# ORIGINAL ARTICLE

# Real-time imputation of missing predictor values improved the application of prediction models in daily practice

Steven Willem Joost Nijman[a,1,*], T. Katrien J. Groenhof[a,1], Jeroen Hoogland[a], Michiel L. Bots[a],
Menno Brandjes[b], John J.L. Jacobs[b], Folkert W. Asselbergs[c,d,e], Karel G.M. Moons[a],
Thomas P.A. Debray[a,d]

[a]*Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht University, Utrecht, the Netherlands*
[b]*LogiqCare, Ortec B.V., Zoetermeer, the Netherlands*
[c]*Division Heart & Lungs, Department of Cardiology, University Medical Center Utrecht, Utrecht University, Utrecht, the Netherlands*
[d]*Institute of Cardiovascular Science, Faculty of Population Health Sciences, University College London, London, UK*
[e]*Health Data Research UK, Institute of Health Informatics, University College London, London, UK*

## Abstract

**Objectives:** In clinical practice, many prediction models cannot be used when predictor values are missing. We, therefore, propose and evaluate methods for real-time imputation.

**Study Design and Setting:** We describe (i) mean imputation (where missing values are replaced by the sample mean), (ii) joint modeling imputation (JMI, where we use a multivariate normal approximation to generate patient-specific imputations), and (iii) conditional modeling imputation (CMI, where a multivariable imputation model is derived for each predictor from a population). We compared these methods in a case study evaluating the root mean squared error (RMSE) and coverage of the 95% confidence intervals (i.e., the proportion of confidence intervals that contain the true predictor value) of imputed predictor values.

**Results:** −RMSE was lowest when adopting JMI or CMI, although imputation of individual predictors did not always lead to substantial improvements as compared to mean imputation. JMI and CMI appeared particularly useful when the values of multiple predictors of the model were missing. Coverage reached the nominal level (i.e., 95%) for both CMI and JMI.

**Conclusion:** Multiple imputations using either CMI or JMI is recommended when dealing with missing predictor values in real-time settings. © 2021 The Authors. Published by Elsevier Inc. This is an open access article under the CC BY license (http://creativecommons.org/licenses/by/4.0/).

*Keywords:* Missing data; Multiple imputations; Real-time imputation; Prediction; Computerized decision support system; Electronic health records

[1] Contributed equally.

* Corresponding author. Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Heidelberglaan 100, 3584 CX Utrecht, the Netherlands. Tel.: +31-(0)88-75 68012; fax: +31(0)887568099.
*E-mail address:* s.w.j.nijman@umcutrecht.nl (S.W.J. Nijman).

**What is new?**

**Key findings**
- Multiple imputation approaches can be adapted without much difficulty to allow for real-time imputation of missing predictor variables.

- Both conditional modeling imputation (CMI??) and joint modeling imputation (JMI??) give more accurate estimates of missing predictor values when compared to mean imputation.

**What this adds to what was known?**
- Imputation of missing predictor values does not require ''live'' access to a source dataset. Simple population characteristics (such as the mean and covariance) can be used to generate imputations that are tailored to a specific individual.

**What is the implication, and what should change now?**
- Real-time multiple imputations using either CMI or JMI should be made available in clinical practice (e.g., via a computerized decision support system) to support guideline-recommended use of prediction models and to be more transparent about uncertainty

- When developing or validating a prediction model, researchers should report the mean and covariance of the study population, as this information can directly be used to impute missing values in routine care.

# 1. Introduction

In present-day medical practice, characterized by an aging population, multimorbidity, and high complexity of diseases, attention has grown toward personalized medicine aiming to administer the most applicable treatment to the individual patient given their risk profile [1—3]. In cardiovascular disease management, guidelines advocate the use of prediction models to assess the patients' risk of developing a certain cardiovascular disease to guide treatment decision making [1]. For integrating risk-guided care in daily practice, technological solutions such as computerized decision support systems (CDSS) are increasingly developed [4,5]. Using predictor values directly extracted from the electronic health record (EHR), CDSS can provide an immediate risk assessment of each encountered patient at a glance [6—8].

The use of prediction models in daily practice in an individual patient requires real-time availability of the patient's values of the predictors in the model. Most prediction models cannot provide a risk estimate in the presence of missing predictor values, which hampers implementation and may ultimately limit guideline adherence [9]. Therefore, predictor values should be measured and registered (e.g., in the Electronic Health Record; EHR) in such a way that they are available in real-time. Yet, routine clinical care data is often incomplete because certain measurements are deemed unnecessary, time-consuming, or expensive, or because they cannot directly be extracted from the EHR (e.g., registered as free text) [10].

Missing data is a well-known challenge in (medical) research, for which several scalable solutions exist [11]. Multiple imputations by chained equations has often been recommended to handle missing data in a research setting where data from multiple patients are available for study analysis purposes [12,13]. This approach, however, is not directly applicable when applying a prediction model in real-time to a single patient in the consulting room. In particular, the models used for imputation cannot be generated ''live'' in clinical practice, and therefore, need to be derived elsewhere and beforehand [14].

One option is to replace missing predictor values by their respective mean/median, which, in turn, is estimated from another data set or training sample [15,16]. While straightforward to implement, mean imputation may be insufficient when the predictor with missing values is a strong predictor or exhibits large variability such that assigning an overall mean may lead to the less predictive accuracy of the prediction model and misinformed treatment decisions. Mean imputation does not distinguish between patients and may, therefore, likely impute values that are unrealistic given the patient's observed predictor values. Also, mean imputation obfuscates any uncertainty about the imputed values.

To address these issues, we expand on two well-known methods that may also be used in real-time imputation of missing predictor values [14]: joint modeling imputation (JMI) [17] and conditional modeling imputation (CMI, known for its common use in multiple imputations by chained equations) [13]. As opposed to mean imputation, these methods are able to incorporate the relationship between multiple patient characteristics, and therefore, allow imputations to be adjusted for observed patient specific characteristics. Similar to mean imputation, these relations can be learned from training data, and in real time, applied on new patients that are not part of the training sample. Additionally, both methods allow for multiple imputations to be estimated, reflecting the uncertainty with respect to the imputed value.

Using a real-world example and empirical data set on cardiovascular risk prediction, we compared the accuracy and usability of three imputation methods (mean imputation, JMI, and CMI) to deal with missing values of predictors in the prediction model in real time. Although mean imputation has been known to be problematic during model development, it was chosen as a comparison due to its common use during model application in routine clinical practice or in decision support [18—21].

## 2. Methods

### 2.1. Imputation methods

For facilitating the live imputation of missing values in routine care, it is essential to obtain information on the distribution of the target population. This summary information can, for instance, be derived in an epidemiologic (e.g., cohort) study and then be utilized for training live imputation models. A key constraint given is that after being trained, all methods are independent and stand-alone, which means that they can directly be used for live imputation in a new, single, patient without requiring the need for any additional procedures.

The three methods under evaluation are mean imputation, joint modeling imputation (JMI), and conditional modeling imputation (CMI) [13,14,17]. All methods were implemented in R and facilitate live imputation of missing values in individual patients. Source code is available from the supplementary information (Appendix D).

#### 2.1.1. Mean imputation

The training sample is used to derive the means of all predictors in the model (Fig. 1). Missing predictor values are then imputed by their respective mean (or proportion in the case of binary variables). This method is relatively straightforward to implement, and can be extended to subgroup-specific means (i.e., creating subdivisions based on certain parameters of a population of which multiple means are calculated).

#### 2.1.2. Joint modeling imputation

The training sample is used to derive the means and covariance of all predictor variables (Fig. 2). It is assumed that all predictor variables of the training sample are normally distributed, such that imputations for an individual patient can directly be generated from the mean and covariance of the training sample and the observed predictor values [14,17]. In contrast to overall mean imputation, the use of covariances between all predictors incorporates the relation between the predictors, and therefore, allows imputations to be tailored to an individual patient's own characteristics. A more detailed description is provided in Appendix A [14].

#### 2.1.3. Conditional modeling imputation

The training sample is used to derive a flexible (e.g., regression) model for each predictor (as dependent variable) with all other predictor variables as independent variables (Fig. 3). These models describe the conditional distribution of each predictor and usually need to be estimated using a Gibbs sampling procedure (as predictor values may also be missing in the training sample). Due to the flexible nature of these conditional models, it is no longer assumed that predictor variables of the training sample are normally distributed (as does JMI). For instance, a logistic regression model can be used to estimate the conditional distribution of a binary predictor variable (e.g., current smoker). Subsequently, when the smoking status for a new patient is unknown, the logistic regression model can be used to generate a probability that they are a current smoker. This probability can directly be used as an imputed value (in case only 1 imputation is needed). Alternatively, if multiple imputations are required, a Bernoulli distribution (with the aforementioned probability) can be used to sample multiple (discrete) values for the patient's current smoking status. If multiple predictor values are missing, the conditional models need to be used successively using an iterative Monte Carlo procedure (Appendix A).

### 2.2. Simulation study

Cardiovascular disease prevention is an example of a setting where risk-guided management of predictors—smoking, blood pressure, cholesterol—is common practice
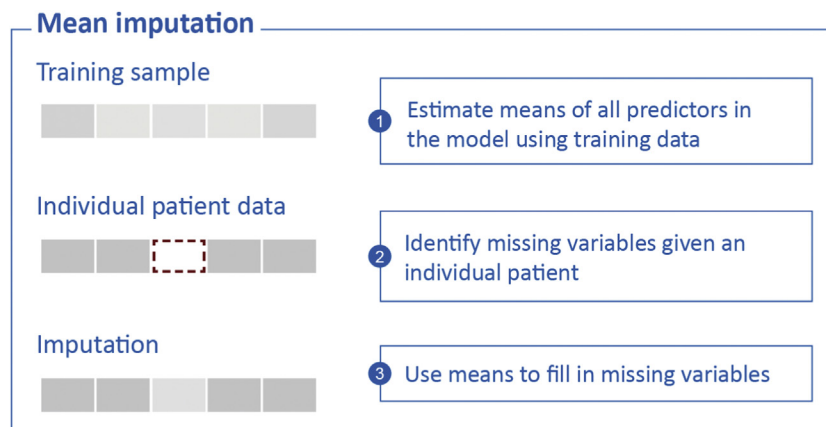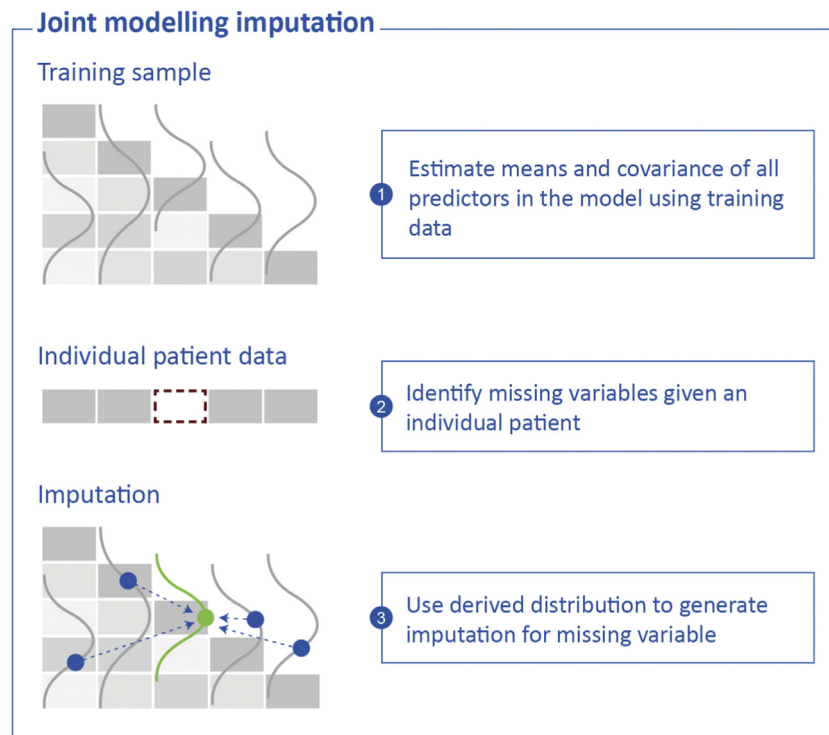


**Fig. 1.** Mean imputation.

**Fig. 2.** Joint modeling imputation.

[22]. Numerous risk prediction models have been developed, and the (international) guidelines advocate the use of risk classification to inform treatment decisions [23,24]. These models are typically implemented in a CDSS, where a patient's characteristics of the predictors can be entered manually or are automatically retrieved from the patient's EHR [4,6,25].

For this study, we used a data set of the ongoing Utrecht Cardiovascular cohort initiative (UCC). This cohort includes all patients who come for a first-time visit to the Center for Circulatory Health at the UMC Utrecht for the evaluation of a symptomatic vascular disease or an asymptomatic vascular condition. A minimum set of predictors, according to the Dutch Cardiovascular Risk Management Guidelines, is collected in all patients. No data on outcomes (i.e., time-to-event data) was recorded. UCC has been approved by the Institutional Review Board of the UMC Utrecht (Biobank Ethics committee). For the present analyses, an anonymized dataset was used of the UCC cohort up to November 2018 [26,22].

The sample consisted of 3,880 patients with information on 23 variables, measured during the patient's visit (Table 1 and Appendix B). For clarity of exposition, we completed this dataset using all 23 variables in k-nearest neighbor imputation, which aggregates the values of the $k$ nearest neighbors to an imputation [27]. In practice, regular multiple imputation techniques can be used in case of incomplete training data.

For evaluating the quality of the three selected imputation methods in individual patients, a leave-one-out-cross-validation (LOOCV) procedure was used in the completed UCC dataset. In LOOCV, all but one patient are used as the training sample from which the overall mean or proportion (method 1) or imputation models (methods 2 and 3) are derived (Fig. 4). In the remaining hold-out patient, missing values are introduced for one or more predictor variables. As we apply each scenario to each patient exactly once, the missing data mechanism is essentially missing-completely-at-random (MCAR) [18]. The summary information from the training sample is then used to impute the missing predictor values in the hold-out patient. For CMI and JMI, we generated 50 imputations for each missing predictor value. This process is repeated until all patients have been taken from the dataset exactly once.

We consider eight scenarios where missing values occur for one predictor variable, and eight scenarios where multiple predictor variables are simultaneously missing (Fig. 5). A detailed description of how the scenarios were selected and of the R code are listed in Appendix C and D, respectively.

### 2.3. Measures of performance

To evaluate the performance of the three imputation methods, we used four performance metrics:

1. We calculated the root-mean-squared error (RMSE) between the average of the multiple imputed

## Conditional modelling imputation



**Fig. 3.** Conditional modeling imputation.

predictor values (i.e., 50 imputations) and the true, original (i.e., before the simulation of missing) predictor value to evaluate the accuracy of the imputations. The RMSE is a performance measure that aggregates error due to bias and variability. Generally, an RMSE of zero means perfect imputation and an increasing RMSE means decreasing performance of the imputation. The clinical relevance of an RMSE depends on the natural range of the predictor. For example, an RMSE of 0.5 is large for LDL-c (mean

3.0 SD 1.3 mmol/L) but not for SBP (mean 143 SD 24 mmHg).

2. For each hold-out patient, we assessed whether the original predictor value was in the 95% confidence interval around the imputed predictor value. Subsequently, we calculated the proportion of confidence intervals that consisted the original value (coverage). For a 95% CI, the coverage should ideally be equal to 95% [28]. A lower coverage translates to imputed predictor values that are too precise (which, in turn,

**Table 1.** Descriptive statistics (after imputation)

| Variables (unit) | Part of missing data scenarios | Mean (SD) or n/total (%)[a] | Original missing % |
|---|---|---|---|
| Age (yr) | No | 61.7 (18.2) | 0.00 |
| Sex (1 = female; 0 = male) | No | 1,987/3,880 (51.2) | 0.00 |
| Smoking (1 = yes; 0 = no) | No | 363/3,880 (9.4) | 24.07 |
| SBP (mmHg) | Yes | 142.8 (24.2) | 10.54 |
| TC (mmol/L) | Yes | 5.1 (1.2) | 24.54 |
| LDL-c (mmol/L) | Yes | 3.1 (1.3) | 26.01 |
| HDL-c (mmol/L) | No | 1.4 (0.4) | 25.39 |
| eGFR (mL/min/1.73 m$^2$) | Yes | 81.8 (24.6) | 15.98 |
| History of CVD (1 = yes; 0 = no) | Yes | 1,971/3,880 (50.8) | 23.45 |
| History of PAD (1 = yes; 0 = no) | No | 335/3,880 (8.6) | 23.45 |
| History of CHD (1 = yes; 0 = no) | No | 591/3,880 (15.2) | 23.45 |
| History of CHF (1 = yes; 0 = no) | No | 284/3,880 (7.3) | 23.45 |
| History of CVA (1 = yes; 0 = no) | No | 579/3,880 (14.9) | 23.45 |
| History of DM (1 = yes; 0 = no) | No | 607/3,880 (15.6) | 23.45 |
| Polyvascular disease | No | 0.6 (0.7) | 23.45 |
| Number of medications | No | 0.8 (1.7) | 27.24 |
| BP lowering medication (1 = yes; 0 = no) | No | 705/3,880 (18.2) | 27.24 |
| Statin (1 = yes; 0 = no) | No | 415/3,880 (10.7) | 27.24 |
| HbA1c (mmol/mol) | No | 40 (10.7) | 26.37 |
| Years after first CVD (yr) | Yes | 4.6 (8.1) | 26.21 |
| Diabetes (1 = yes; 0 = no) | Yes | 755/3,880 (19.5) | 8.12 |
| Diabetes duration (yr) | No | 11.3 (7.3) | 86.11 |
| Pulse pressure (mmHg) | No | 61.7 (18.9) | 10.54 |

*Abbreviations:* SBP, systolic blood pressure; TC, total cholesterol; LDL-c, low-density lipoprotein cholesterol; HDL-c, high-density lipoprotein cholesterol; eGFR, estimated glomerular filtration rate according to the CKD epi formula; CVD, cardiovascular disease; PAD, peripheral artery disease; CHD, coronary heart disease; CHF, chronic heart failure; CVA, cerebrovascular accident; DM, diabetes mellitus; BP, blood pressure; HbA1c, glycated hemoglobin.

[a] After KNN-imputation.

may lead to estimates of predicted risk that are too precise), whereas a coverage above 95% indicates that imputed predictor values are too imprecise [13]. We assessed coverage only for continuous predictor variables.

3. We assessed the effect on treatment decision support for blood pressure in patients with manifest cardiovascular disease (*n* = 1,971) to evaluate the clinical implications of the imputed predictor values. Guidelines indicate that all patients with a history of CVD should receive blood pressure-lowering treatment when their blood pressure is higher than 140/90 mmHg [1,22]. We adopted the LOOCV approach and set values for SBP missing in the hold-out patient. Subsequently, we imputed the missing value and compared the treatment decision for the true value with the treatment decision for the imputed value (SBP < > 140 mmHg). Afterward, we calculated the sensitivity, specificity, positive predictive value, and negative predictive value. Also, we illustrated the importance of reporting confidence intervals based on imputed values to inform the discussion around treatment commencement.

4. We compared the risk predictions that were obtained in the absence of missing values (i.e., in the original data) with the risk predictions that are based on imputations to evaluate the impact of the imputed values on the precision of predicted risk. Ideally, the predictions that are based on imputed values should have a similar distribution as the predictions that are derived from the complete original data. To explore any deviation, we assessed the interquartile range of predicted risk for a single missing predictor scenario and a multiple missing predictor scenario. Rather than developing a new prediction model ourselves, we applied the previously developed SMART prediction model for the risk of 10-year recurrent vascular disease as reported in the original development study [23][.
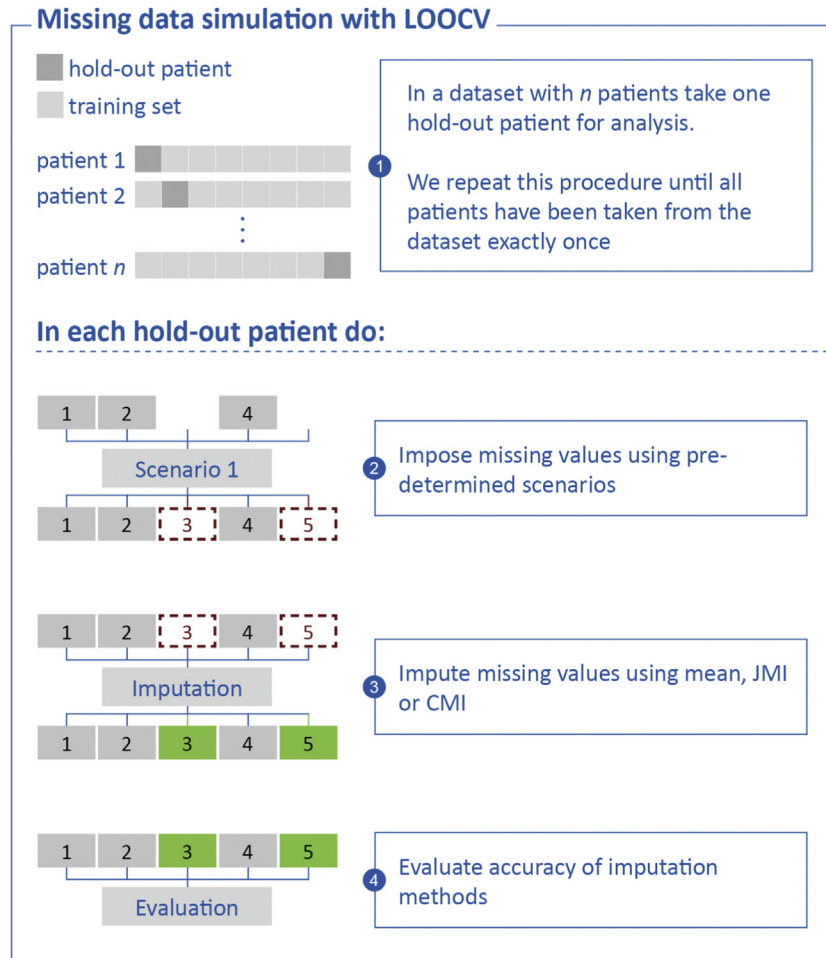
**Fig. 4.** Missing data simulation procedure.

The prediction model includes 11 variables: age, sex, current smoker, SBP, diabetes, history of cerebrovascular disease, aortic aneurysm or peripheral vascular disease, polyvascular disease, HDL-cholesterol, and total cholesterol.

## 3. Results

### 3.1. Root-mean-squared error

With the exception of smoking, all predictor variables in single missing predictor scenarios had a lower RMSE when using JMI or CMI as compared to mean imputation (Table 2). For most multiple missing predictor scenarios, the RMSE is consistently lower when using JMI or CMI as compared to mean imputation. The exceptions being the history of CVD and smoking. Performance diminished as more variables were missing. For example, the RMSEs of years after the 1st CVD event are 6.30 and 6.26 for JMI and CMI respectively when univariately missing, while mean imputation has an RMSE of 8.06. When additional

variables (e.g., SBP, history of CVD, and smoking) are missing, the RMSE for years after the 1st CVD event for both JMI and CMI increases to 7.58 and 7.84, respectively.

### 3.2. Coverage rate

For JMI, the coverage reached nominal levels for all single missing predictor scenarios and multiple missing predictor scenarios (Table 3). For CMI, the coverage reached nominal levels for all single missing predictor scenarios and multiple missing predictor scenarios. For mean imputation, coverage was 0% by definition for all imputed predictors because no uncertainty is taken into account.

### 3.3. Clinical decision accuracy

When assessing the treatment decision for blood pressure management according to the prevailing clinical guidelines (see above), we selected 1,971 out of the total 3,880 patients with manifest cardiovascular disease. We found that 1,134 patients (57.53%) should be treated. However, when blood pressure values were set to miss, the
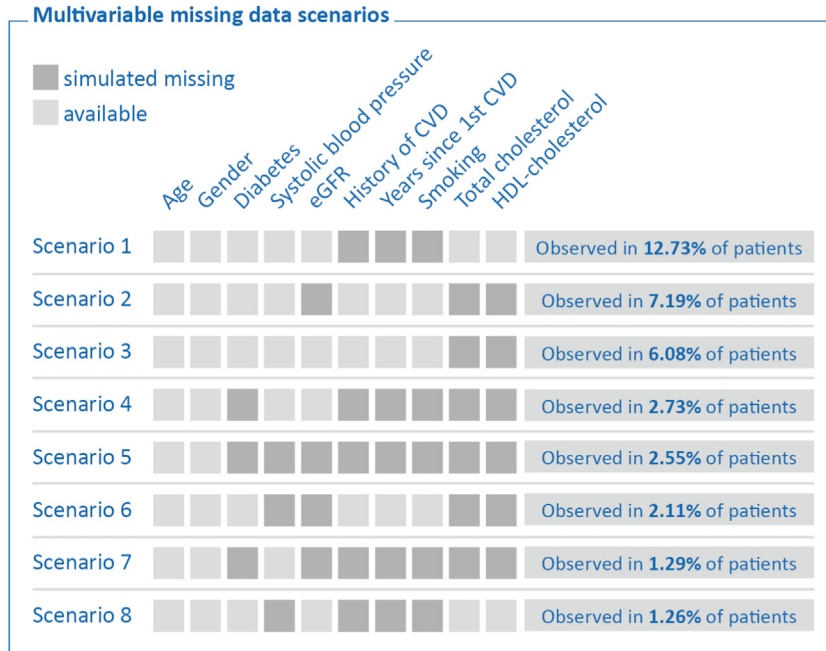
**Fig. 5.** Multivariable missing data scenarios.

overall mean imputed value was 142 mmHg (Table 1), which is just above the treatment threshold of 140 mmHg. As a result, everyone would have been treated when adopting overall mean imputation, such that 837 patients (42.47%) would have been treated unnecessarily. When adopting JMI or CMI, only 16.08% or, respectively, 15.98% of patients would have been treated unnecessarily (Table 4). Hence, the imputation of missing blood pressure values using CMI or JMI was more than adequate than the mean imputation in terms of decision making.

To illustrate the importance of measuring uncertainty, we provided an example in which we compare the use of imputation in a real-life situation (Table 5). In the example, a patient with an imputed SBP of 144 mmHg was given an indication for blood pressure-lowering treatment according to the guidelines [1]. However, given that the uncertainty around the imputed predictor value crosses the treatment line of 140 mmHG (scenario A), there is reasonable doubt this imputation is too uncertain to be used for treatment decision making.

### 3.4. Effect on risk predictions

The predicted risks, given each method, did not seem to deviate much from the originally predicted risk, given the complete data (Table 6). When assessing the single missing predictor scenario, there was a difference between overall mean imputation (median difference of −1.713% to the originally predicted risk) and the combination of JMI and CMI (median difference of respectively 0.301% and 0.399% to the originally predicted risk). Further, we found

that predicted risks for mean imputation were more similar when compared to the complete data (standard deviation = 15.12 vs. the reference of 18.91). In contrast, the standard deviations of JMI and CMI were 17.87 and 17.86, respectively.

In the multiple missing predictor scenario, there was a similar difference between mean imputation (median difference of −2.064% to the originally predicted risk) and JMI and CMI (median difference of respectively 0.375% and 0.390% to the originally predicted risk). With multiple missing predictors, the predicted risks for mean imputation were again more similar than the predicted risk given the complete data (standard deviation = 14.42 vs. the reference of 18.91). The standard deviations of JMI and CMI were 17.67 and 17.68, respectively.

The difference between mean imputation and both JMI and CMI is especially apparent in high-risk patients (i.e., 75% IQR) where mean imputation, as expected, underestimates the risk. This is because mean imputation pulls the risk predictions of patients with missing values toward the prediction for an ''average'' patient. As such, JMI and CMI perform much better with regards to their impact on prediction in high-risk patients when compared to mean imputation.

## 4. Discussion

This project described the development and performance of three imputation methods to handle missing data on an individual patient level in real-life clinical decision making.

**Table 2.** RMSE for each combination of, individual or multiple, missing predictor values

**Single missing variable scenarios**

| Variable name (type of data) | Diabetes (binary) | SBP (continuous) | EGFR (continuous) | History of CVD (binary) | Years after 1st CVD (continuous) | Smoking (binary) | Total cholesterol (continuous) | HDL-cholesterol (continuous) |
|---|---|---|---|---|---|---|---|---|
| Mean imputation[a] | 0.40 | 24.24 | 24.56 | 0.50 | 8.06 | 0.29 | 1.24 | 0.36 |
| JMI | 0.17 | 22.31 | 19.60 | 0.39 | 6.30 | 0.30 | 1.19 | 0.34 |
| CMI | 0.21 | 22.29 | 19.69 | 0.39 | 6.26 | 0.29 | 1.19 | 0.34 |

**Multiple missing variables scenarios**

| Method | Scenario | Diabetes | SBP | eGFR | History of CVD | Years after 1st CVD | Smoking | Total cholesterol | HDL-cholesterol |
|---|---|---|---|---|---|---|---|---|---|
| JMI | 1 | | | | 0.46 | 7.59 | 0.30 | | |
| CMI | 1 | | | | 0.51 | 7.78 | 0.29 | | |
| JMI | 2 | | | 19.68 | | | | 1.19 | 0.35 |
| CMI | 2 | | | 19.69 | | | | 1.20 | 0.35 |
| JMI | 3 | | | | | | | 1.19 | 0.33 |
| CMI | 3 | | | | | | | 1.19 | 0.35 |
| JMI | 4 | 0.17 | | | 0.48 | 7.65 | 0.30 | 1.22 | 0.35 |
| CMI | 4 | 0.20 | | | 0.50 | 7.83 | 0.28 | 1.21 | 0.35 |
| JMI | 5 | 0.17 | 22.62 | 19.86 | 0.47 | 7.66 | 0.30 | 1.23 | 0.35 |
| CMI | 5 | 0.21 | 22.48 | 19.87 | 0.51 | 7.86 | 0.29 | 1.22 | 0.35 |
| JMI | 6 | | 22.45 | 19.61 | | | | 1.19 | 0.34 |
| CMI | 6 | | 22.50 | 19.59 | | | | 1.20 | 0.34 |
| JMI | 7 | 0.17 | | 19.83 | 0.48 | 7.69 | 0.30 | 1.22 | 0.35 |
| CMI | 7 | 0.21 | | 19.75 | 0.50 | 7.84 | 0.29 | 1.23 | 0.35 |
| JMI | 8 | | 22.36 | | 0.46 | 7.58 | 0.30 | | |
| CMI | 8 | | 22.35 | | 0.51 | 7.84 | 0.29 | | |

*Abbreviations:* JMI, joint modeling imputation; CMI, conditional modeling imputation; SBP, systolic blood pressure; eGFR, estimated glomerular filtration rate according to the CKD epi formula; CVD, cardiovascular disease.

The RMSE should ideally be 0. Multiple missing predictor scenarios: (1) history of CVD, years after 1st CVD event & smoking, (2) eGFR, total cholesterol & HDL-cholesterol, (3) total cholesterol & HDL-cholesterol, (4) all variables but SBP & eGFR, (5) all variables, (6) SBP, eGFR, total cholesterol & HDL-cholesterol, (7) all variables but SBP and (8) SBP, history of CVD, years after 1st CVD event and smoking.

[a] Mean imputation is only included in the single missing variable scenarios, as the performance of the model when multiple variables are missing, is equivalent.

As expected, both JMI—using draws from a normal distribution constructed from means and covariance in the training sample and observed values in the patient—and CMI—using a conditional distribution of each variable based on regression models fitted on all other variables, were more accurate and showed better coverage as compared to mean imputation, resulting in fewer inappropriate treatment decisions and lower impact on predicted risk.

The accuracy measures—RMSE, coverage, and clinical decision accuracy—were comparable for JMI and CMI. Hence, both methods can be used for generating live imputations in routine care. Based on usability, we recommend JMI, as its implementation in decision support systems is fairly straightforward and only requires information on the mean and covariance of the target population. Although its assumption of multivariate normality may be unrealistic

for real-life clinical data, simulation studies have demonstrated that this rarely affects the performance of imputation [29–31].

Previous studies on imputation methods to handle missing data on an individual patient level have focused on the impact of missing values on the performance of a prediction model and evaluated the use of mean imputation, as well as the (re)development of a simplified prediction model [15,16]. Mean imputation was recommended due to its applicability in practice and relatively good performance compared to other models but was considered insufficient when strong predictors were missing. For this reason, our proposed multiple imputation models appear particularly relevant when strong or multiple predictors are missing. This was confirmed in our simulation study: RMSE and coverage did not deteriorate much with the increasing number of predictor values that were

**Table 3.** Coverage for each combination of individual or multiple imputations

| Coverage: Single missing variable scenarios | | | | | |
|---|---|---|---|---|---|
| | **SBP** | **eGFR** | **Years after1st CVD** | **Total cholesterol** | **HDL-cholesterol** |
| JMI | 0.945 | 0.948 | 0.952 | 0.952 | 0.950 |
| CMI | 0.945 | 0.948 | 0.954 | 0.953 | 0.948 |

| Coverage: Multiple missing variable scenarios | | | | | | |
|---|---|---|---|---|---|---|
| Method | Scenario | SBP | eGFR | Years after1st CVD | Total cholesterol | HDL-cholesterol |
| JMI | 1 | | | 0.951 | | |
| CMI | 1 | | | 0.948 | | |
| JMI | 2 | | 0.947 | | 0.951 | 0.951 |
| CMI | 2 | | 0.946 | | 0.955 | 0.949 |
| JMI | 3 | | | | 0.949 | 0.951 |
| CMI | 3 | | | | 0.950 | 0.949 |
| JMI | 4 | | | 0.951 | 0.950 | 0.951 |
| CMI | 4 | | | 0.949 | 0.952 | 0.952 |
| JMI | 5 | 0.944 | 0.947 | 0.951 | 0.952 | 0.951 |
| CMI | 5 | 0.948 | 0.948 | 0.946 | 0.953 | 0.953 |
| JMI | 6 | 0.945 | 0.950 | | 0.951 | 0.948 |
| CMI | 6 | 0.948 | 0.948 | | 0.949 | 0.949 |
| JMI | 7 | | 0.950 | 0.951 | 0.951 | 0.951 |
| CMI | 7 | | 0.947 | 0.950 | 0.948 | 0.951 |
| JMI | 8 | 0.945 | | 0.952 | | |
| CMI | 8 | 0.945 | | 0.950 | | |

*Abbreviations:* JMI, joint modeling imputation; CMI, conditional modeling imputation; SBP, systolic blood pressure; eGFR, estimated glomerular filtration rate according to the CKD epi formula; CVD, cardiovascular disease.

The presented values depict the coverage of 95% confidence intervals (hence, the reference value is 0.95). Multiple missing predictor scenarios: (1) history of CVD, years after 1st CVD event & smoking, (2) eGFR, total cholesterol & HDL-cholesterol, (3) total cholesterol & HDL-cholesterol, (4) all variables but SBP & eGFR, (5) all variables, (6) SBP, eGFR, total cholesterol & HDL-cholesterol, (7) all variables but SBP and (8) SBP, history of CVD, years after 1st CVD event and smoking.

simultaneously missing for the individual patient. Because of the way missing data was introduced, it is noted that our simulations were not able to distinguish between various mechanisms by which data can be missing, for example, data that is missing at random (MAR) vs. data that is missing-completely-at-random (MCAR??) [18].

Furthermore, because the described imputation methods can accommodate numerous patient characteristics that are not necessarily disease-specific, they are highly scalable to other settings and populations. However, it is likely that some local tailoring is necessary when imputation models are derived from specific studies or settings that do not fully match the intended target population. For JMI, the means and covariances could, for instance, simply be replaced by their respective values in a local "training" sample. For CMI, the regression coefficients can be revised using recently described updating methods [32]. When the (local) training data are affected by missing predictor values, advanced methods exist to estimate the mean and the covariance [33]. All methods can be potentially incorporated within an EHR based computerized decision support system and generate imputations based on observed data from individual patients extracted from the EHR. Evidently, before implementing imputation models in clinical practice, it is of the utmost importance to assess their validity, likely impact on treatment decisions, patient outcomes, as well as any practical, security, and ethical constraints.

Although multiple imputations offer a computational framework to account for missing values, we always recommend optimizing data collection first and avoid having missing values: clinical decision making should never be based solely on imputed values. However, imputed values can serve as a proxy for prior risk, setting an indication for more (advanced) diagnostic tests. This is especially useful for expensive tests, tests associated with complications, or when tests are unavailable. Additional diagnostic testing should preferably only be performed when it is expected to change treatment, and the potential clinical benefit outweighs the risk of the tests). Note that in this study, we do not take into account the (un)certainty around imputed values when assessing treatment decision support. Additionally, due to the limited data at our disposal, a full evaluation of the impact on predicted risk was not possible.

**Table 4.** 2 × 2 tables of guideline adherence to treatment threshold given the point estimate of each method

| | True value | | |
|---|---|---|---|
| **Mean imputation** | **Treatment advised (≥140 mmHg)** | **Treatment not advised (<140 mmHg)** | **Totals** |
| Point estimate | | | |
|    Treatment advised (>140 mmHg) | 1,134 | 837 | 1,971 |
|    Treatment not advised (<140 mmHg) | 0 | 0 | 0 |
| Totals | 1,134 | 837 | 1,971 |

| | True value | | |
|---|---|---|---|
| **Joint modeling imputation** | **Treatment advised (≥140 mmHg)** | **Treatment not advised (<140 mmHg)** | **Totals** |
| Point estimate | | | |
|    Treatment advised (>140 mmHg) | 946 | 317 | 1,263 |
|    Treatment not advised (<140 mmHg) | 188 | 520 | 708 |
| Totals | 1,134 | 837 | 1,971 |

| | True value | | |
|---|---|---|---|
| **Conditional modeling imputation** | **Treatment advised (≥140 mmHg)** | **Treatment not advised (<140 mmHg)** | **Totals** |
| Point estimate | | | |
|    Treatment advised (>140 mmHg) | 960 | 315 | 1,275 |
|    Treatment not advised (<140 mmHg) | 174 | 522 | 696 |
| Totals | 1,134 | 837 | 1,971 |

Sensitivity 100%, specificity 0%, Positive Predictive Value 58%, Negative Predictive Value (cannot be calculated) %.
Sensitivity 83%, specificity 62%, Positive Predictive Value 75%, Negative Predictive Value 73%.
Sensitivity 85%, specificity 62%, Positive Predictive Value 75%, Negative Predictive Value 75%.

Ideally, uncertainty in imputed values should be propagated to (additional) uncertainty in predicted risk and evaluated with presented confidence intervals. The predicted risk in this paper primarily serves as a way to illustrate how imputations could influence the predicted risk.

In cardiovascular risk management, the decision to start treatment of a risk factor is based on (i) the predicted risk for a cardiovascular disease or patient characteristics that are per definition associated with a high risk for cardiovascular disease and (ii) the absolute value of the risk factor itself. We focused on imputation models to recover the missing value and to quantify its uncertainty. We demonstrated that the choice of imputation method might impact risk predictions and decision making. While the magnitude of this effect was not always substantial, it may vary according to the number of missing predictors and their weight in the decision-making process and should, therefore, be evaluated when applying these models in different settings and populations.

Last, traditional (e.g., regression-based) prediction models assume complete input data, which is often not realistic in routine clinical practice. Although we developed models for imputing the missing values, which can subsequently be used to generate predictions, it is also possible to develop prediction models that do not require complete information on the predictors. Well-known examples are the use of decision trees with surrogate or sparsity-aware splits [34–36], the use of submodels [37], or the use of missing indicator variables [38]. More research is warranted to evaluate whether these methods may offer any improvement in model predictions, as well as facilitate their implementation in routine care.

In summary, this study describes three imputation methods to handle missing values in the context of computerized decision support systems in clinical practice. We found that JMI and CMI provide imputations that are closer to the original value (as compared to mean imputation) and able to reflect uncertainty due to missing data. We,

**Table 5.** Clinical interpretation of imputed SBP values and 95% confidence intervals from a patient with a history of CVD

| | True | Scenario A | Scenario B |
|---|---|---|---|
| SBP (95%CI) | 144 | 144 (138–150) | 144 (142–146) |
| Treatment based on point estimate | >140 mmHg, Start treatment | >140 mmHg, Start treatment | >140 mmHg, Start treatment |
| Treatment based on 95% CI | NA | Uncertain | >140 mmHg, Start treatment |

*Abbreviations:* SBP, systolic blood pressure; 95% CI, 95% confidence interval; A, hypothetical situation where imputed value interval contains treatment threshold; B, hypothetical situation where imputed value interval does not contain treatment threshold.

**Table 6.** Differences in predicted 10-yr risk of CVD for both a single missing predictor scenario and a multiple missing predictor scenario

| Single missing predictor: eGFR | 25% IQR | Absolute risk difference to completed data | Median | Absolute risk difference to completed data | 75% IQR | Absolute risk difference to completed data |
|---|---|---|---|---|---|---|
| Predicted risk complete data | 8.382% | – | 13.711% | – | 28.170% | – |
| Predicted risk (mean) | 7.287% | −1.095% | 11.997% | −1.713% | 23.035% | −5.135% |
| Predicted risk (joint) | 8.767% | 0.385% | 14.012% | 0.301% | 27.734% | 0.435% |
| Predicted risk (conditional) | 8.786% | 0.404% | 14.110% | 0.399% | 27.783% | 0.387% |
| Multiple missing predictors: SBP, TC, LDL-c and eGFR | | | | | | |
| Predicted risk complete data | 8.382% | – | 13.711% | – | 28.170% | – |
| Predicted risk (mean) | 7.473% | −0.909% | 11.647% | −2.064% | 22.692% | −5.478% |
| Predicted risk (joint) | 8.809% | 0.427% | 14.085% | 0.375% | 28.410% | 0.240% |
| Predicted risk (conditional) | 8.786% | 0.404% | 14.100% | 0.390% | 28.267% | 0.097% |

therefore, recommend their implementation in situations where information on relevant predictors is often incomplete due to practical constraints.

## Supplementary data

Supplementary data to this article can be found online at https://doi.org/10.1016/j.jclinepi.2021.01.003. References [39,40] are cited in the appendix section.

## References

[1] Piepoli MF, Hoes AW, Agewall S, Albus C, Brotons C, Catapano AL, et al. 2016 European guidelines on cardiovascular disease prevention in clinical practice: the sixth joint task force of the European society of Cardiology and other societies on cardiovascular disease prevention in clinical practice (constituted by representatives of 10 societies and by invited experts) developed with the special contribution of the European association for cardiovascular prevention & rehabilitation (EACPR). Eur Heart J 2016;37:2315–81.

[2] Hoffman MA, Williams MS. Electronic medical records and personalized medicine. Hum Genet 2011;130(1):33–9.

[3] Ginsburg G. Personalized medicine: revolutionizing drug discovery and patient care. Trends Biotechnol 2001;19(12):491–6.

[4] Groenhof TKJ, Groenwold RHH, Grobbee DE, Visseren FLJ, Bots ML, On Behalf of the UCC-SMART Study Group. The effect of computerized decision support systems on cardiovascular risk factors: a systematic review and meta-analysis. BMC Med Inform Decis Mak 2019;19(1):108.

[5] Bezemer T, de Groot MC, Blasse E, ten Berg MJ, Kappen TH, Bredenoord AL, et al. A Human(e) factor in clinical decision support systems. J Med Internet Res 2019;21(3):e11732.

[6] Groenhof TKJ, Bots ML, Brandjes M, Jacobs JJL, Grobbee DE, van Solinge WW, et al. A computerised decision support system for cardiovascular risk management 'live' in the electronic health record environment: development, validation and implementation—the Utrecht Cardiovascular Cohort Initiative. Neth Heart J 2019;27(9): 435–42.

[7] Kawamoto K, Houlihan CA, Balas EA, Lobach DF. Improving clinical practice using clinical decision support systems: a systematic review of trials to identify features critical to success. BMJ 2005;330: 765.

[8] Sitapati A, Kim H, Berkovich B, Marmor R, Singh S, El-Kareh R, et al. Integrated precision medicine: the role of electronic health records in delivering personalized treatment: Integrated precision medicine. Wiley Interdiscip Rev Syst Biol Med 2017;9(3):e1378.

[9] Kotseva K, Wood D, De Bacquer D, De Backer G, Rydén L, Jennings C, et al. Euroaspire IV: a European Society of Cardiology survey on the lifestyle, risk factor and therapeutic management of coronary patients from 24 European countries. Eur J Prev Cardiol 2016;23(6):636–48.

[10] Rumsfeld JS, Joynt KE, Maddox TM. Big data analytics to improve cardiovascular care: promise and challenges. Nat Rev Cardiol 2016; 13(6):350–9.

[11] Perkins NJ, Cole SR, Harel O, Tchetgen Tchetgen EJ, Sun B, Mitchell EM, et al. Principled approaches to missing data in epidemiologic studies. Am J Epidemiol 2018;187:568–75.

[12] Azur MJ, Stuart EA, Frangakis C, Leaf PJ. Multiple imputation by chained equations: what is it and how does it work?: multiple imputation by chained equations. Int J Methods Psychiatr Res 2011;20(1): 40–9.

[13] Van Buuren S. Flexible imputation of missing data. 2nd ed. Boca Raton, FL: CRC Press; 2018.

[14] Hoogland J, van Barreveld M, Debray TPA, Reitsma JB, Verstraelen TE, Dijkgraaf MGW, et al. Handling missing predictor values when validating and applying a prediction model to new patients. Stat Med 2020;17:3591–607.

[15] Berkelmans G. Dealing with missing patient characteristics when using cardiovascular prediction models in clinical practice. Oxford, England: European Heart Journal, Oxford University Press; 2018: 110–33. https://doi.org/10.1093/eurheartj/ehy565.P1533. Accessed January 4, 2019.

[16] Janssen KJM, Vergouwe Y, Donders ART, Harrell FE, Chen Q, Grobbee DE, et al. Dealing with missing predictor values when applying clinical prediction models. 8. Clinical chemistry, American Association for Clinical Chemistry: Washington, DC. Available at https://doi.org/10.1373/clinchem.2008.115345. Accessed 1 February 2019.

[17] Carpenter JR, Kenward MG. Multiple imputation and its application. Chichester: Wiley; 2013:345. : (Statistics in practice).

[18] Donders ART, van der Heijden GJMG, Stijnen T, Moons KGM. Review: a gentle introduction to imputation of missing values. J Clin Epidemiol 2006;59:1087–91.

[19] Gökçay D, Eken A, Baltacı S. Binary classification using neural and clinical features: an application in fibromyalgia with likelihood based decision level fusion. 10. IEEE Journal of Biomedical and Health Informatics: Piscataway, NJ. Available at https://doi.org/10.1109/JBHI.2018.2844300. Accessed 5 December 2020.

[20] Debédat J, Sokolovska N, Coupaye M, Panunzi S, Chakaroun R, Genser L, et al. Long-term relapse of type 2 diabetes after roux-en-Y Gastric Bypass: prediction and clinical relevance. Diabetes Care 2018;41:2086–95.

[21] Chen R, Stewart WF, Sun J, Ng K, Yan X. Recurrent neural networks for early detection of heart failure from longitudinal electronic health record data. Circ Cardiovasc Qual Outcomes 2019;15:e005114.

[22] Nederlands Huisartsen Genootschap. Multidisciplinaire richtlijn cardiovasculair risicomanagement. Houten: Bohn Stafleu van Loghum; 2011.

[23] Dorresteijn JAN, Visseren FLJ, Wassink AMJ, Gondrie MJA, Steyerberg EW, Ridker PM, et al. Development and validation of a prediction rule for recurrent vascular events based on a cohort study of patients with arterial disease: the SMART risk score. Heart 2013; 99:866–72.

[24] Damen JAAG, Hooft L, Schuit E, Debray TPA, Collins GS, Tzoulaki I, et al. Prediction models for cardiovascular disease risk in the general population: systematic review. BMJ 2016i2416.

[25] Visseren FLJ, Dorresteijn JAN, van der Graaf Y. U-prevent u bent 'in control' [Internet] 2018. Available at https://www.u-prevent.nl/nl-NL. Accessed October 7, 2019.

[26] Asselbergs FW, Visseren FL, Bots ML, de Borst GJ, Buijsrogge MP, Dieleman JM, et al. Uniform data collection in routine clinical practice in cardiovascular patients for optimal care, quality control and research: the Utrecht Cardiovascular Cohort. Eur J Prev Cardiol 2017;24(8):840–7.

[27] Kowarik A, Templ M. Imputation with the R Package VIM. J Stat Softw 2016;74. Available at http://www.jstatsoft.org/v74/i07/. Accessed September 26, 2019.

[28] Burton A, Altman DG, Royston P, Holder RL. The design of simulation studies in medical statistics. Stat Med 2006;25:4279–92.

[29] Hughes RA, White IR, Seaman SR, Carpenter JR, Tilling K, Sterne JA. Joint modelling rationale for chained equations. BMC Med Res Methodol 2014;14:28.

[30] Demirtas H, Hedeker D. An imputation strategy for incomplete longitudinal ordinal data. Stat Med 2008;27:4086–93.

[31] Murray JS. Multiple imputation: a review of practical and theoretical findings. Arxiv180104058 Stat 2018. Available at http://arxiv.org/abs/1801.04058. Accessed September 26, 2019.

[32] Vergouwe Y, Nieboer D, Oostenbrink R, Debray TPA, Murray GD, Kattan MW, et al. A closed testing procedure to select an appropriate method for updating prediction models: method selection to update a prediction model. Stat Med 2017;36:4529–39.

[33] Varadhan R. condMVNorm: conditional multivariate normal distribution. 2015; R package version 2015.2-1. Available at http://CRAN.R-project.org/package=condMVNorm. Accessed April 3, 2019.

[34] Hapfelmeier A. Analysis of missing data with random forests [Internet]. 2012. Available at https://edoc.ub.uni-muenchen.de/15058/1/Hapfelmeier_Alexander.pdf. Accessed September 4, 2019.

[35] Feelders A. Handling missing data in trees: surrogate splits or statistical imputation?. In: Żytkow JM, Rauch J, editors. Principles of data mining and knowledge discovery [Internet]. Berlin, Heidelberg: Springer Berlin Heidelberg; 1999:329–34. Available at http://link.springer.com/10.1007/978-3-540-48247-5_38. Accessed October 2, 2019.

[36] Burgette LF, Reiter JP. Multiple imputation for missing data via sequential regression trees. Am J Epidemiol 2010;172:1070–6.

[37] Hoogland J, van Barreveld M, Debray T, Reitsma J, Verstraelen T, Dijkgraaf M, et al. Handling missing predictor values when validating and applying a prediction model to new patients. Stat Med 2019;39:3591–607. [Under review].

[38] Groenwold RHH, White IR, Donders ART, Carpenter JR, Altman DG, Moons KGM. Missing covariate data in clinical research: when and when not to use the missing-indicator method for analysis. Can Med Assoc J 2012;184(11):1265–9.

[39] Chen M-H. Monte carlo methods in bayesian computation. Berlin, Germany: Springer; 2013.

[40] Genz A, Bretz F, Miwa T, Mi X, Leisch F, Scheipl F, et al. Mvtnorm: multivariate normal and t distributions. J Stat Softw 2018. Available at http://CRAN.R-project.org/package=mvtnorm. Accessed April 3, 2019.