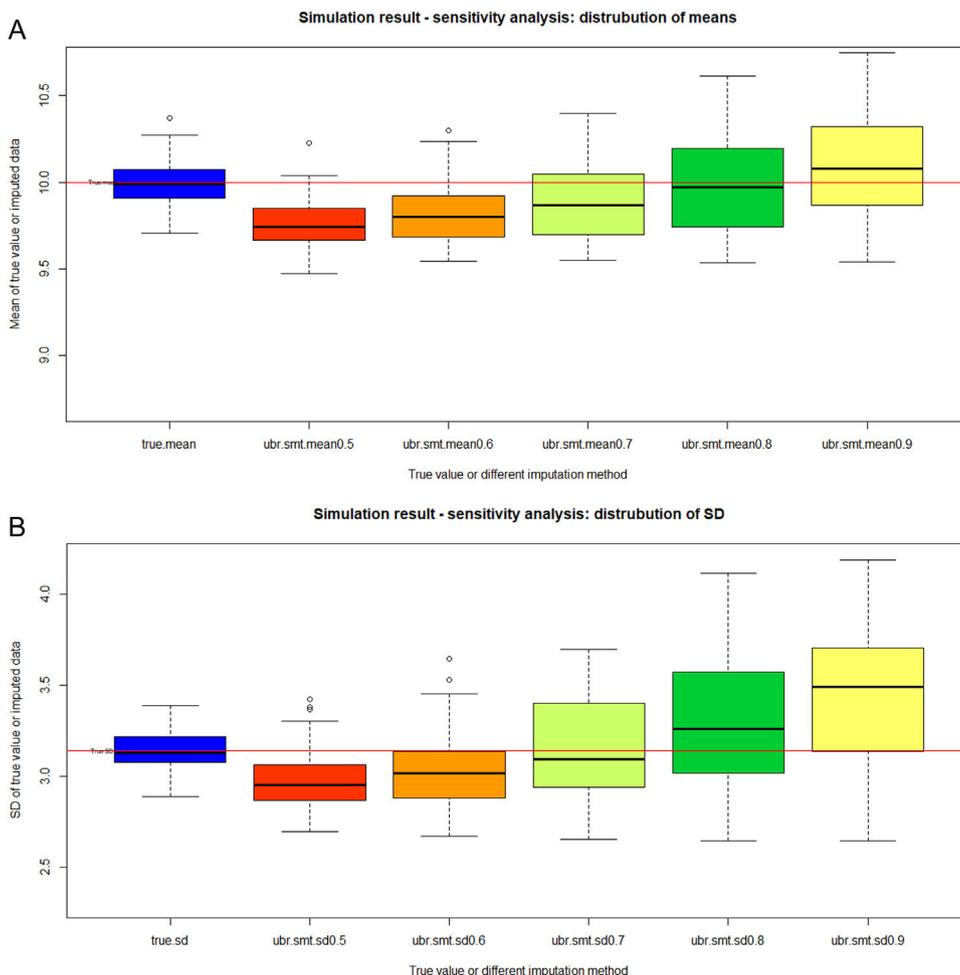


Fig. 6. Simulation result – sensitivity analysis: distribution of means (A) and SDs (B) of imputed data by ubr+smt using different coefficients in the relevance function (0.5, 0.6, 0.7, 0.8, 0.9). ubr = utility-based regression, smt = SMOTER data.

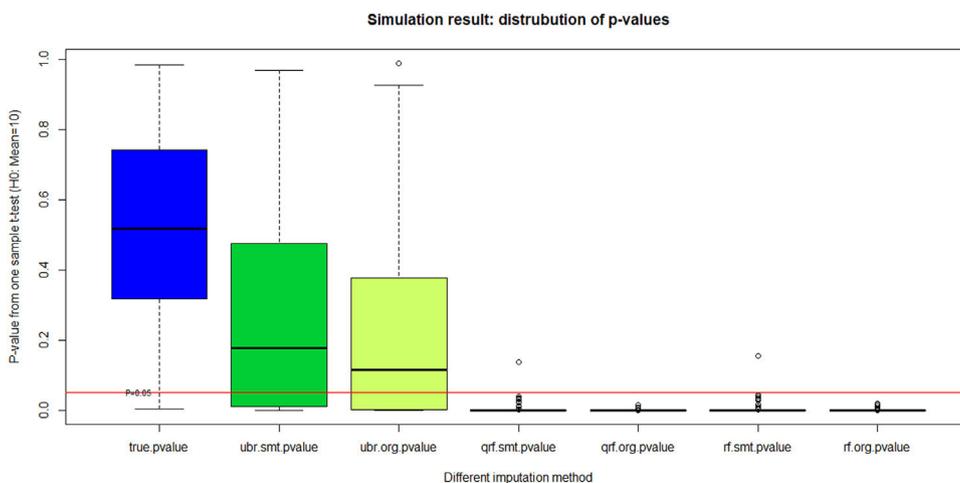


Fig. 7. Simulation result – distribution of p-values from one sample *t*-test on imputed data by different methods. ubr = utility-based regression (coefficient of relevance function = 0.75), smt = SMOTER data, org = original data, qrf = quantile random forests, rf = random forests.

analyze the outcome variable without imputation using ANCOVA. The results from the different approach are presented in forest plot (Fig. 9). The proposed imputation method (i.e., UBR + SMOTER) provided the most conservative estimation for the treatment effect in both treatment groups. There is systematic bias in the results from other methods. This bias is more pronounced in the Control

group, a possible reason could be there are more low responders with missing data in Control group (e.g., may be more MNAR in Control group). In general, comparing with the proposed method, other methods tend to be optimistic, which may lead to aggressive estimation and hence introduce bias in the study conclusion (es-

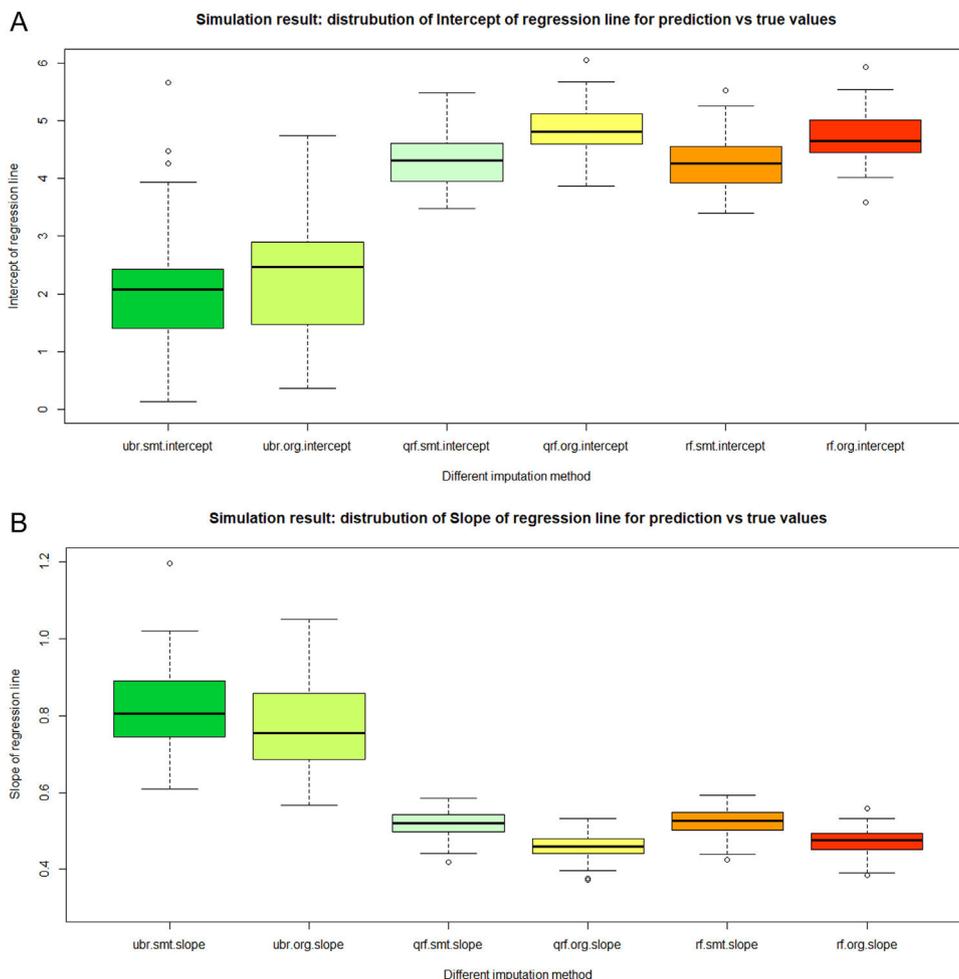


Fig. 8. Simulation result – distribution of the intercepts (A) and slops (B) from simple regression of true data vs. imputed data by different methods. ubr = utility-based regression (coefficient of relevance function = 0.75), smt = SMOTER data, org = original data, qrf = quantile random forests, rf = random forests.

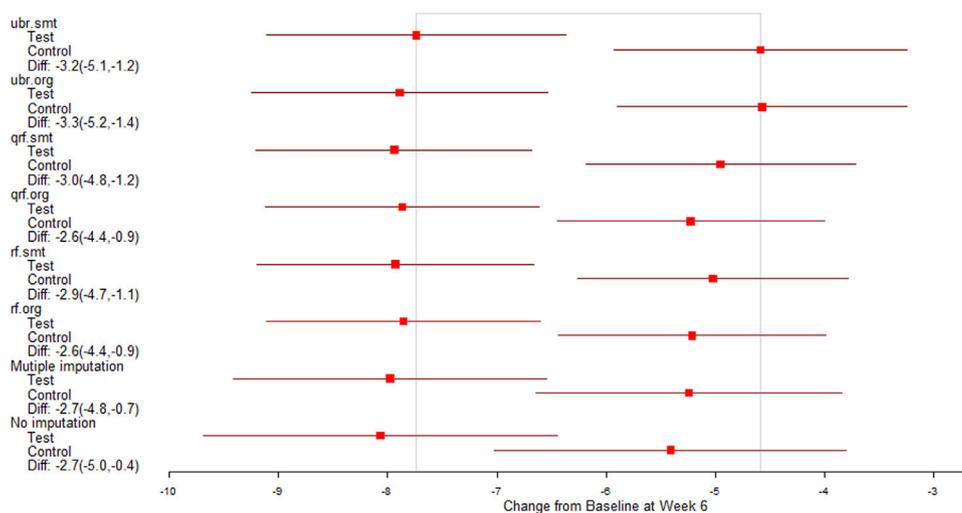


Fig. 9. Real data: forest plot for the analysis results of change from baseline in HAMD17 score at Week 6 using different methods. ubr = utility-based regression, smt = SMOTER data, org = original data, qrf = quantile random forests, rf = random forests.

pecially in the cases when the dropout rate or the efficacy pattern of dropouts are not comparable between treatment group).

5. Discussion

We aim to handle the realistic missing data scenarios (i.e., mixture of MCAR, MAR, and MNAR data) in clinical trials with con-

tinuous outcome variable. We treat MNAR as imbalanced learning task. The standard error measures are not suitable for non-unique cost learning. We propose a hybrid imbalanced learning approach that combines UBR with SMOTER. The UBR takes both the prediction error and relevance of the target variable value into account such that the areas been assigned high relevance get more focus

in the learning process. SMOTER is an effective approach to give more weights on the rare cases and also to avoid the model overfitting problem. The relevance function is a crucial part of the proposed method. The choice of the relevance function and its associated parameters should be based on the specific problem in hand and it should be pre-specified. It is inevitable to define the relevance function according to the distribution of available data, and it is also important to make a plausible assumption on the missing data (i.e., the possible locations of target variable scale where the missing data tend to occur) based on the information collected in the clinical trial. We recommend to perform sensitivity analysis with different relevance functions (and associated parameters) to check the appropriateness and robustness of the primary analysis. We evaluate the performance of proposed method in a comprehensive manner in the simulation study. When assessing the impact of missing data on the aggregated estimation, we recommend to evaluate the performance of imputation method not only in terms of the bias (like mean of imputed data) but also in terms of variance the imputed data, which is also an important element in the decision making (e.g., the decision based on the inferential statistics).

The commonly used imputation methods (like random forests and multiple imputation) do not perform as well as the proposed method and showed systematic bias in the aggregated estimation. Those methods tend to underestimate the mean and SD given the presence of MNAR in the area of high values of the target variable. A similar pattern of the systematic bias is also observed in the real data from an antidepressant clinical trial with a dropout rate of 25%. Overall, our hybrid imbalanced learning approach provides plausible prediction for all the MCAR, MAR and MNAR data and reduced the bias of missing data in the aggregated estimation. Therefore, we encourage the integration of utility-based learning strategies for handling of missing data in the analysis of clinical trials.

Limitations of this study include: (1) The use of some specific technical elements, such as QRF and SMOTER, is based on our current knowledge in this domain, and this can be further improved once new and better methods emerge; (2) To demonstrate the basic idea of utility-based regression, we look at the cross-sectional data only. However, in practice, missing data problem is more often in the longitudinal studies. Therefore, from practical point of view, an extension of the utility-based regression in the longitudinal setting is necessary.

Supporting information

All R programs for the whole workflow, datasets and outputs will be available at the website of Computer Methods and Programs in Biomedicine.

Declaration of Competing Interest

The authors have declared no conflict of interest.

Acknowledgements

The authors thank Prof. Anne-Laure Boulesteix for her valuable contribution to this work. The authors also thank the anonymous reviewers, the Associate Editor and the Editor for their generous and constructive detailed comments that helped us to improve the paper.

Appendix

Not applicable.

References

- [1] D.B. Rubin, *Inference and missing data*, *Biometrika* 63 (3) (1976) 581–592.
- [2] C.K. Enders, in: *Applied Missing Data Analysis*, Guilford Press, New York, 2010, pp. 295–301. Page.
- [3] N.H. Haliduola, F. Bretz, U. Mansmann, Missing data imputation in clinical trials using recurrent neural network facilitated by clustering and oversampling, *Biometrical J.* 64 (5) (2022) 863–882, doi:10.1002/bimj.202000393.
- [4] R.P. Ribeiro, *Utility-based Regression*, Dep. Computer Science, Faculty of Sciences - University of Porto, 2011 PhD thesis.
- [5] N.V. Chawla, K.W. Bowyer, L.O. Hall, W.P. Kegelmeyer, *Smote: synthetic minority over-sampling technique*, *J. Acad. Ind. Res.* 16 (2002) 321–357.
- [6] Torgo, L., Ribeiro, R.P. (2007). *Utility-Based Regression*. 597–604. 10.1007/978-3-540-74976-9_63.
- [7] L. Torgo, R.P. Ribeiro, B. Pfahringer, P. Branco, *Smote for regression*, in: *Progress in Artificial Intelligence*, Springer, 2013, pp. 378–389. pages.
- [8] F.N. Fritsch, R.E. Carlson, *Monotone piecewise cubic interpolation*, *SIAM J. Numer. Anal.* 17 (2) (1980) 238–246.
- [9] N. Meinshausen, *Quantile Regression Forests*, *J. Mach. Learn. Res.* 7 (2006) 983–999.
- [10] Meinshausen, N. (2017). *Quantile regression forests*, a R package available at <https://cran.r-project.org/package=quantregforest>.
- [11] A.L. Boulesteix, R. Wilson, A. Hapfelmeier, *Towards evidence-based computational statistics: lessons from clinical research on the role and design of real-data benchmark studies*, *BMC Med. Res. Methodol.* 17 (1) (2017) 138 2017Published 2017 Sep 9, doi:10.1186/s12874-017-0417-2.
- [12] A.L. Boulesteix, H. Binder, M. Abrahamowicz, W. Sauerbrei, *On the necessity and design of studies comparing statistical methods*, *Biom J.* 60 (2017), doi:10.1002/bimj.201700129.
- [13] P. Branco, L. Torgo, R. Ribeiro, *A survey of predictive modelling under imbalanced distributions*, *ACM Comput. Surv.* 1 (1) (2016) Article 1, Publication date: January 1.
- [14] M.M. Rau, S. Seitz, F. Brimiouille, E. Frank, O. Friedrich, D. Gruen, B. Hoyle, *Accurate photometric redshift probability density estimation – method comparison and application*, *Mon. Not. R. Astron. Soc.* 452 (4) (2015) 3710–3725 01 October 2015, Pages, doi:10.1093/mnras/stv1567.
- [15] Branco, P., Ribeiro, R.P., Torgo, L. (2017). *UBL: an R package for utility-based learning*.
- [16] L. Breiman, *Random Forests*, *Mach Learn* 45 (1) (2001) 5–32, doi:10.1023/A:1010933404324.
- [17] A. Liaw, M. Wiener, in: *Package “randomForest”: Breiman and Cutler’s random Forests for Classification and Regression*, 4, R Development Core Team, 2018, pp. 6–10.
- [18] R. Koenker, *Quantile Regression (Econometric Society Monographs)*, Cambridge University Press, Cambridge, 2005, doi:10.1017/CBO9780511754098.
- [19] M.S. Santos, J.P. Soares, P.H. Abreu, H.J. Araujo, *Cross-validation for imbalanced datasets: avoiding overoptimistic and overfitting approaches*, *IEEE Comput. Intell. Mag.* (2018).
- [20] H. Demirtas, B. Doganay, *Simultaneous generation of binary and normal data with specified marginal and association structures*, *J. Biopharm. Stat.* 22 (2) (2012) 223–236, doi:10.1080/10543406.2010.521874.
- [21] Amatya, A., Demirtas, H., Gao, R. (2020). *BinNor: an R package for con-current generation of binary and normal data*.
- [22] D.B. Rubin, *Multiple Imputation For Nonresponse in Surveys*, Wiley, New York, 1987.
- [23] J. Siddique, T.R. Belin, *Multiple imputation using an iterative hot-deck with distance-based donor selection*, *Stat. Med.* 27 (1) (2008) 83–102.
- [24] S. van Buuren, K. Groothuis-Oudshoorn, et al., *mice: multivariate imputation by chained equations in R*, *J. Stat. Softw.* 45 (3) (2011) 1–67 <https://www.jstatsoft.org/v45/i03/>.
- [25] London School of Hygiene and tropical medicine (2017). (<https://missingdata.lshtm.ac.uk/2017/04/28/example-dataset-from-an-antidepressant-clinical-trial/>).
- [26] D.J. Goldstein, Y. Lu, M.J. Detke, C. Wiltse, C. Mallinckrodt, M.A. Demitrack, *Duloxetine in the treatment of depression: a double-blind placebo-controlled comparison with paroxetine*, *J. Clin. Psychopharmacol.* 24 (2004) 389–399.

References

- Rubin, D. B. (1976). Inference and missing data. *Biometrika*. 63(3): 581–592.
- Heckman, J. J., (1976). The Common Structure of Statistical Models of Truncation, Sample Selection and Limited Dependent Variables and a Simple Estimator for Such Models. *Annals of Economic and Social Measurement* 5 (4): 475–92.
- Little, R. J. A., Rubin D. B. (1987). *Statistical Analysis with Missing Data*. New York: Wiley. XIV+278 pp.
- Little, R. J. A. (1993). Pattern-mixture models for multivariate incomplete data. *J Am Stat Assoc*. 88: 125–134.
- Little, R. J. A. (1994). A Class of Pattern-Mixture Models for Normal Incomplete Data. *J Biometrika*. 81(3): 471-483. doi:10.2307/2337120
- Little, R. J. A. (1995). Modeling the drop-out mechanism in repeated-measures studies. *J Am Stat Assoc*. 90: 1113–1121.
- Enders, C. K. (2010). *Applied Missing Data Analysis*. New York, Guilford Press. Page 295-301.
- Boulesteix A. L., Wilson R., Hapfelmeier A. (2017). Towards evidence-based computational statistics: lessons from clinical research on the role and design of real-data benchmark studies. *BMC Med Res Methodol*. 2017;17(1): 138. Published 2017 Sep 9. doi:10.1186/s12874-017-0417-2
- Boulesteix, A. L., Binder, H., Abrahamowicz, M., Sauerbrei, W. (2017). On the necessity and design of studies comparing statistical methods. *Biometrical Journal*. 60. 10.1002/bimj.201700129.
- Genolini, C., Falissard, B. (2011). KmL: A Package To Cluster Longitudinal Data. *Computer Methods and Programs in Biomedicine*. 104(3): e112–21. doi: 10.1016/j.cmpb.2011.05.008.
- Falbel, D., Allaire, J., Chollet, F., RStudio, Google, Tang, Y., Van Der Bijl, W., Studer, M., Keydana, S. (2015). Keras. GitHub. Retrieved from <https://github.com/fchollet/keras>
- European Medicines Agency. (2011). *Guideline on Missing Data in Confirmatory Clinical Trials*.
- National Research Council of the National Academies. (2010). *The Prevention and Treatment of Missing Data in Clinical Trials*. Washington D.C.: National Academies Press.
- Breiman, L. (2001). Statistical Modeling: The Two Cultures. *Statistical Science*, Vol. 16, No. 3, 199-215. doi:10.1214/ss/1009213726.
- Weiss, G. M. (2013). *Imbalanced Learning: Foundations, Algorithms and Applications*. Wiley-IEEE Press.
- Gower, J. C. (1971). A General Coefficient of Similarity and Some of Its Properties. *Biometrics*, vol. 27, no. 4, pp. 857–871.
- Rumelhart, D. E., Hinton, G. E., Williams, R. J.. (1986). Learning representations by back-propagating errors. *Nature*. 323 (6088): 533–536.

Schmidhuber, J. (1993). *Netzwerkarchitekturen, Zielfunktionen und Kettenregel* (Network architectures, objective functions, and chain rule). Habilitation (postdoctoral thesis - qualification for a tenure professorship), Institut für Informatik, Technische Universität München, 1993.

Claesen, M., De Moor, B. (2015). *Hyperparameter Search in Machine Learning*. arXiv:1502.02127 [cs.LG].

Torgo, L., Ribeiro, R. P. (2007). Utility-Based Regression. 597-604. 10.1007/978-3-540-74976-9_63.

Ribeiro, R. P. (2011). *Utility-based Regression*. PhD thesis, Dep. Computer Science, Faculty of Sciences - University of Porto.

Torgo, L., Ribeiro, R. P., Pfahringer, B., and Branco, P. (2013). Smote for regression. In *Progress in Artificial Intelligence*, pages 378-389. Springer.

Chawla, N.V., Bowyer, K.W., Hall, L.O., Kegelmeyer, W.P. (2002). Smote: Synthetic minority over-sampling technique. *JAIR* 16 (2002) 321-357.

Amatya, A., Demirtas, H., Gao, R. (2020). BinNor: An R package for con-current generation of binary and normal data.

Branco, P., Torgo, L., Ribeiro, R. (2016). A Survey of Predictive Modelling under Imbalanced Distributions. *ACM Computing Surveys*, Vol. 1, No. 1, Article 1, Publication date: January 1.

Branco, P., Ribeiro, R.P., Torgo, L. (2017). UBL: an R package for Utility-based Learning.

van Buuren S., Groothuis-Oudshoorn K, et al. (2011). "mice: Multivariate Imputation by Chained Equations in R." *Journal of Statistical Software*, 45(3), 1-67. <https://www.jstatsoft.org/v45/i03/>.

Liaw, A., Wiener, M. (2018). Package "randomForest": Breiman and Cutler's random forests for classification and regression. R Development Core Team. 4. 6-10.

Meinshausen, N. (2006). Quantile Regression Forests. *Journal of Machine Learning Research*. 7. 983-999.

Meinshausen, N. (2017). Quantile regression forests, a R package available at <https://cran.r-project.org/package=quantregforest>.

Siddique, J., Belin, T.R. (2008). Multiple imputation using an iterative hot-deck with distance-based donor selection. *Statistics in medicine*, 27, 1, 83–102.

Acknowledgements

First of all, I would like to thank my main supervisor Prof. Dr. Ulrich Mansmann. I was deeply impressed by Prof. Mansmann's openness to new ideas and his broad and deep knowledge in biostatistics and related fields, which fundamentally made this research possible (since this research is kind of "bridging" the biostatistics and machine learning fields for the missing data problem). I learnt a lot from Prof. Mansmann on how to conduct scientific research, e.g., how to think through the overall research paradigm, how to design a simulation study, how to generalize an idea into a new approach etc. I would like to say a big THANK YOU for his great supports in the past few years. Without his detailed guidance and advice in each individual step over the whole research process, this research work will never be completed.

I would like to say "Thank You" to my second supervisor Prof. Dr. Anne-Laure Boulesteix (from IBE at LMU Munich) for her valuable contribution to this work, specifically for her advice on the simulation study. Inspired by her "evidence-based computational statistics" concept, our proposed methods are evaluated in the practically relevant simulation data, which made our approaches even more convincing.

I would like to thank my third supervisor Prof. Dr. Frank Bretz (from the Medical University of Vienna), for his valuable contribution to this work, specifically during the scientific paper writing process. Prof. Bretz made me aware of the gap between the readers of the two fields (i.e., biostatistics and machine learning) and he suggested to take steps back to explain machine learning technicalities using a language that are more acceptable to the readers in the field of biostatistics. He provided constructive inputs on the structure and language of the published papers.

Of course, none of these could have happened without the support of my family. My father Haliduola and my mother Kulash, who offered their encouragement through phone calls every week. A special big thanks to my wife Gulzat, who was willing to selflessly take care of the main household chores and take good care of our two children (Alex and Eldos) in the past years.