

BEST : A decision tree algorithm that handles missing values

Cédric Beaulac ¹ Jeffrey S. Rosenthal ²

July 31, 2019

Abstract

The main contribution of this paper is the development of a new decision tree algorithm. The proposed approach allows users to guide the algorithm through the data partitioning process. We believe this feature has many applications but in this paper we demonstrate how to utilize this algorithm to analyse data sets containing missing values. We tested our algorithm against simulated data sets with various missing data structures and a real data set. The results demonstrate that this new classification procedure efficiently handles missing values and produces results that are slightly more accurate and more interpretable than most common procedures without any imputations or pre-processing.

Keywords : Classification and Regression Tree, Missing Data, Applied Machine Learning, Interpretable Models, Variable Importance Analysis

¹University of Toronto
Department of Statistical Sciences
Sidney Smith Hall, 100 St George St, Toronto, ON M5S 3G3
E-mail : beaulac.cedric@gmail.com
ORCID : 0000-0002-6050-5313

²University of Toronto
Department of Statistical Sciences
Sidney Smith Hall, 100 St George St, Toronto, ON M5S 3G3
E-mail : jeff@math.toronto.edu

1 Introduction

Machine learning algorithms are used in many exciting real data applications, but may have problems handling predictors with missing values. Imputation techniques were designed to handle data with missing value under the missing completely at random (MCAR) assumption. Since this is a restrictive assumption we propose a solution to missing values that uses the tree structure of Classification and Regression Trees (CART) to deal in an intuitive manner with observations that are missing in patterns which are not completely at random.

Our proposed new tree construction procedure was inspired by a data set where the missing pattern of one subset of predictors could be perfectly explained by another subset (see Section 5.1). A typical decision tree is an algorithm that partitions the predictor space based upon a predictor value, splitting it into two subspaces and repeats this process recursively.

Our proposed algorithm is different as it allows the researcher to impose a structure on the variables available for the partitioning process. By doing so, we construct Branch-Exclusive Splits Trees (BEST). When a predictor X_j contains missing values, we can use other predictors to identify the region where the predictor X_j contains no missing value. Therefore we can use the proposed algorithm to consider splitting on a predictor only when it contains no missing value based on previous partitioning. BEST can be easily adapted to any splitting rule and any forest forming procedure [4, 5, 10]. BEST also has other applications; it can be used by researchers that would like to utilize some knowledge they have on the data generating distribution in order to guide the algorithm in selecting a more accurate and more interpretable classifier.

In this article we will briefly discuss the classification problem and its notation and we will explain how classification trees solve that problem. We will then do a quick review of the literature on missing value techniques currently being used. Afterwards, we will introduce the proposed algorithm and some motivating examples before explaining in

detail how the algorithm functions. Finally, some tests will be performed on simulated data sets and on the real data that inspired this new algorithm.

2 The classification problem

In a typical supervised statistical learning problem we are interested in understanding the relationship between a response variable \mathbf{Y} and an associated m -dimensional predictor vector $\mathbf{X} = (X_1, \dots, X_m)$. When the response variable is categorical and takes k different possible values, this problem is defined as a k -class classification problem. In that set up, an interesting challenge is to use a data set $S = \{(y_i, x_{i,1}, \dots, x_{i,m}); i = 1, \dots, n\}$ in order to construct a classifier h . Most of the time, it is assumed that the observations within our data set were drawn independently from the same unknown and true distribution \mathcal{D} , i.e. $\mathbf{X} \times \mathbf{Y} \sim \mathcal{D}$. A classifier is built to emit a class prediction for any new data point X that belongs in the predictor space $\mathcal{X} = \mathcal{X}_1 \times \dots \times \mathcal{X}_m$. Therefore a classifier divides the predictor space \mathcal{X} into k disjoint regions R_1, \dots, R_k , one per class, such that $\cup_{q=1}^k R_q = \mathcal{X}$, i.e. $h(x) = \sum_{q=1}^k q \mathbf{1}\{x \in R_q\}$.

2.1 Classification and Regression Trees

A classification tree [3] is an algorithm that forms regions in the predictor space by recursively dividing it, more precisely, this procedure performs recursive binary partitioning. Beginning with the entire predictor space, the algorithm selects the variable to split upon and the location of the split that minimize some impurity measure. Then the resulting two regions are each split into two more regions until some stopping rule is applied. The classifier will label each region with one of the k possible classes.

The traditional labelling process goes as follows; let $p_{rq} = \frac{1}{n_r} \sum_{x_i \in R_r} \mathbf{1}\{y_i = q\}$, the proportion of the class q in the region r where n_r is the number of observations contained in region r . Then, the label of the region r is the majority class in that region, i.e. if $x \in R_r$, $h_S(x) = \operatorname{argmax}_q(p_{rq})$. For regression trees, the output mean within a leaf region

is used as prediction for observations that belong in that region. The impurity measure function for region r is defined as Q_r and can take many forms such as the *Gini index*, the *deviance* or the *misclassification error*. For regression trees, the *mean squared error* is one possible region impurity measure.

When splitting a region R_p into two new regions R_r and R_t the algorithm will compute the total impurity of the new regions ; $n_r Q_r + n_t Q_t$ and will pick the split variable X_j and split location s that minimizes that total impurity. If the predictor X_j is continuous, the possible splits are of the form $X_j \leq s$ and $X_j > s$ which usually results in $n_p - 1$ possible splits. For a categorical predictor having c possible values, we usually consider all of the $2^{c-1} - 1$ possible partitions.

The partitioning continues until a stopping rule is applied. In some cases, the algorithm stops whenever every terminal node of the tree contains less than β observations, in other cases it stops when all observations within a region belong to the same class. To prevent overfitting, a deep tree is built and then the tree can be pruned. Tree-pruning is a cost-complexity procedure that relies on considering that each leaf, region, is associated with a cost α . The procedure begins by collapsing leaves that produce the smallest increase in total impurity and this technique will collapse leaves as long as the increase in impurity is less than the cost α of the additional leaf. The α parameter can be determined by cross-validation or with the use of a validation set.

3 Missing values

Let us now introduce the definition of missing data we are using in this article. As described in the previous section, a standard assumption in data analysis is that all observations are distributed according to the true data generation distribution \mathcal{D} . We could think of the missingness itself as a random variable \mathbf{M} also of dimension m that is distributed according to some missingness generating distribution which is a part of \mathcal{D} , i.e. $\mathbf{X} \times \mathbf{M} \times \mathbf{Y} \sim \mathcal{D}$. Formally, if \mathbf{M} represents the missingness of the vector of predictors \mathbf{X} it means that $M_j = 1$ if X_j is observed and $M_j = 0$ if X_j is missing.

Three different relationships between \mathbf{M} and \mathbf{X} were defined by Rubin [19] and by Little and Rubin [15]. Seaman [21] later untangled the many definition inconsistencies of these relationships. In this article, we will rely on simple definitions for an easy understanding of the structure we will consider. First, missing completely at random (MCAR) is the simplest structure we consider: $\mathbf{M} \perp \mathbf{X}$. Here, we consider the data is MCAR if the set of missing patterns M is independent of the set of predictors.

Missing at random (MAR) is much more complicated; it essentially means that the missingness \mathbf{M} is independent of missing observations but can still depend on observed predictors. More rigorously, we define $X^o = \{x_{ij} \in S | x_{ij} \text{ is observed}\}$ as the set of all observed predictors value, and $X^{na} = \{x_{ij} \in S | x_{ij} \text{ is missing}\}$ as the set of missing predictors value. We say that data is MAR if the distribution of the missingness is conditionally independent of missing values given observed values : $\mathbf{M} \perp X^{na} | X^o$. As pointed by Seaman [21], MAR has not always been used consistently and the definition above is the one we settled on for this project. Note that MCAR implies MAR.

Finally, if the missing data is neither MCAR nor MAR, we say that the data is missing not at random (MNAR). We will see that the relationship between \mathbf{M} and \mathbf{X} has a considerable effect on the efficiency of the missing values techniques that exist. In the sections to come we will test our algorithms under these various missing data structures.

4 Related work

Decision trees are well-established and a wide variety of solutions have already been proposed to handle missing values. In this section, we will position our contribution within the current literature. We will briefly introduce the current missing value techniques that are paired with decision trees and we will establish the main differences between these techniques and our proposed algorithm. To do so, we will use recent surveys [20, 27, 6, 9] that defined and compared these techniques using various simulated and real data sets.

Predictive value imputation (PVI) methods are popular approaches to deal with miss-

ing value. They estimate and impute the missing values within both the training and the test set. The simplest imputation consists of replacing the missing values with the mean for numerical predictors or the mode for categorical predictors. More advanced prediction models have also been proposed, such as linear model, k-nearest neighbours or expectation-maximization (EM).

Since these methods use known predictors to impute values for the missing ones, if the predictors are uncorrelated these approaches will have no predictive power. This will lead to poor imputation and it is a major drawback noted by Saar-Tsechansky and Provost [20] and Gavankar [9]. Nonetheless, Twala [27] demonstrated using simulated data sets the great performances of expectation-maximization multiple imputations (EMMI) [14]. This imputation technique produces multiple different imputations based on expectation-maximization and then aggregates the results.

Our proposed algorithm differs from imputation methods as it only uses known information to build the classifier instead of using potentially unreliable prediction to replace missing values.

The surrogate variable (SV) approach [3] is a special case of predictive value imputation. As explained by Hastie et al. [11], during the training process, when considering a predictor for a split, only the observations for which that predictor is not missing are used. After the primary predictor and split point have been selected, a list of surrogate predictors and split points is constructed. The first surrogate split is the predictor and split point pair that best mimic the split of the training data achieved by the primary split. Then the second surrogate split is determined among the leftovers predictors and so on. When splitting the training set during the tree-building procedure or when sending an observation down the tree during prediction, the surrogate splits are used in order if the primary splitting predictor value is missing.

Many articles [7, 20, 27, 6] showed that the results are not satisfactory in many cases and Kim and Loh [12] noted the variable selection biased caused by this approach. Our

proposed approach is much more computationally efficient and utilized the missing pattern as a predictor instead of ignoring it.

Reduced-feature models are suggested by Saar [20] when missing values appear only in the prediction process. This technique relies on using only known predictors of the new observation we are trying to classify. A tree is built using only the known predictors of the new observation. If multiple observations contain different missing patterns then multiple trees are built to classify the various observations. It shares a great deal of similarities with lazy decision trees [8] as both models tailor a classifier to a specific observation and uses only known predictors to do so.

A major drawback of this technique is that it only manages missing values during prediction while our proposed technique can handle missing values for both training and prediction. BEST also differs from reduced-feature models as it not only uses the known values but also utilizes the fact that we know some predictors are missing instead of discarding this information.

The popular C4.5 implementation [16] has its own way to manage missing data, defined as a distribution-based imputation (DBI). When selecting the predictor to split upon, only the observations with known values are considered. After choosing the best predictor to split upon, observations with known values are split as usual. Observations with unknown values are distributed among the two child nodes proportionately to the split on observed values. Similarly, for prediction, a new test observation with missing value is split into branches according to the portions of training examples falling into those branches. The prediction is then based upon a weighted prediction among possible leaves.

This technique is computationally slow and offers poor performances in terms of prediction accuracy according to some of the surveys we mentioned [20, 28]. Our technique should be faster, more accurate and more interpretable.

As described by Ding and Simonoff [6], the Separate Class (SC) method replaces the missing value with a new value or a new class for all missing observations. For categorical

predictors we can simply define *missing value* as a category on its own and for continuous predictors any value out of the interval of observed value can be used.

This technique has the best performances according to Ding and Simonoff [6] when there are missing values in both the training and the test set and when observations are missing not at random (MNAR). Twala et al. [28] also came up with similar results with a generalization of the separate class method coined Missing Incorporated in Attribute (MIA). These techniques are by far the closest in spirit to BEST. As BEST, MIA and SC allow for similar data partitioning we do not expect BEST to offer a drastic improvement in terms of accuracy. On the other hand, BEST identifies the missing pattern using other predictors rather than including the missingness information within the predictor containing missing values. Doing so, our approach should offer highly interpretable results and a more accurate variable importance analysis.

Finally, there exist many other articles discussing decision trees in the context of missing values. Some introduce ways to use decision trees and random forests as missing values imputation techniques such as Rahman and Islam [18] and others to identify the missingness structure such as Tierney et al. [26]. Our approach is different since we are not using decision trees to pre-process data with missing value or to identify the missingness structure but to perform a decision tree analysis of a data set that contains missing values.

5 Branch-Exclusive Splits Trees (BEST)

We now introduce the proposed algorithm, BEST. The purpose of BEST is to utilize the tree structure itself in order to manage some missing data or some special structure among predictors.

As we explained in section 2.1, a classification tree aims at partitioning the predictor space and labelling the resulting regions. CART does so by looking through all the possible splits and selecting the one that minimizes some pre-specified error measure.

When using BEST some predictors are available to split upon only within some regions of the predictor space. These regions are defined according to other predictors values, for example, predictor X_l could be only available for the partitioning process in the region defined by $X_j < t$. We say that X_j is a *gating variable* for X_l or that X_l is *gated* by X_j . The variable X_l will not be available for the partitioning process until the gating variable allows it. The table below illustrates when BEST can partition the data using X_j based on a previous partitioning where BEST selects X_j as the splitting variable and s as the splitting value :

	$s < t$		$s > t$	
	Region $X_j < s$	Region $X_j \geq s$	Region $X_j < s$	Region $X_j \geq s$
X_l is available	Yes	No	No	No

Table 1: Availability of X_l if X_j is previously selected as splitting variable and s as splitting value

Doing so, predictors with missing values can be handled easily as BEST will partition the data according to that predictor only in regions where it does not contain missing value. Similarly, some insight on the data structure can be used to force some variable to be partitioned upon before others which is another application of BEST not described in this article. The result is a tree-structured classification model where some split variables are branch-exclusives. Even though we do not further mention it, the construction described below could be used for regression trees as well.

5.1 Motivating Example

Let us now explain which data structure BEST is suitable for by introducing the motivating data set. It contains information regarding the academic performances of students. The data set was provided to us by the University of Toronto and was first introduced and analysed by Bailey et al. [1]. It was later analysed by Beaulac and Rosenthal [2] where the goal was to predict whether or not a student would complete its program. The predictors represent the number of credits and grades obtained in all the departments

during the first two semesters. Understanding the importance of these predictors was also a question raised by the authors. Here is what the data would look like with a reduced number of departments :

Student ID	Credits Math	Grade Math	Credits Econ	Grade Econ	Credits Hist	Grade Hist
101	2	72	3	88	0	NA
208	0	NA	0	NA	5	78

Table 2: An example of the motivating data set for 2 students and 3 departments

A student has no grade in many departments as they can only register to a limited number of courses within a year. In this situation, many grade variables are missing for every student. BEST handles that problem by considering the averaged grade obtained in a department only for students who took courses in that department. For example, BEST will force the classification tree algorithm to split upon the *Number of credits* predictor to begin. Then, suppose *Number of credits in Statistics* is selected and 2 is the split point for the partitioning, BEST will then allow splits on the *Grade in Statistics* predictor for the group of students in the region defined by *Number of credits in Statistics* > 2 . Therefore, the *Number of credits* variables are used to define the region where the respective *Grade* variables are available for the partitioning process. The many *Grade* variables are gated by their respective *Number of credits* variables. Figure 1 illustrates how a BEST decision tree uses *Number of credits* as a gating variable for *Grade*.

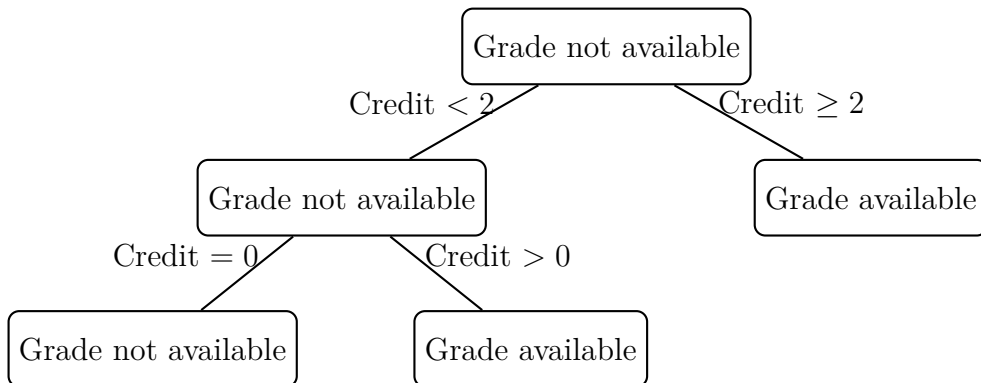


Figure 1: An example of BEST decision tree partitioning upon the Credit variable and the availability of the Grade variable within each of the produced regions

Other real data sets with similar problems are surveys. Many surveys have questions

that are only relevant based on previous answers. Suppose question #1 is yes/no and is followed by : *If you answered no, please go to question #3*. This is quite typical and in that situation, BEST can use the yes/no question as a gating variable for the question #2.

If a data set contains missing values on predictor X_j but no predictor can help define the region with no missing value, we can add a new predictor X_{m+1} to the model as our gating variable. This new predictor is a dummy variable such that $X_{i,m+1} = 0$ if $X_{i,j}$ is missing and 1 if not. Doing so, we effectively add M_j as defined in section 3, as a predictor in the model and thus will be defined as follows in the rest of the text. Then, BEST will only consider splitting on X_j in the subspace defined by $M_j = 1$ as illustrated in figure 2. Multiple dummy variables are added to the model if multiple predictors contain missing values. Doing so also allows us to analyse the individual importance of the missing patterns M .

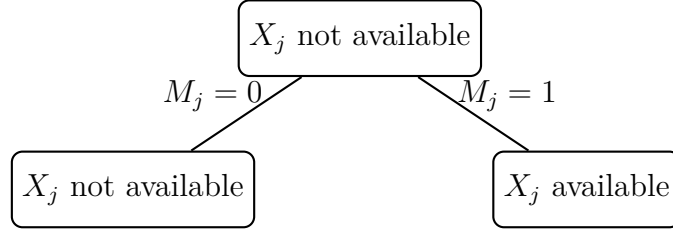


Figure 2: An example of BEST decision tree partitioning upon the dummy variable M_j that defines the availability of X_j

5.2 Theoretical intuition

Formally, the loss of a classifier h is defined as :

$$L_{\mathcal{D}}(h) = \mathbf{P}_{\mathcal{D}}[h(x_i) \neq y_i], \quad (1)$$

which is the probability under the true data generating distribution \mathcal{D} that the classifier h misclassifies an observation x_i . Since the data generating distribution \mathcal{D} is unknown, then the empirical loss computed with the data set S is typically used as an estimator of

the true loss :

$$L_S(h) = \frac{|\{i \in [n] : h(x_i) \neq y_i\}|}{n}, \quad (2)$$

which is the proportion of misclassified observations in the training set S . Usually a set of classifiers \mathcal{H} is selected in advance and most learning algorithms are trying to identify the classifier $h \in \mathcal{H}$ that minimizes the empirical loss $L_S(h)$. The set \mathcal{H} was named the hypothesis class by Shai & Shai [22]. The true loss can be decomposed in a manner to observe a bias-complexity tradeoff. Suppose $h_S = \operatorname{argmin}_{h \in \mathcal{H}} L_S(h)$, then :

$$\begin{aligned} L_D(h_S) &= \min_{h \in \mathcal{H}} L_D(h) + (L_D(h_S) - \min_{h \in \mathcal{H}} L_D(h)). \\ &= e_{\text{app}}(\mathcal{H}) + e_{\text{est}}(h_S). \end{aligned} \quad (3)$$

The approximation error, $\min_{h \in \mathcal{H}} L_D(h) = e_{\text{app}}$, is the minimum achievable loss within the hypothesis class. The second term, $(L_D(h_S) - \min_{h \in \mathcal{H}} L_D(h)) = e_{\text{est}}$, is the estimation error and is caused by the use of the empirical loss instead of the true loss when selecting the best classifier h . Since the goal is to minimize the total loss a natural tradeoff emerges from equation 3. A vast, large and complex hypothesis class \mathcal{H} leads to a wider choice of functions and therefore reduces e_{app} , but the classifier is more prone to overfitting, which increases e_{est} . Inversely, a small hypothesis class \mathcal{H} reduces e_{est} but increases e_{app} .

Our proposed algorithm aims at obtaining a better classifier by restricting the hypothesis class to a smaller one without increasing the approximation error. Suppose \mathcal{H}_T is defined as the set of all tree-structured classifiers. Then, BEST is a new algorithm that aims to find the best classifier in a new hypothesis class \mathcal{H}_B that contains all the tree-structured classifiers that respect a set of conditions regarding the order that variables can be partitioned upon. Therefore, we have $\mathcal{H}_B \subset \mathcal{H}_T$. Since the complexity of \mathcal{H}_B is smaller than the complexity of \mathcal{H}_T the estimation error of BEST will be smaller.

Next, let us take a look at the approximation error : $\min_{h \in \mathcal{H}_B} L_{\mathcal{D}}(h)$. When using BEST, we make multiple assumptions on how the partitioning should be processed. For example, we assume it is better to partition the data using the missing indicator M_j before partitioning the data using X_j . Doing so, we assume that the best tree-structured classifier among all classification tree \mathcal{H}_T is contained within the set of tree-structured classifiers that respect the partition ordering that defines \mathcal{H}_B . In other words, we assume $\operatorname{argmin}_{h \in \mathcal{H}_T} L_{\mathcal{D}}(h) \in \mathcal{H}_B$. Suppose S is a data set, $h_S(\mathcal{H}_T)$ is the classifier that minimizes the empirical loss on \mathcal{H}_T and $h_S(\mathcal{H}_B)$ is the classifier that minimizes the empirical loss on \mathcal{H}_B , then :

$$\begin{aligned}
L_D(h_S(\mathcal{H}_T)) &= \min_{h \in \mathcal{H}_T} L_{\mathcal{D}}(h) + e_{\text{est}}(h_S(\mathcal{H}_T)). \\
&= \min_{h \in \mathcal{H}_B} L_{\mathcal{D}}(h) + e_{\text{est}}(h_S(\mathcal{H}_T)). \\
&\geq \min_{h \in \mathcal{H}_B} L_{\mathcal{D}}(h) + e_{\text{est}}(h_S(\mathcal{H}_B)). \\
&= L_D(h_S(\mathcal{H}_B)),
\end{aligned} \tag{4}$$

which implies that the under the assumption we have made we would not only naturally manage missing values but also reduce the loss. If our assumption $\operatorname{argmin}_{h \in \mathcal{H}_T} L_D(h) \in \mathcal{H}_B$ is false, we might increase the loss, and the assumption itself is impossible to verify. Therefore, the behaviour of the algorithm under multiple scenarios will be tested in section 6 with simulated data.

5.3 Algorithm

Let us now explain how the algorithm functions. BEST takes as input the full data set S , the tuning parameter β and a list containing the predictor availability structure V . First S is set as the root node, the first set of observations to go through the following steps. The algorithm verifies a set of conditions before proceeding with the partitioning process. The first condition (C1) is that the region contains more than β observations, this is the main stopping rule. Then, the next condition (C2) is that the observations in the region have different labels; this condition makes sure that the algorithm has a

reason to partition the data. Finally, the last condition (C3) is that at least one of the available predictors takes different values among the observations in the region; this is to guarantee that the algorithm can actually partition the data.

If at least one condition is false, then the region is defined as a leaf, a label is assigned to that leaf for prediction purposes and the partitioning process is stopped. Usually the class that represents the majority in a leaf is selected as the label for that region, but one could define different label assignment rules.

If all conditions (C1, C2 and C3) are respected then the partitioning process begins. The algorithm will go through all available predictors. For a predictor j , the algorithm will go through all possible partitions s of the region with respect to the predictor j and will compute the total impurity of the resulting two regions $n_{r_1}Q_{r_1} + n_{r_2}Q_{r_2}$. Any region impurity measure Q can be used. BEST then selects the predictor j and the split s that minimize the total impurity and creates two children regions by splitting the data according to s .

The last step is to update the list of available predictors for the children regions. There exist multiple possible structures that can contain this information but within the R-language [17] we have settled on a list V where $V[0]$ represents the set of predictors available for the partitioning process in the root node and $V[i]$ for $i \in 1, \dots, m$ represent the update for the predictors available on the children node after the data has been partitioned using the i th predictor. More specifically $V[0]$ is a vector of size m where $V[0][j] = 1$ if the j th predictor is available to be split upon in the root node and $V[0][j] = 0$ otherwise. Suppose now that the j th predictor is selected as the splitting predictor. If j is a continuous predictor, $V[j][0]$ contains a threshold value, if the splitting point s is greater (or less) than the threshold, then the child node containing the observations such that $X_j > s$ ($X_j \leq s$) has their available predictors updated according to the vector contained in $V[j][1]$ for example. If j is a categorical predictor, then $V[j]$ contains a matrix where the i th line represents the update needed on the node containing the observations such that $X_j = i$ after the partitioning. This formatting was easy to implement in R [17] since

lists are very flexible, but other structure could have been used.

Here is a pseudo-code of the proposed algorithm :

Algorithm : BEST(S, β, V)
<ol style="list-style-type: none"> 1. Starting with the entire data set S and the set of available predictors $V[0]$. 2. Check conditions (C1, C2 and C3). 3. if (any condition is false) : <ul style="list-style-type: none"> Assign a label to the node and exit. else : <ul style="list-style-type: none"> for (j in all available predictors): <ul style="list-style-type: none"> for (s in all possible splits) : Compute total impurity measure. Select the variable j and the split s with minimum impurity measure. Split the node into two child nodes. Update the available predictors for both children nodes using $V[j]$. Repeat steps 2 & 3 on the two children nodes.

In addition, the resulting tree can be pruned, and constructed with any splitting rule, any stopping rule and any label assignment rule. Since one of the goals of this new algorithm is to produce accurate but also interpretable models we did not discuss forests so far, but the proposed tree construction procedure can be used to build any type of forest as well.

6 Experiments : Simulated data sets

All of our experiments are done using the R-language [17]. In the following set of simulations, we will compare 6 methods; (1) BEST, our proposed approach, (2) the Distribution Based Imputation (DBI) proposed by Quinlan [16] implemented in the *C5.0* package [13], (3) a simple single variable imputation (SVI), either the mode for a categorical predictor or the mean for a numerical one, (4) a refined predictive value imputation (PVI) using known predictors; EM for numerical predictors and multinomial logistic regression for

categorical ones [29], (5) the separate class (SC) approach and finally (6) the surrogate variable technique introduced by Breiman et al. [3] implemented in the *rpart* package [25]. Since the Reduced-Feature Model was the least accurate in every single experiment we have done, we decided not to include it in the following plots to improve readability.

We will generate data sets containing 4 predictors; X_1 and X_2 are binary predictors and X_3 and X_4 are continuous predictors on $(0, 1)$. The response is categorical and can take up to 8 different values. The binary predictors are generated according to a Bernoulli distribution and the continuous predictors are generated according to a Uniform distribution. The response labels are assigned according to the tree in figure 3 but a proportion of the responses labels are randomly assigned.

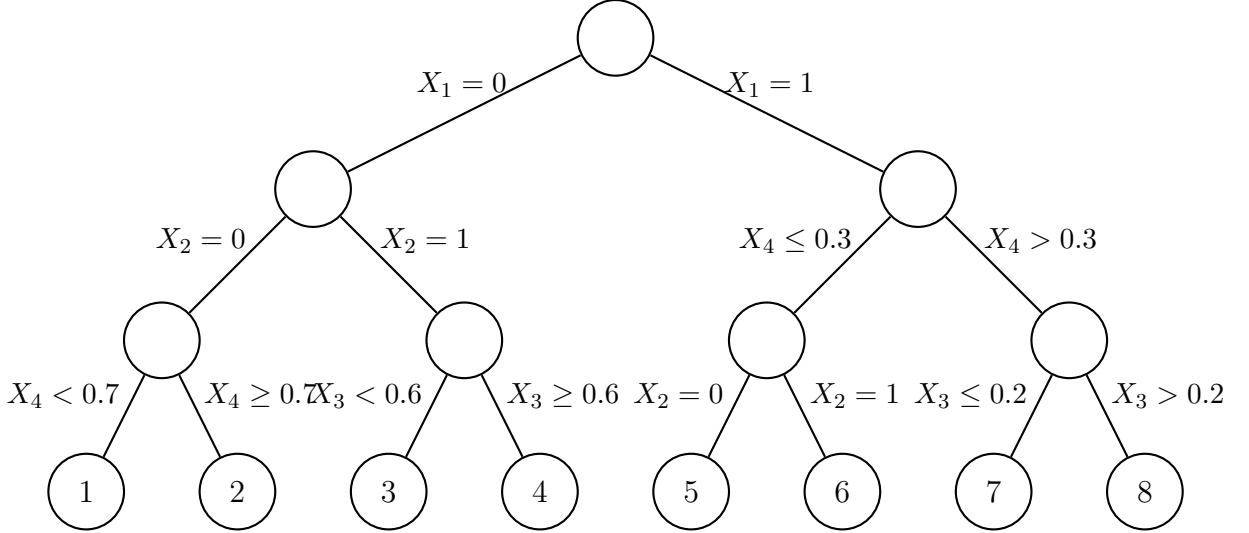


Figure 3: The decision tree used to generate our simulation data sets

Let us describe the experimental procedure. We begin by generating a data set as we described in the previous paragraph. We will apply a missing pattern to the data set; details are included in the respective subsections. Then we will fit a pruned decision tree using each of the six methods mentioned earlier in the section. We will adjust various parameters such as the number of training observations, the variable containing missing values and the type of missing patterns. We will compare the technique performances using the prediction accuracy on the test set where we define the accuracy as the proportion of correctly classified observations.

6.1 MAR : Missingness depends on observed predictors

This first experimental set up is meant to test the missing pattern structure BEST was designed for. In this set up, the missing pattern of a predictor is fully explained by another, fully observed, predictor. The binary predictor X_1 was designated as gating variable for a randomly sampled predictor, either X_2 , X_3 or X_4 . In our first experiment we randomly sample either 0 or 1 and the gated variable is missing if X_1 equals the randomly sampled value. This procedure will be repeated 200 times for this experiment.

We will present our results using Sina Plots [23, 30]. This will allow to better visualize the distribution of the performances of the different techniques for different predictors with missing values.

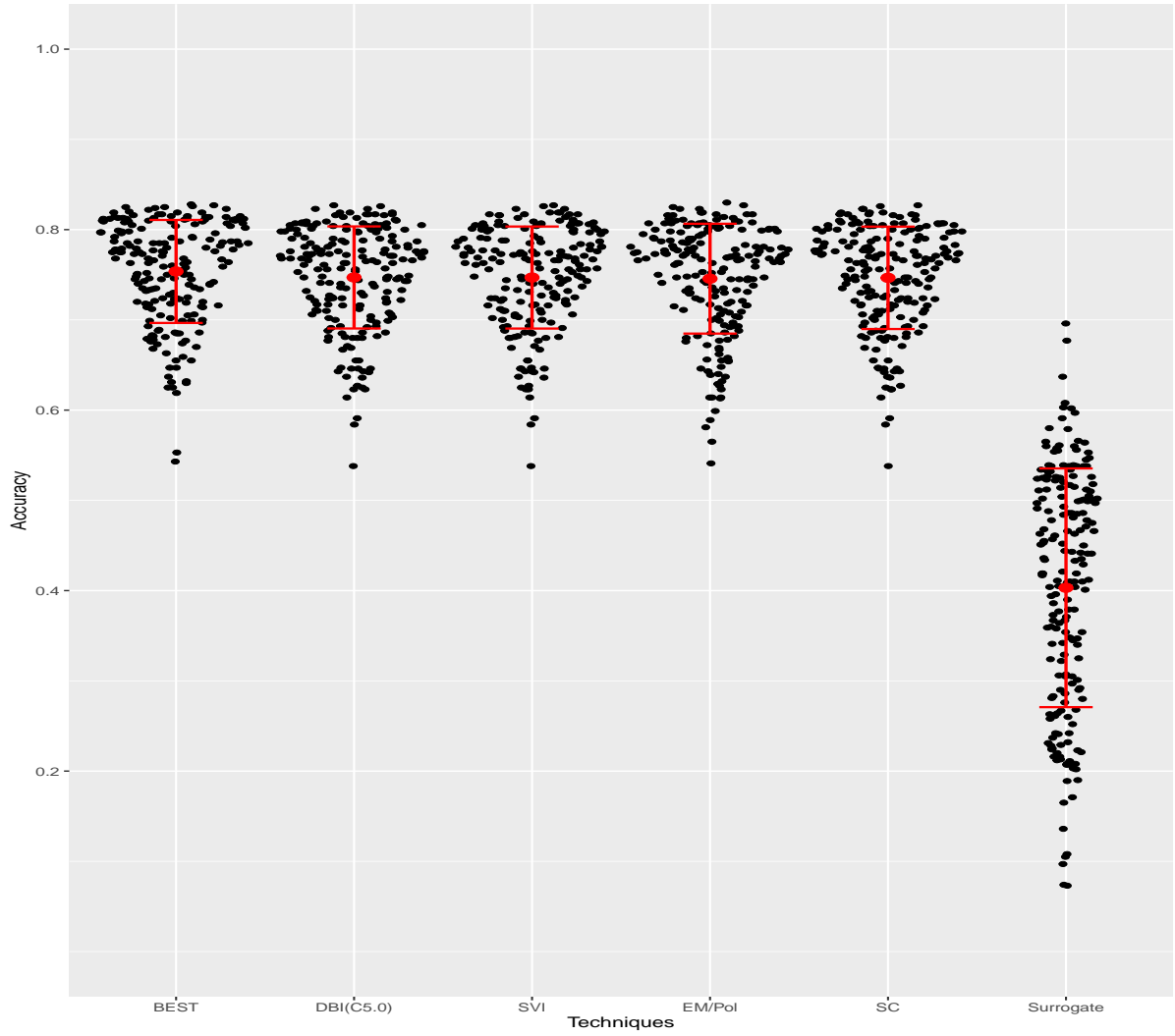


Figure 4: Accuracy by techniques when the missing pattern depends on other predictors with 200 training observations

The figure 4 contains the results of our simulation when we the techniques used 200 training observations to build the classifiers. Figure 5 contains the results of our simulation obtained when we used 1000 training observations.

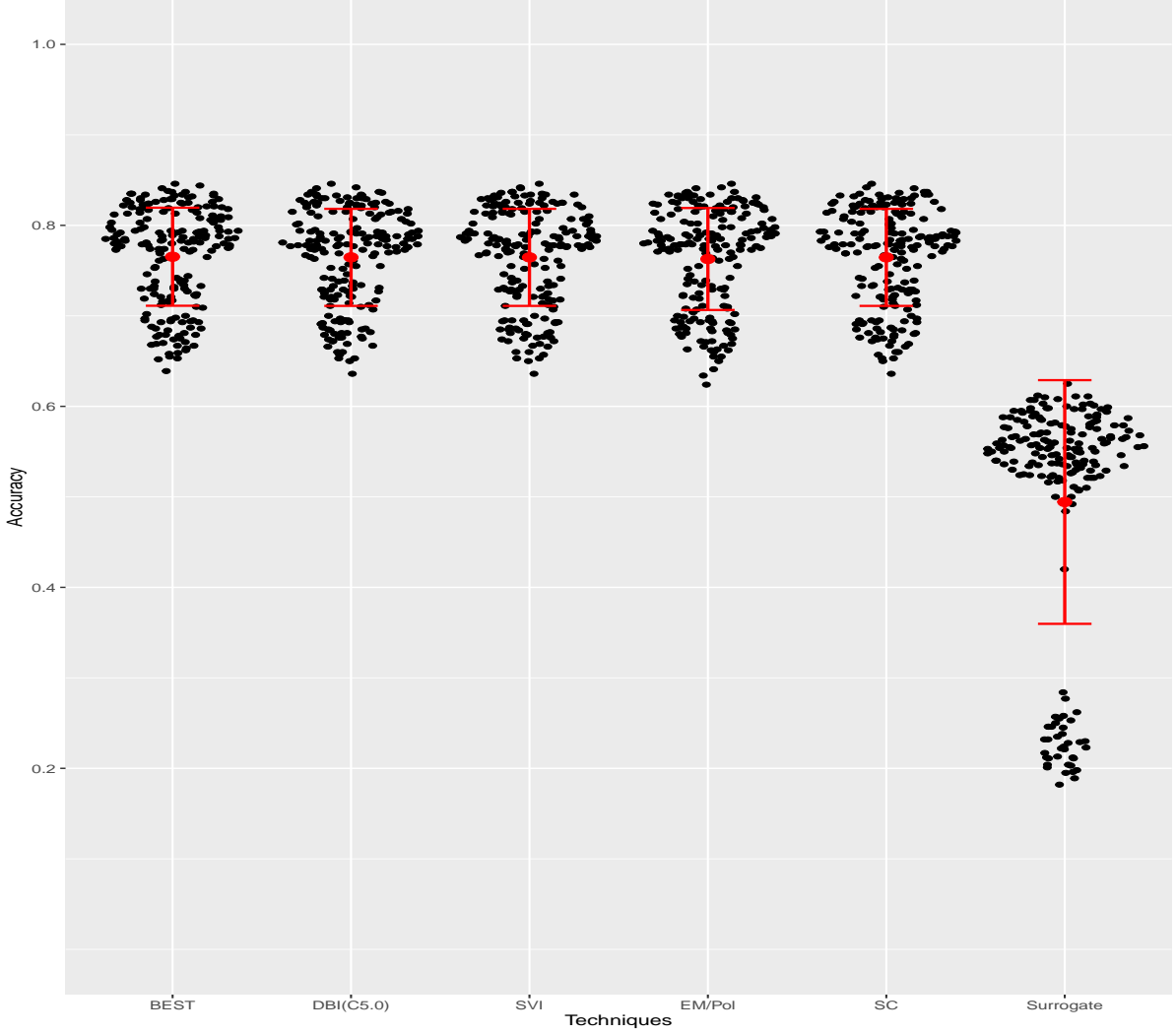


Figure 5: Accuracy by techniques when the missing pattern depends on other predictors with 1000 training observations

Both of the plots convey similar information. When the missing pattern depends on other predictors, the performance of BEST is similar to many competitors but BEST still leads to more interpretable decision trees and does not require any imputation. There is no notable differences between the results obtained with 200 observations and 1000 observations in that precise experiment. Going forward we will keep the number of training observations fixed to 200 but we will vary other parameters.

In the next experiment we will once again designate X_1 as the gating variable and will again sample one of the other predictors as the gated predictor. Once again, we randomly sample either 0 or 1 but this time, we will generate as less extreme missing pattern. The gated variable will be missing with a probability of 50% when X_1 equals the randomly

sampled value.

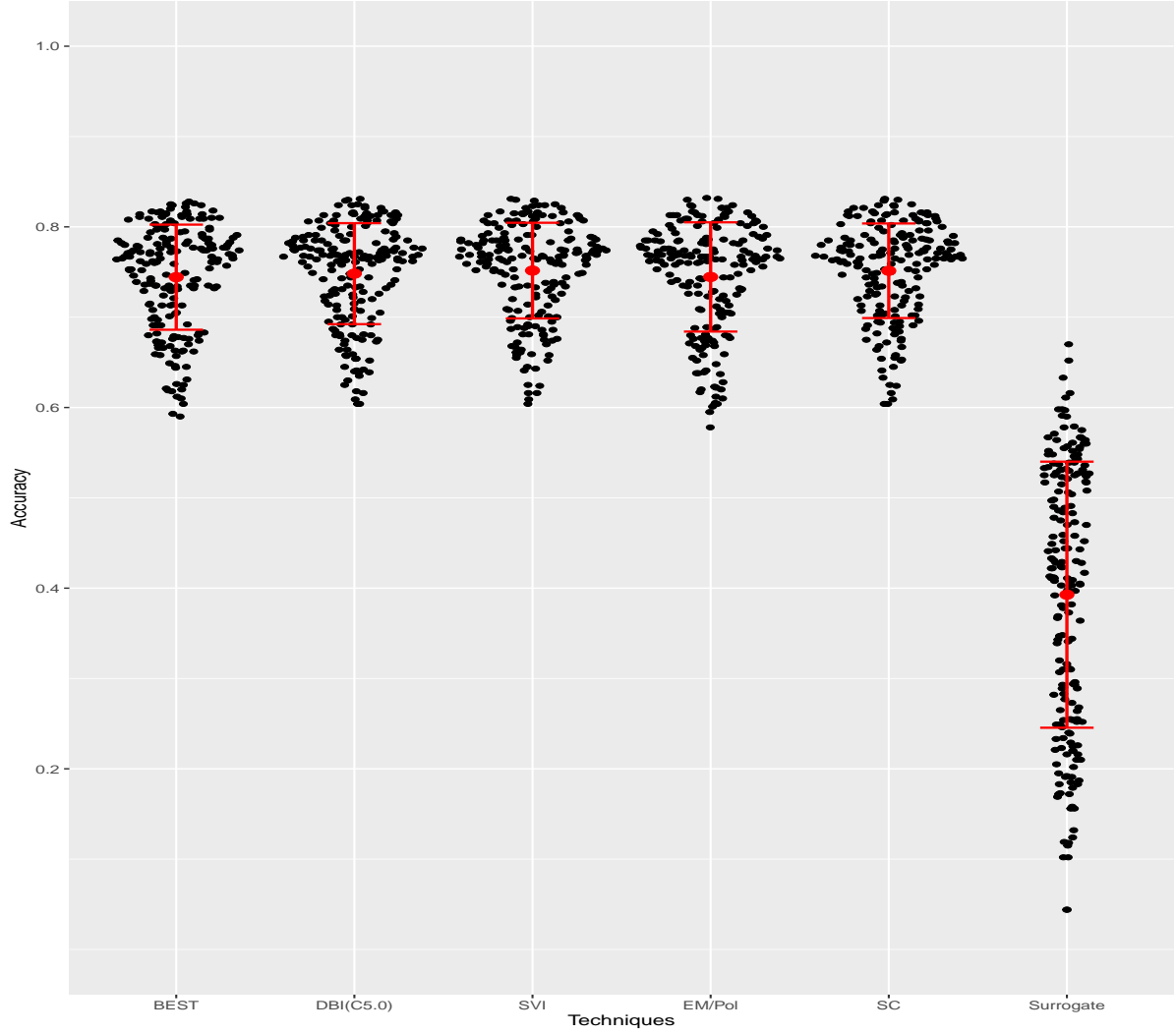


Figure 6: Accuracy by techniques with 50% of the gated variable missing when the condition of the gating variable is respected

We can observe in figure 6 that a less extreme missing pattern does not affect one technique more than others.

Overall, BEST performances when the data is missing at random given other predictors is as good as other tested techniques.

6.2 MAR : Missingness depends on the response

According to Ding et al. [6], the relationship between the missing pattern and the response variable has a great effect on the results obtained from different missing value treatments.

In this simulation, one of the predictors is randomly selected, let us say X_j , every iteration and the censoring process is then applied. The censoring process goes as follows; one of the eight response labels is randomly selected, and X_j is missing for all observations with that selected label. In this experiment we have used a dummy variable as the gating variable for the BEST algorithm. This procedure will be repeated 100 times for this experiment.

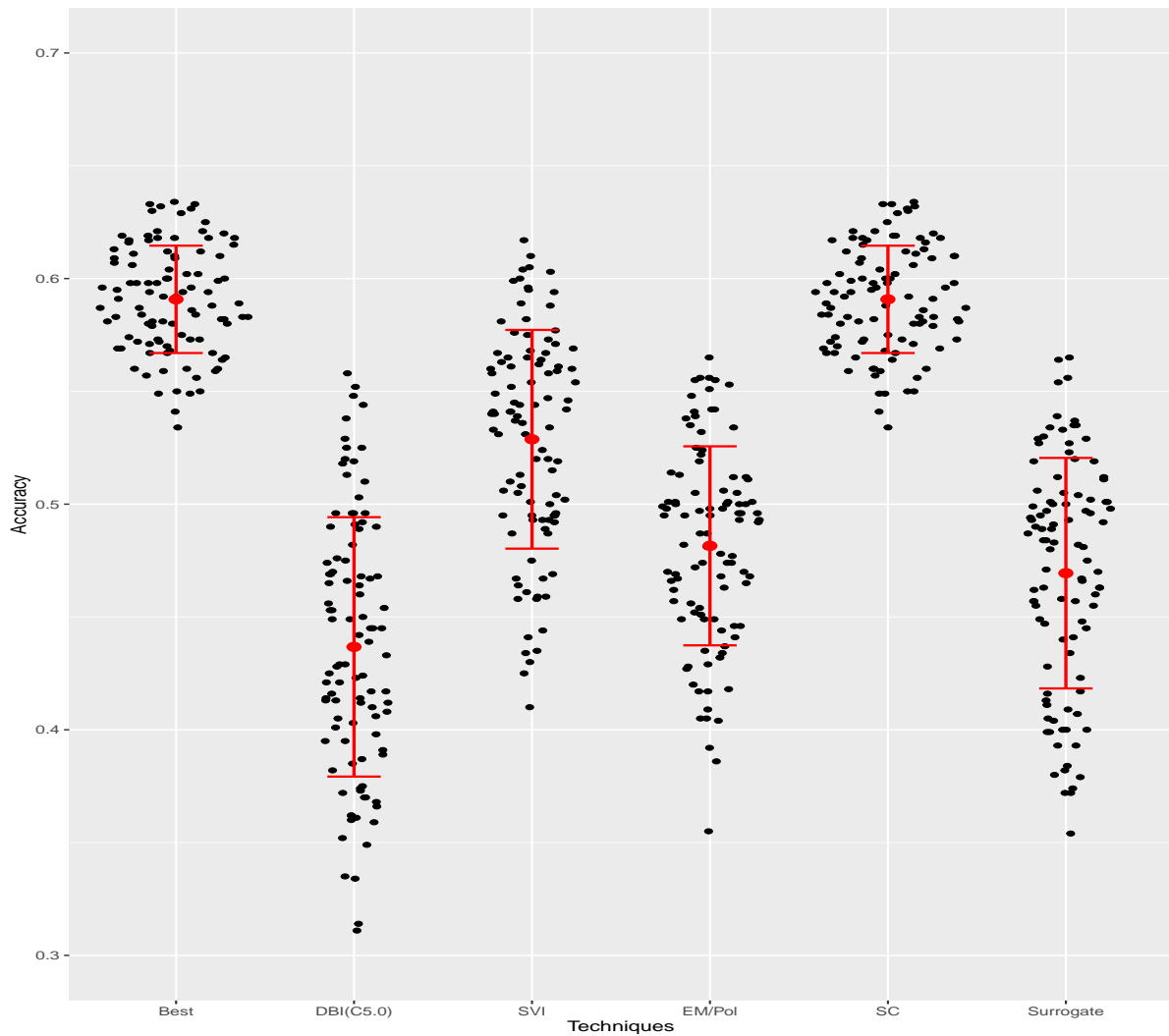


Figure 7: Accuracy by techniques when the missing pattern depends on the response.

In this experiment, the missing pattern is actually a variable with predictive power and therefore, models like BEST and SC shine as they utilize the fact that there is missing values instead of trying to impute them. BEST and SC approaches have similar results and their performances is higher than any other techniques.

It is interesting to notice the high performances of the simple single value imputation in some cases. Our experiments revealed that when the predictor containing missing value is continuous, replacing missing values with the mean is equivalent to creating a separate class because only the missing values will exactly take the value of the mean. If the predictor with missing value is categorical, replacing missing values with the mode will make the observation with missing value undistinguishable from observations that truly belong to that class which leads to the poor results.

Next, we want to understand the impact of the indicator gating variable on the success of BEST. To do so, we added the indicator variable to the data set and we applied the various missing value techniques as well.

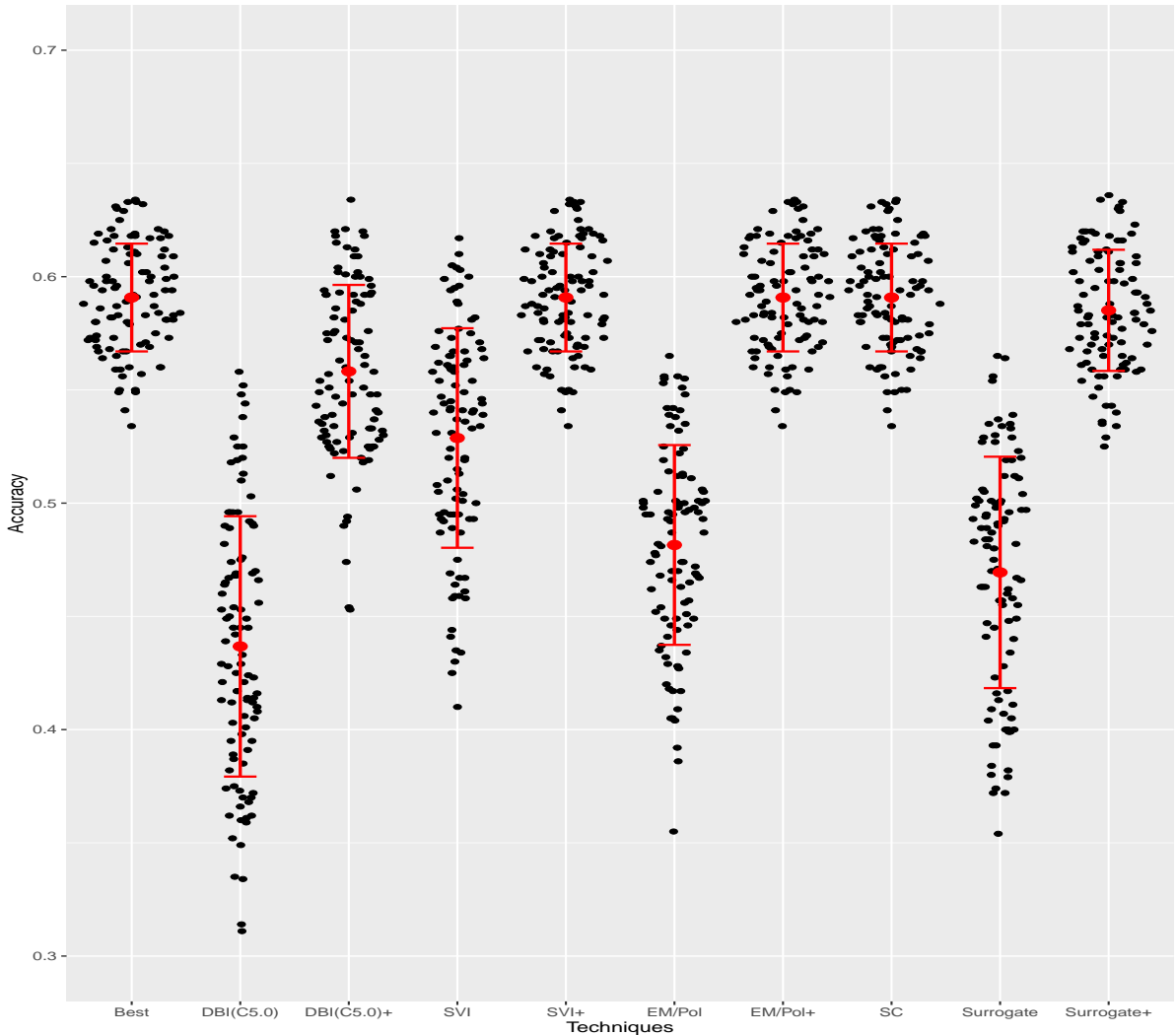


Figure 8: Accuracy by techniques when the missing pattern depends on the response

In figure 8 we can see that the indicator variable is important if the missingness depends on the response variable. In the figure above, the plus sign represents the result obtained by the technique with a data set that also includes the indicator variable. For example, for SVI+, we both imputed the missing values using SVI and we included the indicator variable in the data set. Adding this variable to the data set lead to improved performances for all of the tested techniques. BEST offers great results considering it does not need to impute the missing values. If the missing indicator variable is to be added to the data set, it is counter-intuitive to also impute missing values. We argue in sections 7.1 and 6.4 that BEST also produces trees that are more interpretable and that BEST leads to a more reliable variable importance analysis than other algorithms. Thus, we believe BEST is the best approach when the pattern is missing at random given the response.

6.3 MNAR : Missingness depends on missing values

Let us now proceed with an experiment when the predictors are missing not at random. If a continuous predictor, let us say X_j , is randomly selected, then a random value t which serves as threshold within the domain of X_j is also drawn at random. Finally, a Bernoulli variable b is drawn and if $b = 0$, then if $X_j < t$ X_j is set missing, otherwise if $b = 1$ then X_j is missing if its value is greater than t .

If one of the binary predictors is selected, then a Bernoulli variable is drawn and X_j is missing if X_j equals the Bernoulli variable. Since the missingness of X_j depends on the value of X_j itself, then this is MNAR. The procedure was repeated 100 times.

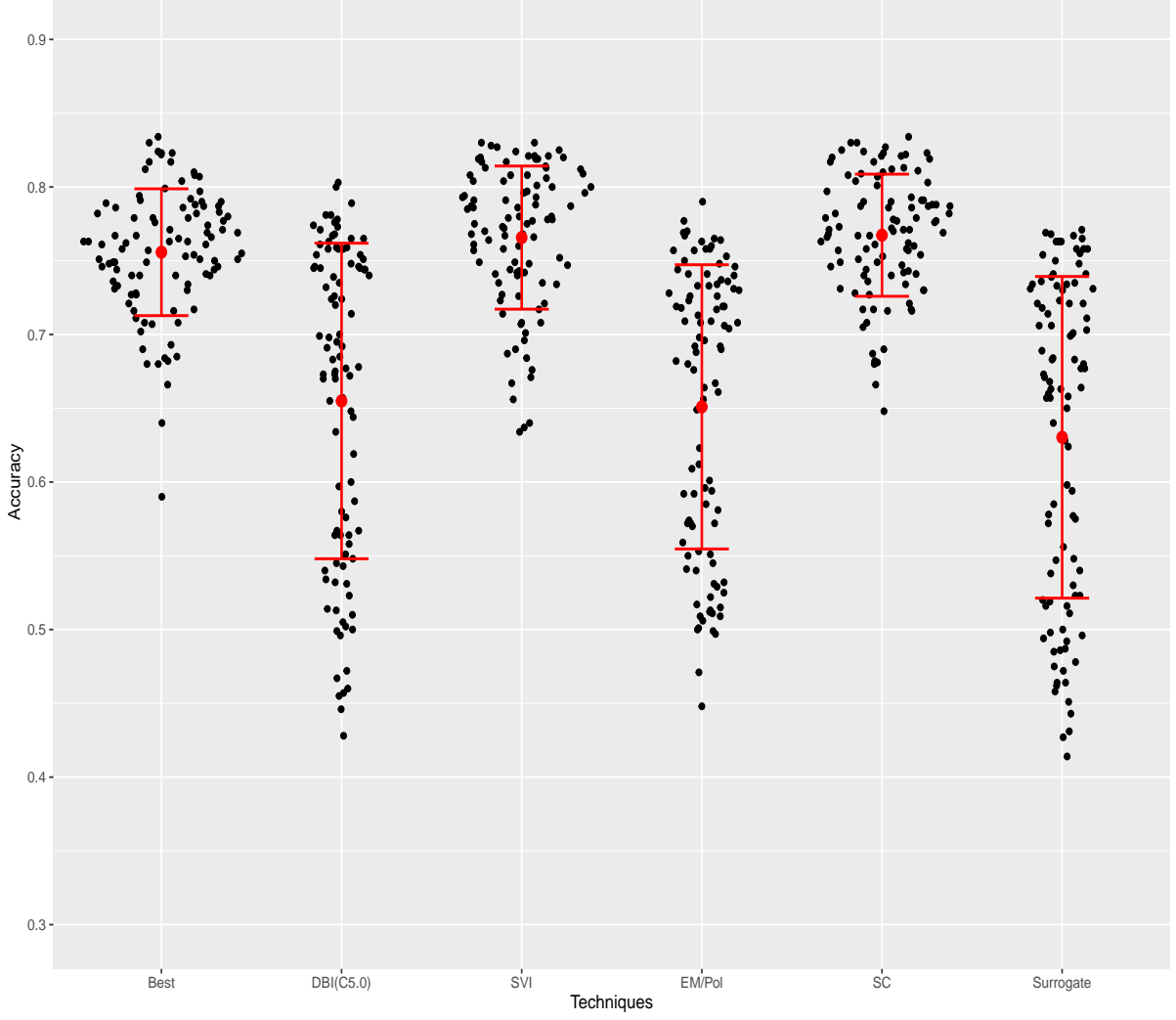


Figure 9: Accuracy by techniques when the missing pattern depends on missing values

In figure 9 we observe that BEST outperforms DBI and multiple imputations. The performances of BEST are comparable to those of SVI but are slightly under those of the SC approach.

6.4 Random forests and variable importance

Let us now build a small example where random forests are used to analyse the variable importance. Random forests are popular in exploratory analysis [24] as the variable importance tools that were developed for this model became popular.

As we have seen in the previous experiments, when BEST performs well, so does the SC approach usually as both of these techniques produce similar trees. Here we will

quickly discuss how BEST produces more accurate variable importance computations than the SC approach. We have created an example where the missing pattern depends on the response, used either the SC approach or BEST to handle missing values and we have built a forest of those trees.

When the values for a predictor are conditionally missing at random given the response, the missing pattern is itself a good predictor. We believe it is important that a variable importance analysis distinguishes between the importance of the predictor with missing value, say X_j , from the importance of its missing pattern M_j . A random forest of trees built under the SC approach would fail to distinguish between the effect of the observed value for that predictor and the effect of the missing pattern. Since BEST actually uses a variable to define the region with missing values, either with another predictor or a user-created dummy variable, then this gating variable importance will better represent the predictive power of the missing pattern.

Data	X_1	X_2	X_3	X_4	M_2
Complete	112.37705	26.76542	158.55069	102.91328	-
BEST	94.53228	10.80073	155.39201	78.14924	156.27593
SC	103.62706	171.70715	147.89474	80.51999	-
SC+.	94.36789	81.83499	128.78712	55.14764	84.45193

Table 3: Variable Importance table computed using the GINI decrease importance

We have built a random forest using the complete data set and computed the GINI decrease importance. Then we have randomly selected one of the eight labels, and the predictor X_2 , a predictor of low importance according to the GINI decreases under the complete data set, is rendered missing depending on the value of the response. Since the SC approach uses the predictor containing missing value to identify observations containing missing value then it identifies X_2 as the most important predictor. Even when we include the missing pattern M_2 to the data set the SC approach uses, the variable importance for X_2 is still inflated.

Using BEST, we can easily observe that the missing pattern, M_2 is the important pre-

dictor and that X_2 is actually of low importance when observed as it should be according to the complete data variable importance. We believe the variable importance analysis produce by BEST reflects better the true predictive power of the predictors and that this is a great benefit from using BEST over the SC approach.

6.5 Simulations takeaways and limitations

Throughout various simulation experiments we have been able to highlight the success of BEST under various scenarios. Regarding the SC approach, since it does not impute the missing values then it can create the same partitioning that BEST creates but the SC approach does not need to first isolate the missing value in order to partition upon the variable containing missing values which can sometimes be valuable.

Our simulation revealed that BEST suffers from a weakness when the gating variable is of low importance. This can happen if only a small proportion of data is missing, if the missing pattern is simply non-informative or in some cases when data is MCAR. In that case, BEST will never partition upon the gating variable and thus will never partition upon the branch-exclusive variable which will almost surely reduce the accuracy of the resulting tree. This weakness is intrinsic to the algorithm as it is caused by the greedy nature of decision trees. Since it can only see the reduction in impurity gained with a single partition, a classic decision tree approach cannot perceive the accuracy gained by the combination of two successive partitioning. Considering pairs of consecutive splits would be a great improvement and would negate this limitation but would drastically increase the run time of the algorithm.

A limitation of our simulations is the absence of the MIA algorithm [28] who shares a lot of similarities with BEST. No implementation or package of MIA was available at the time we ran this experiment and thus an advantage of BEST is that a R package is available for researcher looking for an algorithm ready to use off-the-shelf. Regarding prediction accuracy, we do not expect BEST to outperform MIA as they can produce similar partitioning but we expect BEST to be slightly faster as MIA greatly increases

the number of operations when building the tree. Finally, as discussed earlier we do believe BEST produces more interpretable trees than the SC approach, the same argument holds when comparing BEST to MIA. The interpretability argument is expanded in section 7.1.

Finally, we did not test the algorithm run time for two reasons. First, our code is entirely written in R and is not optimized. Second and more importantly, there is no theoretical reason to believe that BEST is either slower or faster than a traditional decision tree algorithm as the number of operations is of the same order.

7 Experiments : grades data set

The data set mentioned in section 5.1 was analysed using BEST. Once again, the accuracy of the proposed algorithm is compared to other techniques that handle missing values. To begin, we predict if a student completes its program using its first year of courses and results. The data set contains 38842 observations. Our set of predictors consists of the number of credits attempted in all the departments and the average grade obtained in those respective departments. The number of credits is a numerical variable that serves as the gating variable for the respective average grade. If the number of credits attempted in a department is greater than 0 for every observation in a region then BEST acquires access to the grade variable. We have randomly sampled training sets of different sizes and used all the remaining observations to assess the accuracy. This process was repeated 15 times and we have averaged the results. We did not include the single value imputation because we expect this technique to produce the same result as the SC approach since all predictors are numerical. We did not include the imputation produced by the mice package [29] as the package was incapable of handling the data set. Finally, we have used tables to show our results because with only 10 trials the Sina plots were not informative.

# of obs	5000		10000		15000		20000	
Methods	Mean	S.D.	Mean	S.D.	Mean	S.D.	Mean	S.D.
DBI	0.7223	0.0035	0.7336	0.0038	0.7385	0.0092	0.7402	0.0095
SC	0.7307	0.0066	0.7387	0.0044	0.7427	0.0036	0.7462	0.0033
Surrogate	0.7291	0.0053	0.7301	0.0046	0.7311	0.0034	0.7300	0.0032
BEST	0.7333	0.0062	0.7424	0.0045	0.7457	0.0037	0.7479	0.0033

Table 4: Mean accuracy and standard deviation when predicting program completion

In table 4, we observe that BEST produces the most accurate decision trees for that data set for all training data sizes.

Another reason why we might prefer using BEST in this analysis is its ability to rightfully identify the variable importance. As we discussed in section 5.1, the researchers were interested in the importance of the predictors. Therefore BEST is an improvement as it truly identifies the importance of the gating variables, the number of credits, as we have shown in section 6.4. In this case we were able to distinguish the importance of the number of credits in a department of the importance of the grade obtained in that department.

In the second analysis, we will look at the 26488 students who completed their program. Using the same set of explanatory variables, we will try to predict the department they majored in :

# of obs	5000		10000		15000		20000	
Methods	Mean	S.D.	Mean	S.D.	Mean	S.D.	Mean	S.D.
DBI	0.3581	0.0055	0.376	0.0043	0.3871	0.0039	0.3936	0.0057
SC	0.3992	0.0079	0.4172	0.0069	0.4262	0.0047	0.4332	0.0044
Surrogate	0.3639	0.0121	0.3661	0.011	0.3723	0.012	0.3742	0.012
BEST	0.3952	0.0074	0.4164	0.0051	0.4265	0.0043	0.4319	0.0043

Table 5: Mean accuracy and standard deviation when predicting the major

In table 5 we observe closer results from the two best-performing algorithms. BEST and SC have almost indistinguishable performances and are the top performers. Even though the results are a lot closer between BEST and SC, our proposed algorithm still

produces trees that are more interpretable and could be used to produce a non-biased variable importance analysis.

7.1 Improved interpretability

We have mentioned throughout this article that we believe BEST leads to more interpretable decision trees than the SC approach. The experiment on this data set provides a good example for that.

For the SC approach we have replaced missing grades by a value outside of the domain, 101. Frequently in this experiment, the tree constructed under the SC approach partitions upon the grade in departments before the number of credits. For example, if *Grades in Mathematics* is the first split variable selected and 60 is selected as the split point, then the SC approach produces the partitioning illustrated in figure 10.

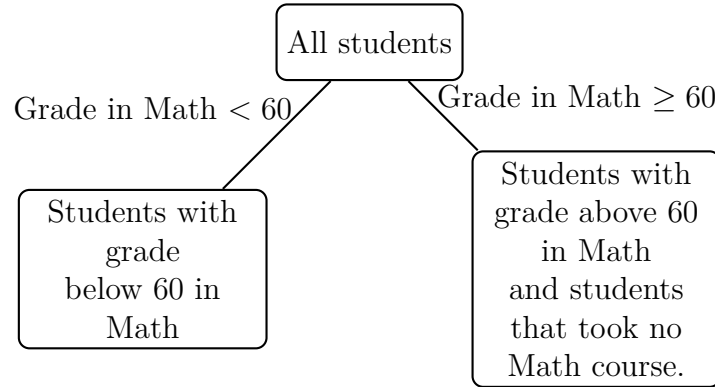


Figure 10: An example of a decision tree partitioning produced by the SC approach and the associated regions

It is definitely hard to extract interpretable information out of this tree. Does this partition implies students with no experience in Mathematics behave similarly to students with good results in Mathematics ? BEST achieves similar or higher accuracy while keeping the partitions logical and interpretable. BEST will begin by partitioning students who attempted at least 1 credit in Mathematics from those who did not. Then, among students who attempted at least 1 credit in Mathematics, BEST will partition them according to their grades, which leads to a more interpretable sequence of partitions as illustrated in figure 11.

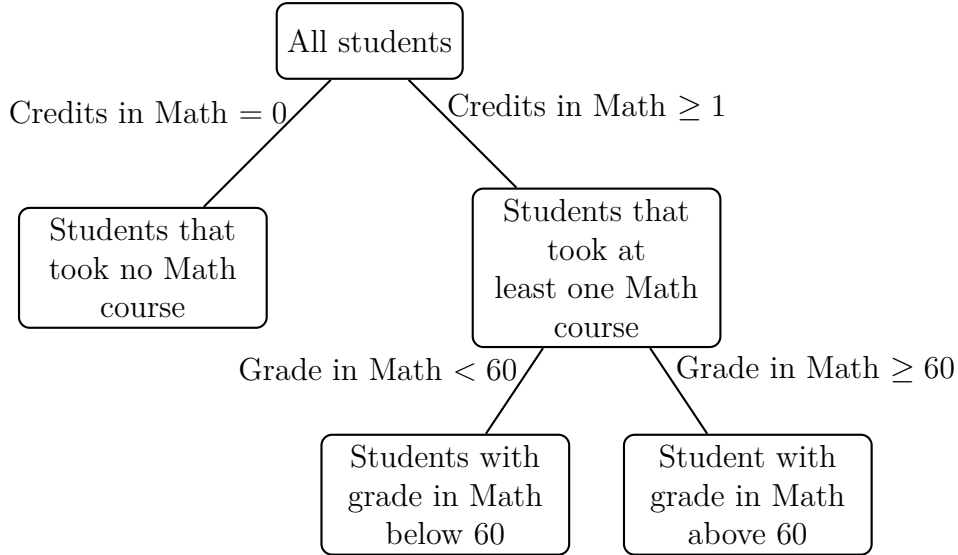


Figure 11: An illustration of a decision tree partitioning produced by BEST and the associated regions

If interpretability is considered a strength of decision trees, then BEST is better than the SC approach at preserving this strength.

7.2 Real-world data set experiment takeaways

Even though we have performed experiments on a single real-world data, the results are extremely positive. BEST has higher or similar performance than the other tested algorithms. BEST produces more accurate variable importance analysis and more interpretable trees than the SC approach, BEST's closest competitor. Finally, BEST did not require any imputations which is another reason why we prefer BEST.

8 Conclusion

We have constructed a modified tree-building algorithm that lets the users decide the regions of the predictor space where variables are available for the data partitioning process. We have focused on using this feature to manage missing values. BEST has the elegant property of analysing a variable only when values are known without assuming any missingness structure. It produces highly interpretable trees and achieves comparable accuracy to most missing value handling techniques in cases we have identified using

simulated data sets. Even though BEST shares similarities with the separate class technique, BEST leads to a more accurate variable importance analysis and produces more interpretable and intuitive trees.

BEST suffers from a weakness when the gating variable has no predictive power. In those cases, the algorithm will never choose to split upon the gating variable and thus will never be allowed to use the branch-exclusive variable. This problem can lead to a decrease in accuracy in some simple cases where the data is MCAR. Fortunately, as we have previously discussed, there already exist multiple techniques to handle data MCAR and we can count on cross-validation in order to help us select the best missing data handling technique. Nonetheless, in the simulated experiments we have performed, results were mostly positive as BEST outperforms some other techniques when data is MAR and MNAR.

The results produced by BEST were also satisfactory when the algorithm was used on the real motivating *grades data set*. We were able to achieve higher accuracy than with most other techniques while obtaining a more interpretable classifier. Since variable importance was a concern in the *grades data set* analysis, BEST was an improvement as it answers that research question by providing a more reliable variable importance analysis than the separate class approach previously used [2].

Acknowledgements

We are very grateful to Glenn Loney and Sinisa Markovic of the University of Toronto for providing us with the anonymised students grade data. The authors also gratefully acknowledge the financial support from NSERC of Canada.

References

- [1] Michael A. Bailey, Jeffrey S. Rosenthal, and Albert H. Yoon. Grades and incentives: assessing competing grade point average measures and postgraduate outcomes. *Stud-*

- ies in Higher Education*, 41(9):1548–1562, 2016.
- [2] C. Beaulac and J. S. Rosenthal. Predicting University Students’ Academic Success and Choice of Major using Random Forests. *ArXiv e-prints*, February 2018.
 - [3] L. Breiman, J. Friedman, R. Olshen, and C. Stone. *Classification and Regression Trees*. Wadsworth and Brooks, Monterey, CA, 1984. new edition [?]?
 - [4] Leo Breiman. Bagging predictors. *Machine Learning*, 24(2):123–140, 1996.
 - [5] Leo Breiman. Random forests. *Machine Learning*, 45(1):5–32, 2001.
 - [6] Yufeng Ding and Jeffrey S. Simonoff. An investigation of missing data methods for classification trees applied to binary response data. *J. Mach. Learn. Res.*, 11:131–170, March 2010.
 - [7] A. J. Feelders. Handling missing data in trees: Surrogate splits or statistical imputation. In *PKDD*, 1999.
 - [8] Jerome Friedman, Ron Kohavi, and Yeogirl Yun. Lazy decision trees. 1, 09 1997.
 - [9] S. Gavankar and S. Sawarkar. Decision tree: Review of techniques for missing values at training, testing and compatibility. In *2015 3rd International Conference on Artificial Intelligence, Modelling and Simulation (AIMS)*, pages 122–126, Dec 2015.
 - [10] Pierre Geurts, Damien Ernst, and Louis Wehenkel. Extremely randomized trees. *Machine Learning*, 63(1):3–42, Apr 2006.
 - [11] T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning*. Springer, 2 edition, 2009.
 - [12] H. Kim and W.-Y. Loh. Classification trees with unbiased multiway splits. *Journal of the American Statistical Association*, 96:589–604, 2001.
 - [13] Max Kuhn and Ross Quinlan. *C50: C5.0 Decision Trees and Rule-Based Models*, 2018. R package version 0.1.2.

- [14] Joseph L. Schafer and Maren K. Olsen. Multiple imputation for multivariate missing-data problems: A data analyst’s perspective. 33, 07 2000.
- [15] Roderick J. A. Little and Donald B. Rubin. *Statistical Analysis with Missing Data*. John Wiley & Sons, Inc., 2002.
- [16] J. Ross Quinlan. *C4.5: Programs for Machine Learning*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1993.
- [17] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2018.
- [18] Md. Geaur Rahman and Md Zahidul Islam. Missing value imputation using decision trees and decision forests by splitting and merging records: Two novel techniques. *Knowledge-Based Systems*, 53:51 – 65, 2013.
- [19] Donald B. Rubin. Inference and missing data. *Biometrika*, 63(3):581–592, 1976.
- [20] Maytal Saar-Tsechansky and Foster Provost. Handling missing values when applying classification models. *J. Mach. Learn. Res.*, 8:1623–1657, December 2007.
- [21] Shaun Seaman, John Galati, Dan Jackson, and John Carlin. What is meant by missing at random? *Statist. Sci.*, 28(2):257–268, 05 2013.
- [22] Shai Shalev-Shwartz and Shai Ben-David. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, New York, NY, USA, 2014.
- [23] Nikos Sidiropoulos, Sina Hadi Sohi, Nicolas Rapin, and Frederik Otzen Bagger. Sinaplot: an enhanced chart for simple and truthful representation of single observations over multiple classes. *bioRxiv*, 2015.
- [24] Carolin Strobl, Anne-Laure Boulesteix, Achim Zeileis, and Torsten Hothorn. Bias in random forest variable importance measures: Illustrations, sources and a solution. *BMC Bioinformatics*, 8(1):25, 2007.

- [25] Terry Therneau and Beth Atkinson. *rpart: Recursive Partitioning and Regression Trees*, 2018. R package version 4.1-13.
- [26] Nicholas J Tierney, Fiona A Harden, Maurice J Harden, and Kerrie L Mengersen. Using decision trees to understand structure in missing data. *BMJ Open*, 5(6), 2015.
- [27] Bhekisipho Twala. An empirical comparison of techniques for handling incomplete data using decision trees. *Appl. Artif. Intell.*, 23(5):373–405, May 2009.
- [28] Bhekisipho Twala, M.C. Jones, and David Hand. Good methods for coping with missing data in decision trees. 29:950–956, 05 2008.
- [29] Stef van Buuren and Karin Groothuis-Oudshoorn. mice: Multivariate imputation by chained equations in r. *Journal of Statistical Software*, 45(3):1–67, 2011.
- [30] Hadley Wickham. *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York, 2016.