

No Imputation without Representation

Oliver Urs Lenz^{1,2}[0000–0001–9925–9482], Daniel Peralta³[0000–0002–7544–8411],
and Chris Cornelis²[0000–0002–6852–4041]

¹ Leiden Institute of Advanced Computer Science, Leiden University
`o.u.lenz@liacs.leidenuniv.nl`

² Research Group for Computational Web Intelligence,
Department of Applied Mathematics, Computer Science and Statistics,
Ghent University

`chris.cornelis@ugent.be` <https://cwi.ugent.be>

³ Department of Information Technology, Ghent University
`daniel.peralta@ugent.be`

Abstract By filling in missing values in datasets, imputation allows these datasets to be used with algorithms that cannot handle missing values by themselves. However, missing values may in principle contribute useful information that is lost through imputation. The missing-indicator approach can be used in combination with imputation to instead represent this information as a part of the dataset. There are several theoretical considerations why missing-indicators may or may not be beneficial, but there has not been any large-scale practical experiment on real-life datasets to test this question for machine learning predictions. We perform this experiment for three imputation strategies and a range of different classification algorithms, on the basis of twenty real-life datasets. In a follow-up experiment, we determine attribute-specific missingness thresholds for each classifier above which missing-indicators are more likely than not to increase classification performance. And in a second follow-up experiment, we evaluate numerical imputation of one-hot encoded categorical attributes. We reach the following conclusions. Firstly, missing-indicators generally increase classification performance. Secondly, with missing-indicators, nearest neighbour and iterative imputation do not lead to better performance than simple mean/mode imputation. Thirdly, for decision trees, pruning is necessary to prevent overfitting. Fourthly, the thresholds above which missing-indicators are more likely than not to improve performance are lower for categorical attributes than for numerical attributes. Lastly, mean imputation of numerical attributes preserves some of the information from missing values. Consequently, when not using missing-indicators it can be advantageous to apply mean imputation to one-hot encoded categorical attributes instead of mode imputation.

Keywords: Missing data · Missing-indicators · Imputation · Classification · Data-centric machine learning.

1 Introduction

Missing values are a frequent issue in real-life datasets, and the subject of a large body of ongoing research. Some implementations of machine learning algorithms can handle missing values natively, requiring no further action by practitioners. But whenever this is not the case, a common general strategy is to replace the missing value with an estimated value: imputation. An advantage of imputation is that we obtain a complete dataset, to which we can apply any and all algorithms that make no special provision for missing values. However, missing values may be informative, and a disadvantage of imputation is that it removes this information.

The missing-indicator approach [12] is an old proposal to represent and thereby preserve the information encoded by missing values. For every original attribute, it adds a new binary ‘indicator’ or ‘dummy’ attribute that takes a value of 1 if the value for the original attribute is missing, and 0 if not (Figure 1).⁴ The missing-indicator approach is often presented as an alternative to imputation, but since it does not resolve the missing values in the original attributes, it can only be used in addition to, not instead of imputation.

a_1	a_2	a_3	i_1	i_2	i_3
0.92	aap	2.50	0	0	0
?	aap	1.00	1	0	0
8.42	?	3.00	0	1	0
2.23	noot	0.05	0	0	0
?	?	?	1	1	1
0.41	mies	?	0	0	1
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots

Figure 1: Illustrative example of a dataset with three attributes (a_1, a_2, a_3) with missing values (‘?’), and corresponding missing-indicators (i_1, i_2, i_3).

It is an open question whether missing-indicators should be used for predictive tasks in machine learning [75]. Both imputation and the missing-indicator approach originate in the statistical literature. While imputation strategies have been the subject of a rich body of research, the missing-indicator approach has not received a large amount of attention, and is often dismissed or disregarded in overviews of approaches towards missing values.

In the context of machine learning, the effect of missing-indicators can be framed as follows. On the one hand, the addition of missing-indicators results in a more complete, higher-dimensional representation of the data. On the other hand, their omission corresponds to a form of dimensionality reduction, which

⁴ Some authors use the opposite convention, letting the indicator express non-missingness.

may increase the efficiency and effectiveness of a dataset by eliminating redundancy.

To determine whether this trade-off is useful, a key question is to which extent missing values in a given dataset are informative. If they are not, the phrase “missing at random” (MAR) [66] is used to indicate that the distribution of missing values is dependent on the known values, while the stricter phrase “missing completely at random” (MCAR) denotes values that are distributed truly randomly. In contrast, informative missing values are often denoted as “missing not at random” (MNAR).

In this respect, it is often argued that one should distinguish between missing values that could in principle have been obtained, and missing values that fundamentally do not exist, like attributes related to pregnancy tests for male subjects.⁵ In the latter case, the missing values are definitely informative. However, such clear-cut cases may be comparatively rare. Moreover, it does not follow that the missing values in the former case are definitely non-informative. In fact, for real-life datasets, unless we have specific knowledge about the process responsible for the missing values, we have to assume some degree of informativeness in principle.⁶

Nonetheless, it has been argued that in practice, the attributes of a dataset can be sufficiently redundant that one can get away with assuming that its missing values are MAR [69]. This means that most of the information contained by the missing values should in principle be recoverable through imputation. But even if this is so, imputation may not always perform optimally, in which case missing-indicators may still prove useful for machine learning.

A more subtle point is that even when missing values are informative, the information they encode need not be lost completely through imputation. This is particularly evident in the case of numerically encoded binary attributes (e.g. 0 and 1), where imputation can represent missing values as a third, intermediary value (e.g. 0.5). More generally, Le Morvan et al. [49] have recently observed that almost all deterministic imputation functions map records with missing values to distinct manifolds in the attribute space that can in principle be identified by sufficiently powerful algorithms. Nevertheless, missing-indicators can potentially make this learning task easier.

In light of these conflicting theoretical arguments, the usefulness of missing-indicators for real-life machine learning problems is an interesting empirical question. However, previous experiments in this direction have been limited in scope and number. These limitations include the use of only one or a handful of datasets, the use of datasets from which values have been removed artificially, and not comparing the same imputation strategies with and without missing-indicators.

⁵ We are grateful to an anonymous reviewer for this example.

⁶ This is acknowledged by authors working under the assumption of MAR, e.g. “When data are missing for reasons beyond the investigator’s control, one can never be certain whether MAR holds. The MAR hypothesis in such datasets cannot be formally tested unless the missing values, or at least a sample of them, are available from an external source.” [69]

The purpose of the present paper is straightforward. On the basis of twenty real-life classification problems with naturally occurring missing values, we measure the performance of a range of popular classification algorithms, using three common types of imputation, with and without missing-indicators. This allows us to evaluate the effect of using missing-indicators, as well as the choice of imputation strategy.

Moreover, we conduct three follow-up experiments to gain a better understanding of when and why missing-indicators can be useful. In the first, we determine whether this is influenced by the type (categorical or numerical) and the amount of missing values of a given attribute. In the second, we test the hypothesis that numerical imputation partially preserves the information from missing values. And in the third follow-up experiment, we compare missing-indicators to two model-specific approaches to missing values for nearest neighbour and decision tree classifiers.

In Section 2, we provide a brief overview of the existing literature on missing-indicators, including previous experimental evaluations. In Section 3, we describe our experimental setup. We report our results in Section 4 and conclude in Section 5.

2 Background

We start with a brief discussion of the origins and reception of the missing-indicator approach, as well as previous experimental evaluations of the use of missing-indicators in prediction tasks.

2.1 Origins and Reception

The missing-indicator approach originates in the literature on linear regression. It dates back to at least Cohen [12], who pointed out that values in real-life datasets are typically not missing completely at random, and that the distribution of missing values may in particular depend on the values of the attribute that is to be predicted. He proposed that each attribute could be said to have two ‘aspects’, its value, and whether that value is present to begin with, which should be encoded with a pair of variables. For missing attribute values, the first of these variables was to be filled in with the mean of the known values, although other applications might call for different values. Cohen’s proposal was subsequently expanded in [13], but received only limited recognition in the following years [46,76,11,42,4,57].

Cohen’s proposal was subjected to a formal analysis by Jones [44], who showed that, if one assumes that missing values are MAR, and the true linear regression model does not contain any terms related to missingness, it produces biased estimates of the regression coefficients (unless the sample covariance between independent variables is zero). However, these assumptions run directly counter to the position set out in [13] that a priori, the missingness of each attribute is a possible explanatory factor, that it is safer not to assume that missing

values are distributed randomly, and that the usefulness of missing-indicators is ultimately an empirical question.

Allison [2], motivated by [44] and working under the general assumption of MAR, dismissed missing-indicators as “clearly unacceptable”, before conceding that they in fact produce optimal estimates when the missing value is not just missing, but cannot exist, such as the marital quality of an unmarried couple. However, this semantic distinction may not always be clear-cut in practice, and the more pertinent question may be whether missing values are informative. Allison [3] later acknowledged that missing-indicators may lead to better predictions and their use for that purpose was acceptable. Missing-indicators have also been dismissed in [61,70,36,5], and are frequently omitted in overviews of missing data strategies [69,25,23,32,16].

2.2 Previous Experiments

Only a handful of experimental comparisons of missing data approaches have included the missing-indicator approach, and these have been limited in scope. [81] and [56] only use a single dataset with randomly removed values, and base their evaluation on the performance of a single algorithm (respectively a neural network and linear regression). The authors of [59] use three classification algorithms and 22 datasets, but again with randomly removed values, explicitly assuming an MCAR context. They conclude that imputation outperforms missing-indicators, but the comparison is not like-for-like, since it involves several forms of imputation but only combines indicator attributes with zero imputation. The authors of [41] compare missing-indicators with zero imputation against several other forms of imputation without missing-indicators on one real dataset, for logistic regression. However, they do not evaluate predictive performance.

Ding & Simonoff [18] conduct a more extensive investigation, using insights from a series of Monte Carlo simulations to systematically remove values from 36 datasets to simulate different forms of missingness. They use these datasets to compare zero imputation⁷ with indicator attributes against mean/mode imputation without, as well as a number of other missing data approaches, for logistic regression. In addition, the authors evaluate a related representation of missing values⁸ on the same set of 36 datasets, and on one real-life dataset with missing values, for decision trees. They find that there is strong evidence that representing missing values is the best approach when they are informative; when this is not the case their results show no strong difference.

The comparison by Grzymala-Busse & Hu [38] is based on 10 datasets with naturally occurring missing values. However, the setting is purely categorical

⁷ Presumably, they use one-hot encoding for categorical attributes, in which case zero imputation is equivalent to treating missing values as a separate category, but they do not state this explicitly.

⁸ For categorical values, encoding missing values as a separate category; for numerical values, encoding missing values as an extremely large value that can always be split from the other values.

— all attributes are transformed into categorical attributes — the only form of imputation is mode imputation, and the missing value approaches are evaluated on the basis of the LERS classifier (Learning from Examples based on Rough Sets [37]).

Marlin [52] compares zero imputation with missing-indicators (*augmentation with response indicators*) against several forms of imputation without, for logistic regression and neural networks, on the basis of an extensive series of simulations, one dataset with artificially removed values, and three real datasets. For the real datasets, there is no strong difference in performance between the different approaches.

Most recently, building on earlier experiments with simulated regression datasets [45,49], Perez-Lebel et al. [60] compare four different imputation techniques with and without missing-indicators (*missingness mask*) on seven prediction tasks derived from four real medical datasets, and conclude that missing-indicators consistently improve performance for gradient boosted trees, ridge regression and logistic regression.

We point out that the Missingness in Attribute (MIA) proposal [80] for decision trees and decision tree ensembles can be understood as an implicit combination of missing-indicators with automatic imputation, and has also been shown to outperform imputation without missing-indicators in small-scale experimental studies [45,60].

Finally, even experimental comparisons of missing data that do not feature the missing-indicator approach generally do not involve more than a handful of real-life datasets with naturally occurring missing values. We have only found the connected works [50,51], which feature 21 datasets from the UCI repository, but 12 of these are problematic.⁹

⁹ The target column of the *echocardiogram* dataset ('alive-at-1') is supposed to denote whether a patient survived for at least one year, but it doesn't appear to agree with the columns from which it is derived, that denote how long a patient (has) survived and whether they were alive at the end of that period. The *audiology* dataset has a large number of small classes with complex labels and should perhaps be analysed with multi-label classification. In addition, it has ordinal attributes where the order of the values is not entirely clear, and three different values that potentially denote missingness ('?', 'unmeasured' and 'absent'), and it is not completely clear how they relate to each other. The *house-votes-84* dataset contains '?' values, but its documentation explicitly states that these values are not unknown, but indicate different forms of abstention. The *ozone* dataset is a time-series problem, while the task associated with the *sponge* and *water-treatment* datasets is clustering, with no obvious target for classification among their respective attributes. Finally, the *breast-cancer* (9), *cleveland* (7), *dermatology* (8), *lung-cancer* (5), *post-operative* (3) and *wisconsin* (16) datasets contain only very few missing values, and any performance difference between missing value approaches on these datasets may to a large extent be coincidental.

3 Experimental Setup

To evaluate the effect of the missing-indicator approach on classification performance, we conduct a series of experiments, using the Python machine learning library *scikit-learn* [58].

3.1 Questions

The aim of our experiments is to answer the following questions:

- Do missing-indicators increase performance, and does it matter which imputation strategy they are paired with?
- When do missing-indicators start to become useful in terms of missingness?
- Does using mean imputation instead of mode imputation allow for more information to be learned from missing categorical values?
- How do missing-indicators compare to model-specific approaches to missing values?

3.2 Evaluation

We preprocess datasets by standardising numerical attributes and one-hot encoding categorical attributes (as required by the implementations in *scikit-learn*).

We measure classification performance by performing stratified five-fold cross-validation, repeating this for five different random states (which determine both the dataset splits and the initialisation of algorithms with a random component), and calculating the mean area under the receiver operator curve (AUROC). For multi-class datasets, we use the extension of AUROC defined in [40].

To compare two alternatives A and B, we consider the p -value of a one-sided Wilcoxon signed-rank test [82] on the mean AUROC scores for our selection of datasets. When we compare A vs B, a score below 0.5 means that A increased performance on our selection of datasets; the lower the scores, the more confident we can be that this generalises to other similar datasets. Conversely, a score higher than 0.5 means that A decreased performance on our selection of datasets.

3.3 Imputation Strategies

We consider the following three imputation strategies:

- *Mean/mode imputation* replaces missing values of numerical and categorical attributes by, respectively, the mean and the mode of the non-missing values.
- *Nearest neighbour imputation* [79] replaces missing values of numerical and categorical attributes by, respectively, the mean and the mode of the 5 nearest non-missing values, with distance determined by the corresponding non-missing values for the other attributes.

- *Iterative imputation*, as implemented in scikit-learn, based on [8], predicts missing values of one attribute on the basis of the other attribute values using a round-robin approach. For numerical attributes, this uses Bayesian ridge regression [77], initialised with mean imputation, while for categorical attributes, we use logistic regression, initialised with mode imputation.

The scikit-learn implementations of nearest neighbour and iterative imputation can currently only impute numerical features, so we had to adapt them for categorical imputation. In all other aspects, we follow the default settings of scikit-learn.¹⁰

3.4 Classification Algorithms

We consider the classification algorithms listed in Table 1, as implemented in scikit-learn. Hyperparameters take their default values, except for SVM-L, LR and MLP, where we increase the maximum number of iterations to 10 000 to increase the probability of convergence.

For a number of these algorithms, specific ways have been proposed to handle missing values: e.g. NN-2-D [19], SVM-G [72], MLP [78,73,43] and CART [63,80]. The purpose of the present experiment is to evaluate the general approach of using imputation with missing-indicators when these solutions have not been implemented, as is the case in scikit-learn.

Table 1: Classification algorithms.

Name	Description
NN-1	Nearest neighbours [28] with (Bosovich) 1-distance
NN-2	Nearest neighbours with (Euclidean) 2-distance
NN-1-D	Nearest neighbours with 1-distance, distance-weighted [21]
NN-2-D	Nearest neighbours with 2-distance, distance-weighted
SVM-L	Soft-margin Support Vector Machine [14] with linear kernel
SVM-G	Soft-margin Support Vector Machine with Gaussian kernel
LR	Multinomial logistic regression [15]
MLP	Multilayer perceptron [65] with ReLu activation [31], Glorot initialisation [34] and Adam optimisation [47]
CART	Classification and Regression Tree [7]
RF	Random Forest [6]
ERT	Extremely Randomised Trees [33]
ABT	Ada-boosted trees [29] with SAMME (stagewise additive modeling using a multi-class exponential loss function) [83]
GBM	Gradient Boosting Machine [30]

¹⁰ For the *nomao* dataset, iterative imputation diverged, so we had to restrict imputation to the interval $[-100, 100]$.

3.5 Datasets

We use twenty real-life datasets with naturally occurring missing values from the UCI repository for machine learning [20] (Table 2). These datasets are quite varied — they cover a number of different domains and contain between 155 and 76 000 records, between 4 and 590 attributes, between 2 and 21 decision classes and missing value rates between 0.0032 and 0.43.

We have preprocessed these datasets in the following manner. We have removed attributes that were non-informative according to the accompanying documentation, as well as identifiers and alternative target values. When it was clear from the description that an attribute was categorical, we have treated it as such, even if it was originally represented with numerals. Conversely, where the possible values of an attribute admitted a semantic order, we have encoded them numerically. We have left binary attributes in their original encoding (categorical or numerical). To enable 5-fold cross-validation, we have removed classes with fewer than 5 records.

Table 2: Real-life classification datasets with missing values from the UCI repository for machine learning.

Dataset	Records	Classes	Attributes			Missing value rate			Source
			Num	Cat	Total	Num	Cat	Total	
adult	48842	2	5	8	13	0.0	0.017	0.010	[48]
agaricus-lepiota	8124	2	1	21	22	0.0	0.015	0.014	[71]
aps-failure	76000	2	170	0	170	0.083		0.083	[27]
arrhythmia	443	10	279	0	279	0.0032		0.0032	[39]
bands	540	2	19	15	34	0.054	0.054	0.054	[26]
ckd	400	2	14	10	24	0.14	0.059	0.11	[67]
crx	690	2	6	9	15	0.0060	0.0068	0.0065	[62]
dress-sales	500	2	3	9	12	0.20	0.19	0.19	
exasens	399	4	7	0	7	0.43		0.43	[74]
hcc	165	2	49	0	49	0.10		0.10	[68]
heart-disease	1611	2	13	1	14	0.18	0.0	0.17	[17]
hepatitis	155	2	19	0	19	0.057		0.057	[22]
horse-colic	368	2	19	1	20	0.25	0.39	0.26	[54]
mammographic-masses	961	2	2	2	4	0.042	0.041	0.042	[24]
mi	1700	8	111	0	111	0.085		0.085	[35]
nomao	34465	2	89	29	118	0.38	0.37	0.38	[9]
primary-tumor	330	15	16	1	17	0.029	0.20	0.039	[10]
secom	1567	2	590	0	590	0.045		0.045	[53]
soybean	683	19	22	13	35	0.099	0.096	0.098	[55]
thyroid0387	9172	18	7	16	23	0.22	0.0021	0.069	[64]

Table 3: One-sided p -values, imputation with missing-indicators versus without.

Classifier	Imputation strategy		
	Mean/mode	Neighbours	Iterative
NN-1	0.0088	0.0015	0.0017
NN-2	0.015	0.0024	0.00048
NN-1-D	0.0045	0.0019	0.0011
NN-2-D	0.0019	0.0031	0.00027
SVM-L	0.13	0.27	0.099
SVM-G	0.0032	0.0027	0.0021
LR	0.079	0.063	0.068
MLP	0.0027	0.0063	0.0056
CART	0.44	0.39	0.40
RF	0.038	0.051	0.17
ERT	0.28	0.0099	0.026
ABT	0.089	0.078	0.47
GBM	0.17	0.012	0.36

4 Results and Discussion

Using the experimental setup detailed in the previous section, we now try to answer the questions listed in Subsection 3.1.

4.1 Do Missing-Indicators Increase Performance, and Does It Matter Which Imputation Strategy They Are Paired With?

The p -values obtained by comparing imputation with and without missing-indicators are displayed in Table 3. Missing-indicators generally lead to increased performance — with the notable exception of CART, to which we return below. The more complicated imputation strategies do not result in much better results than mean/mode imputation when we pair imputation with missing-indicators (Table 4). At best, nearest neighbour and iterative imputation only lead to a modest improvement, and for many classifiers, they actually decrease performance. Therefore, we focus on mean/mode imputation for the remainder of this section.

A possible reason for the failure of missing-indicators to increase performance with CART, is that by default, the scikit-learn implementation of this classifier does not perform pruning, making it prone to overfitting. To test this hypothesis, we repeat our experiment for CART and mean imputation, but this time we apply cost complexity pruning ($\alpha = 0.01$). This clearly improves performance ($p = 0.0069$ without missing-indicators, $p = 0.015$ with missing-indicators), and now missing-indicators have a slight advantage ($p = 0.23$).

We have also taken a closer look at ERT and GBM, for which the performance increase from missing-indicators is not very significant. For ERT, this may be due to underfitting. If we increase the number of trees from the default 100 to

Table 4: One-sided p -values, missing-indicators with iterative and nearest neighbour versus mean/mode imputation.

Classifier	Imputation strategy	
	Neighbours	Iterative
NN-1	0.94	0.15
NN-2	0.78	0.19
NN-1-D	0.97	0.55
NN-2-D	0.84	0.23
SVM-L	0.53	0.61
SVM-G	0.47	0.94
LR	0.40	0.83
MLP	0.30	0.55
CART	0.69	0.79
RF	0.61	0.86
ERT	0.61	0.64
ABT	0.33	0.78
GBM	0.93	0.85

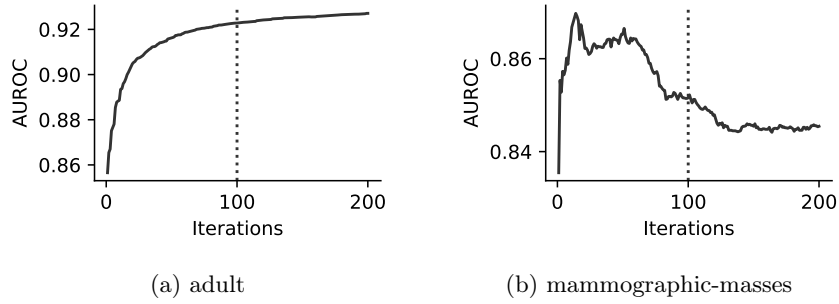


Figure 2: GBM test AUROC for two illustrative datasets, using mean/mode imputation without missing-indicators, for one random state and one cross-validation fold. The default hyperparameter value of 100 iterations leads to under- (a) and overfitting (b).

1000, this improves performance ($p = 0.0011$ without missing-indicators, $p = 0.0032$ with missing-indicators), and makes the advantage of missing-indicators somewhat clearer ($p = 0.092$).

For GBM, the default choice of 100 iterations of gradient descent can lead to both under- or overfitting, depending on the dataset (Fig. 2). We believe that it is generally preferable to continue training until an early-stopping criterion is met. However, applying the same criterion as with MLP¹¹ does not improve performance over the default of 100 ($p = 0.81$ without missing-indicators, $p = 0.85$ with missing-indicators) and does not change the relative advantage due to missing-indicators ($p = 0.20$).

4.2 When Do Missing-Indicators Start to Become Useful in Terms of Missingness?

The theoretical motivation for representing missing values through missing-indicators is that this allows classifiers to learn the information encoded in their distribution. In principle, this should be easier when there are more examples to learn from. We can use this principle to obtain a better understanding of when missing-indicators might be useful on a per-attribute level.

The challenge that we have to overcome is that we would like to study individual attributes, but classification performance is measured on the dataset level. We tackle this by studying datasets with only one attribute with missing values, allowing us to investigate the relation between the properties of the attribute and classification performance on the dataset.

We conduct the following experiment. For each attribute with missing values in each dataset, we reduce the original dataset by removing all other attributes with missing values. We thus obtain 1148 reduced datasets with only one attribute with missing values, onto which we apply each of our classifiers (with pruning for CART, 1000 trees for ERT and early-stopping for GBM) and consider whether missing-indicators increase or decrease AUROC (we dismiss ties). Finally, for each classifier we fit a logistic regression model with cluster robust covariance (clustered by the originating dataset), with the following potential parameters: categoricalness (whether the attribute is categorical) and either the number of missing values (log-transformed) or the missing rate. We use the Akaike information criterion [1] to decide whether to select these parameters.

We find that for most classifiers, either the absolute or the relative number of missing values is an informative parameter with positive coefficient. For MLP, neither parameter is informative, while for RF, the number of missing values is an informative parameter with negative coefficient, for which we have no explanation at present. For every classifier, categoricalness is an informative parameter with positive coefficient, meaning that missing-indicators are more beneficial for categorical than for numerical attributes.

¹¹ Setting aside 10% of the data for validation, stopping when validation loss has not decreased by at least 0.0001 for ten iterations, with a maximum of 10 000 iterations.

The fitted logistic regression models allow us to calculate attribute-specific thresholds above which missing-indicators are more likely than not to increase AUROC, for all classifiers except MLP and RF (Table 5). In many cases, these thresholds are 1 or 0.0, indicating that missing-indicators are always likely to increase AUROC.

Table 5: Thresholds above which missing-indicators are more likely than not to increase AUROC, in terms of the absolute number of missing values or the missing rate.

Classifier	Missing values		Missing rate	
	Cat	Num	Cat	Num
NN-1	1	302		
NN-2	2	130		
NN-1-D	1	291		
NN-2-D	1	73		
SVM-L			0.0	0.0
SVM-G			0.0	0.40
LR			0.0	0.0
CART			0.0	0.12
ERT			0.0	1.0
ABT	1	23200		
GBM			0.0	0.0

4.3 Does Using Mean Imputation Instead of Mode Imputation Allow for More Information to Be Learned from Missing Categorical Values?

As indicated above, missing-indicators are generally more likely to increase performance for categorical than for numerical attributes. A potential explanation for this is the fact that the mode of a categorical attribute is one of the non-missing values, whereas the mean of a numerical attribute is generally not equal to one of the non-missing values. Therefore, categorical imputation renders missing values truly indistinguishable from non-missing values, whereas numerical imputation does not — the information expressed by missing values may be partially recoverable, as argued by Le Morvan et al. [49] and discussed in the Introduction.

We can achieve a similar partial representation of missing categorical values by changing the order in which we perform imputation and one-hot encoding, i.e. by performing numerical imputation on one-hot encoded categorical attributes with missing values. For imputation without missing-indicators, this indeed leads

to better performance for some classifiers, while in combination with missing-indicators, it does not make much of a difference (Table 6)¹².

Table 6: One-sided p -values, mean imputation after one-hot encoding versus mode imputation of missing categorical values.

Classifier Without — With missing-indicators		
NN-1	0.020	0.077
NN-2	0.14	0.031
NN-1-D	0.016	0.12
NN-2-D	0.16	0.031
SVM-L	0.43	0.57
SVM-G	0.17	0.56
LR	0.81	0.057
MLP	0.16	0.60
CART	0.44	0.30
RF	0.046	0.57
ERT	0.030	0.95
ABT	0.48	0.62
GBM	0.077	0.54

4.4 How Do Missing-Indicators Compare to Model-Specific Approaches to Missing Values?

While not the primary focus of this paper, we may also wonder how the missing-indicator approach compares to model-specific approaches to missing values. For CART and RF, we consider the proposal by [80], that a decision tree should evaluate two variants of each split, with missing values sent to either side. This has been implemented in the latest version of scikit-learn (1.4.0), which was released after the previous experiments in this section were conducted. In addition, we have modified the implementation of the nearest neighbour classifier in scikit-learn to obtain the approach labelled as ‘normal’ in [19]. This calculates the distance between two records by linearly extrapolating the distance calculated only on the basis of all non-missing feature values. We note that every model-specific approach is different — we expect that their effect on classification performance will differ from case to case — so our evaluation of these two approaches only serves an illustrative purpose.

We find (Table 7, Test 1) that the model-specific approach for the nearest neighbour classifiers performs significantly worse than mean/mean imputation with missing-indicators. In contrast, there is no difference for CART, and the model-specific approach appears to perform better for RF. We can also ask whether these model-specific approaches benefit from adding missing-indicators

¹² LR is an exception here, we have no explanation for this.

— here this only appears to be the case for the nearest neighbour classifiers (Table 7, Test 2), i.e. when the model-specific approach performs badly. However, even with missing-indicators the model-specific approach for the nearest neighbour classifiers does not perform better than mean/mean imputation with missing indicators (Table 7, Test 3).

Table 7: One-sided p -values, model-specific missing value approaches. Test 1: Mean/mean imputation with missing-indicators vs model-specific approach without; Test 2: Model-specific approach with missing-indicators vs model-specific approach without; Test 3: Model-specific approach vs mean/mean imputation, both with missing-indicators.

Classifier	Test 1	Test 2	Test 3
NN-1	0.00036	0.00017	0.56
NN-2	0.00074	0.00015	0.94
NN-1-D	0.00020	0.00015	0.89
NN-2-D	0.00023	0.00011	0.93
CART	0.50	0.86	0.66
RF	0.92	0.60	0.092

5 Conclusion

We have presented the first large-scale experimental evaluation of the effect of the missing-indicator approach on classification performance, conducted on real datasets with naturally occurring missing values, paired with three different imputation techniques. The central question was whether, on balance, more benefit can be derived from the additional information encoded in a representation of missing values, or from the lower-dimensional projection of the data obtained by omitting missing-indicators.

On the whole, missing-indicators increase performance for the classification algorithms that we considered. An exception was CART, which suffers from overfitting in its default scikit-learn configuration. When pruning is applied, missing-indicators do increase performance. For ERT, the advantage of missing-indicators becomes more significant when underfitting is controlled.

We also found that, in the presence of missing-indicators, nearest neighbour and iterative imputation do not significantly increase performance over simple mean/mode imputation. This is a useful finding, because implementations of more sophisticated imputation strategies may not always be available to practitioners working in different frameworks, or easy to apply.

In a follow-up experiment, we determined attribute-specific missingness thresholds, above which missing-indicators are more likely than not to increase performance. For categorical attributes, this threshold is generally very low, while

for numerical attributes, there is more variation among classifiers, in particular as to whether this threshold is absolute or relative to the total number of records.

The greater usefulness of missing-indicators for categorical than for numerical attributes can be explained by the fact that the mean of a numerical attribute is not generally identical to any of the non-missing values, and that mean imputation therefore preserves some of the information of missing values. This is supported by the results of a further experiment, which showed that, in the absence of missing-indicators, applying mean imputation to one-hot encoded categorical attributes results in somewhat better performance than mode imputation.

While we have mainly considered the use of missing-indicators with imputation, there also exist model-specific solutions for missing values, that can in turn be combined with missing-indicators. Whether missing-indicators outperform these model-specific approaches has to be determined on a case-by-case basis. This was illustrated by our third follow-up experiment for nearest neighbour and decision tree classifiers.

We conclude that the combination of mean/mode imputation with missing-indicators is a safe default approach towards missing values in classification tasks. While over- or underfitting is a concern for certain classifiers, it is a concern for these classifiers with or without missing-indicators. However, practitioners may want to omit missing-indicators when the classification algorithm to be used has a special provision for missing values, when the missingness thresholds that we determined are not met, or on the basis of specific information about the distribution of missing values in the dataset. The use of missing-indicators can also be combined with dimensionality reduction algorithms to increase the information density of the resulting dataset.

The problem of missing data has been the subject of a rich body of theoretical literature. We hope to have contributed with this paper to the practical evaluation of some of that theory. In particular, we are happy to have identified twenty real-life datasets with missing values, and hope that in the future, more such datasets will be collected.

Data and code. Datasets and the code to reproduce our experiments are available at <https://cwi.ugent.be/~oulenz/code/lenz-2024-no.tar.gz>.

Acknowledgments. The research reported in this paper was conducted with the financial support of the Odysseus programme of the Research Foundation – Flanders (FWO). This publication is part of the project Digital Twin with project number P18-03 of the research programme TTW Perspective, which is (partly) financed by the Dutch Research Council (NWO). We would like to express our thanks to Geert van der Heijden for answering a question about [41].

Disclosure of Interests. The authors have no competing interests to declare that are relevant to the content of this article

References

1. Akaike, H.: Information theory and an extension of the maximum likelihood principle. In: *Proceedings of the 2nd International Symposium on Information Theory*. pp. 267–281. Akadémiai Kiadó (1971)
2. Allison, P.D.: *Missing Data*. Sage Publications, Thousand Oaks, California (2001)
3. Allison, P.D.: Missing data. In: Marsden, P.V., Wright, J.D. (eds.) *Handbook of Survey Research*, chap. 20, pp. 631–657. Emerald Group Publishing, Bingley, England, second edn. (2010)
4. Anderson, A.B., Basilevsky, A., Hum, D.P.J.: Missing data: A review of the literature. In: Rossi, P.H., Wright, J.D., Anderson, A.B. (eds.) *Handbook of Survey Research*, chap. 12, pp. 415–494. Quantitative Studies in Social Relations, Academic Press, New York (1983)
5. Aste, M., Boninsegna, M., Freno, A., Trentin, E.: Techniques for dealing with incomplete data: a tutorial and survey. *Pattern Analysis and Applications* **18**(1), 1–29 (2015)
6. Breiman, L.: Random forests. *Machine Learning* **45**(1), 5–32 (2001)
7. Breiman, L., Friedman, J.H., Olshen, R.A., Stone, C.J.: *Classification and Regression Trees*. The Wadsworth statistics/probability series, Wadsworth, Monterey, California (1984)
8. van Buuren, S., Groothuis-Oudshoorn, K.: mice: Multivariate imputation by chained equations in R. *Journal of Statistical Software* **45**(3), 1–67 (2011)
9. Candillier, L., Lemaire, V.: Design and analysis of the Nomao Challenge: Active learning in the real-world. In: *ECML-PKDD 2012: Active Learning in Real-world Applications Workshop* (2012)
10. Cestnik, B., Kononenko, I., Bratko, I.: ASSISTANT 86: A knowledge-elicitation tool for sophisticated users. In: *EWSL 87: Proceedings of the 2nd European Working Session on Learning*. pp. 31–45. Sigma Press (1987)
11. Chow, W.K.: A look at various estimators in logistic models in the presence of missing values. Tech. Rep. N-1324-HEW, Rand Corporation, Santa Monica, California (1979)
12. Cohen, J.: Multiple regression as a general data-analytic system. *Psychological Bulletin* **70**(6), 426–443 (1968)
13. Cohen, J., Cohen, P.: *Applied Multiple Regression/Correlation Analysis for the Behavioral Sciences*, chap. 7. Missing Data, pp. 265–290. Lawrence Erlbaum Associates, Hillsdale, New Jersey (1975)
14. Cortes, C., Vapnik, V.: Support-vector networks. *Machine Learning* **20**(3), 273–297 (1995)
15. Cox, D.R.: Some procedures connected with the logistic qualitative response curve. in (fn david, ed.) *research papers in statistics: Essays in honour of j. neyman's 70th birthday*. In: David, F.N. (ed.) *Research Papers in Statistics: Festschrift for J. Neyman*, pp. 55–71. John Wiley & Sons, London (1966)
16. Das, S., Datta, S., Chaudhuri, B.B.: Handling data irregularities in classification: Foundations, trends, and future challenges. *Pattern Recognition* **81**, 674–693 (2018)
17. Detrano, R., Janosi, A., Steinbrunn, W., Pfisterer, M., Schmid, J.J., Sandhu, S., Guppy, K.H., Lee, S., Froelicher, V.: International application of a new probability algorithm for the diagnosis of coronary artery disease. *The American Journal of Cardiology* **64**(5), 304–310 (1989)

18. Ding, Y., Simonoff, J.S.: An investigation of missing data methods for classification trees applied to binary response data. *Journal of Machine Learning Research* **11**(1), 131–170 (2010)
19. Dixon, J.K.: Pattern recognition with partly missing data. *IEEE Transactions on Systems, Man, and Cybernetics* **9**(10), 617–621 (1979)
20. Dua, D., Graff, C.: UCI machine learning repository (2019), <http://archive.ics.uci.edu/ml>
21. Dudani, S.A.: The distance-weighted k-nearest-neighbor rule. *IEEE Transactions on Systems, Man, and Cybernetics* **6**(4), 325–327 (1976)
22. Efron, B., Gong, G.: Statistical theory and the computer. In: *Computer Science and Statistics: Proceedings of the 13th Symposium on the Interface*. pp. 3–7. Springer (1981)
23. Eirola, E.: Machine learning methods for incomplete data and variable selection. Ph.D. thesis, Aalto University, Espoo (2014)
24. Elter, M., Schulz-Wendtland, R., Wittenberg, T.: The prediction of breast cancer biopsy outcomes using two CAD approaches that both emphasize an intelligible decision process. *Medical physics* **34**(11), 4164–4172 (2007)
25. Enders, C.K.: *Applied Missing Data Analysis*. Methodology in the Social Sciences, The Guilford Press, New York (2010)
26. Evans, B., Fisher, D.: Overcoming process delays with decision tree induction. *IEEE Expert* **9**(1), 60–66 (1994)
27. Ferreira Costa, C., Nascimento, M.A.: IDA 2016 industrial challenge: Using machine learning for predicting failures. In: *IDA 2016: Proceedings of the 15th International Symposium on Intelligent Data Analysis*. Lecture Notes in Computer Science, vol. 9897, pp. 381–386. Springer (2016)
28. Fix, E., Hodges, Jr, J.: Discriminatory analysis — nonparametric discrimination: Consistency properties. Tech. Rep. 21-49-004, USAF School of Aviation Medicine, Randolph Field, Texas (1951)
29. Freund, Y., Schapire, R.E.: A decision-theoretic generalization of on-line learning and an application to boosting. In: *EuroCOLT '95: Proceedings of the Second European Conference on Computational Learning Theory*. Lecture Notes in Computer Science, vol. 904, pp. 23–37. Springer (1995)
30. Friedman, J.H.: Greedy function approximation: A gradient boosting machine. *The Annals of Statistics* **29**(5), 1189–1232 (2001)
31. Fukushima, K.: Visual feature extraction by a multilayered network of analog threshold elements. *IEEE Transactions on Systems Science and Cybernetics* **5**(4), 322–333 (1969)
32. García, S., Luengo, J., Herrera, F.: *Data Preprocessing in Data Mining*, Intelligent Systems Reference Library, vol. 72, chap. 4. Dealing with Missing Values. Springer, Cham, Zug (2015)
33. Geurts, P., Ernst, D., Wehenkel, L.: Extremely randomized trees. *Machine learning* **63**(1), 3–42 (2006)
34. Glorot, X., Bengio, Y.: Understanding the difficulty of training deep feedforward neural networks. In: *AISTATS 2010: Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*. Proceedings of Machine Learning Research, vol. 9, pp. 249–256. JMLR Workshop and Conference Proceedings (2010)
35. Golovenkin, S.E., Bac, J., Chervov, A., Mirkes, E.M., Orlova, Y.V., Barillot, E., Gorban, A.N., Zinovyev, A.: Trajectories, bifurcations, and pseudo-time in large clinical datasets: applications to myocardial infarction and diabetes data. *GigaScience* **9**(11), giaa128 (2020)

36. Graham, J.W.: Missing data analysis: Making it work in the real world. *Annual Review of Psychology* **60**, 549–576 (2009)
37. Grzymala-Busse, J.W.: Knowledge acquisition under uncertainty—a rough set approach. *Journal of Intelligent and Robotic Systems* **1**(1), 3–16 (1988)
38. Grzymala-Busse, J.W., Hu, M.: A comparison of several approaches to missing attribute values in data mining. In: *RSCTC 2000: Proceedings of the Second International Conference on Rough Sets and Current Trends in Computing. Lecture Notes in Artificial Intelligence*, vol. 2005, pp. 378–385. Springer (2000)
39. Güvenir, H.A., Acar, B., Demiröz, G., Çekin, A.: A supervised machine learning algorithm for arrhythmia analysis. In: *Proceedings of the 24th Annual Meeting of Computers in Cardiology. Computers in Cardiology*, vol. 24, pp. 433–436. IEEE (1997)
40. Hand, D.J., Till, R.J.: A simple generalisation of the area under the ROC curve for multiple class classification problems. *Machine learning* **45**(2), 171–186 (2001)
41. van der Heijden, G.J.M.G., Donders, A.R.T., Stijnen, T., Moons, K.G.M.: Imputation of missing values is superior to complete case analysis and the missing-indicator method in multivariable diagnostic research: A clinical example. *Journal of Clinical Epidemiology* **59**(10), 1102–1109 (2006)
42. Hutcheson, Jr, J.D., Prather, J.E.: Interpreting the effects of missing data in survey research. *Southeastern Political Review* **9**(2), 129–143 (1981)
43. Ipsen, N., Mattei, P.A., Frellsen, J.: How to deal with missing data in supervised deep learning? In: *Artemiss 2020: First ICML Workshop on the Art of Learning with Missing Values* (2020)
44. Jones, M.P.: Indicator and stratification methods for missing explanatory variables in multiple linear regression. *Journal of the American statistical association* **91**(433), 222–230 (1996)
45. Josse, J., Chen, J.M., Prost, N., Varoquaux, G., Scornet, E.: On the consistency of supervised learning with missing values. *Statistical Papers* (2024)
46. Kim, J.O., Curry, J.: The treatment of missing data in multivariate analysis. *Sociological Methods & Research* **6**(2), 215–240 (1977)
47. Kingma, D.P., Ba, J.L.: Adam: A method for stochastic optimization. In: *ICLR 2015: 3rd International Conference on Learning Representations* (2015)
48. Kohavi, R.: Scaling up the accuracy of naive-bayes classifiers: a decision-tree hybrid. In: *KDD-96: Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*, pp. 202–207. AAAI Press (1996)
49. Le Morvan, M., Josse, J., Scornet, E., Varoquaux, G.: What’s a good imputation to predict with missing values? In: *NeurIPS 2021: Proceedings of the Thirty-fifth Conference on Neural Information Processing Systems. Advances in neural information processing systems*, vol. 34, pp. 11530–11540. NIPS Foundation (2021)
50. Luengo, J., García, S., Herrera, F.: On the choice of the best imputation methods for missing values considering three groups of classification methods. *Knowledge and Information Systems* **32**(1), 77–108 (2012)
51. Luengo, J., Sáez, J.A., Herrera, F.: Missing data imputation for fuzzy rule-based classification systems. *Soft Computing* **16**(5), 863–881 (2012)
52. Marlin, B.M.: Missing Data Problems in Machine Learning. Ph.D. thesis, University of Toronto (2008)
53. McCann, M., Li, Y., Maguire, L., Johnston, A.: Causality challenge: benchmarking relevant signal components for effective monitoring and process control. In: *NIPS 2008: Proceedings of Workshop on Causality. Proceedings of Machine Learning Research*, vol. 6, pp. 277–288. JMLR Workshop and Conference Proceedings (2008)

54. McLeish, M., Cecile, M.: Enhancing medical expert systems with knowledge obtained from statistical data. *Annals of Mathematics and Artificial Intelligence* **2**(1–4), 261–276 (1990)
55. Michalski, R.S., Chilausky, R.L.: Learning by being told and learning from examples: An experimental comparison of the two methods of knowledge acquisition in the context of developing an expert system for soybean disease diagnosis. *International Journal of Policy Analysis and Information Systems* **4**(2), 125–161 (1980)
56. Ng, C.G., Yusoff, M.S.B.: Missing values in data analysis: Ignore or impute? *Education in Medicine Journal* **3**(1) (2011)
57. Orme, J.G., Reis, J.: Multiple regression with missing data. *Journal of Social Service Research* **15**(1–2), 61–91 (1991)
58. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, É.: Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* **12**(85), 2825–2830 (2011)
59. Pereira Barata, A., Takes, F.W., van den Herik, H.J., Veenman, C.J.: Imputation methods outperform missing-indicator for data missing completely at random. In: *ICDM 2019: Proceedings of the Workshops*. pp. 407–414. IEEE (2019)
60. Perez-Lebel, A., Varoquaux, G., Le Morvan, M., Josse, J., Poline, J.B.: Benchmarking missing-values approaches for predictive models on health databases. *GigaScience* **11**(1), giac013 (2022)
61. Pigott, T.D.: A review of methods for missing data. *Educational Research and Evaluation* **7**(4), 353–383 (2001)
62. Quinlan, J.R.: Simplifying decision trees. *International Journal of Man-Machine Studies* **27**(3), 221–234 (1987)
63. Quinlan, J.R.: Unknown attribute values in induction. In: *Proceedings of the Sixth International Workshop on Machine Learning*. pp. 164–168. Morgan Kaufmann (1989)
64. Quinlan, J.R., Compton, P.J., Horn, K.A., Lazarus, L.: Inductive knowledge acquisition: a case study. In: *Proceedings of the Second Australian Conference on Applications of Expert Systems*. pp. 157–173. Turing Institute Press (1986)
65. Rosenblatt, F.: Principles of neurodynamics — perceptrons and the theory of brain mechanisms. Tech. Rep. VG-1196-G-8, Cornell Aeronautical Laboratory, Buffalo, New York (1961)
66. Rubin, D.B.: Inference and missing data. *Biometrika* **63**(3), 581–592 (1976)
67. Rubini, L.J., Eswaran, P.: Generating comparative analysis of early stage prediction of chronic kidney disease. *International Journal of Modern Engineering Research* **5**(7), 49–55 (2015)
68. Santos, M.S., Abreu, P.H., García-Laencina, P.J., Simão, A., Carvalho, A.: A new cluster-based oversampling method for improving survival prediction of hepatocellular carcinoma patients. *Journal of biomedical informatics* **58**, 49–59 (2015)
69. Schafer, J.L.: *Analysis of Incomplete Multivariate Data*, Monographs on Statistics and Applied Probability, vol. 72. Chapman & Hall, London (1997)
70. Schafer, J.L., Graham, J.W.: Missing data: Our view of the state of the art. *Psychological Methods* **7**(2), 147–177 (2002)
71. Schlimmer, J.C.: *Concept Acquisition Through Representational Adjustment*. Ph.D. thesis, University of California, Irvine (1987)
72. Śmieja, M., Struski, Ł., Tabor, J., Marzec, M.: Generalized RBF kernel for incomplete data. *Knowledge-Based Systems* **173**, 150–162 (2019)

73. Śmieja, M., Struski, Ł., Tabor, J., Zieliński, B., Spurek, P.: Processing of missing data by neural networks. In: *NeurIPS 2018: Proceedings of the Thirty-second Annual Conference on Neural Information Processing Systems*. Advances in neural information processing systems, vol. 31, pp. 689–696. NIPS Foundation (2018)
74. Soltani Zarrin, P., Röckendorf, N., Wenger, C.: In-vitro classification of saliva samples of COPD patients and healthy controls using machine learning tools. *IEEE Access* **8**, 168053–168060 (2020)
75. Sperrin, M., Martin, G.P., Sisk, R., Peek, N.: Missing data should be handled differently for prediction than for description or causal explanation. *Journal of Clinical Epidemiology* **125**, 183–187 (2020)
76. Stumpf, S.A.: A note on handling missing data. *Journal of Management* **4**(1), 65–73 (1978)
77. Tipping, M.E.: Sparse Bayesian learning and the relevance vector machine. *Journal of Machine Learning Research* **1**, 211–244 (2001)
78. Tresp, V., Neuneier, R., Ahmad, S.: Efficient methods for dealing with missing data in supervised learning. In: *NIPS-94: Proceedings of the Eighth Annual Conference on Neural Information Processing Systems*. Advances in neural information processing systems, vol. 7, pp. 689–696. MIT Press (1994)
79. Troyanskaya, O., Cantor, M., Sherlock, G., Brown, P., Hastie, T., Tibshirani, R., Botstein, D., Altman, R.B.: Missing value estimation methods for DNA microarrays. *Bioinformatics* **17**(6), 520–525 (2001)
80. Twala, B.E., Jones, M., Hand, D.J.: Good methods for coping with missing data in decision trees. *Pattern Recognition Letters* **29**(7), 950–956 (2008)
81. Vamplew, P., Adams, A.: Missing values in a backpropagation neural net. In: *ACNN '92: Proceedings of the Third Australian Conference on Neural Networks*. pp. 64–66. Sydney University Electrical Engineering (1992)
82. Wilcoxon, F.: Individual comparisons by ranking methods. *Biometrics Bulletin* **1**(6), 80–83 (1945)
83. Zhu, J., Zou, H., Rosset, S., Hastie, T.: Multi-class AdaBoost. *Statistics and Its Interface* **2**(3), 349–360 (2009)