

The Effects of Missing Data Characteristics on the Choice of Imputation Techniques

Oyekale Abel Alade^{*,†,||}, Ali Selamat^{*,‡,§,¶,***} and Roselina Sallehuddin^{**,††}

**School of Computing, Faculty of Engineering
Universiti Teknologi Malaysia
81310 Johor Bahru, Johor, Malaysia*

*†Department of Computer Science
Federal Polytechnic Bida, Nigeria*

*‡Media & Games Center of Excellence (MAGICX)
Faculty of Engineering, Universiti Teknologi Malaysia
81310 Johor Bahru, Johor, Malaysia*

*§Malaysia-Japan International Institute of Technology (MJIIT)
Universiti Teknologi Malaysia
Jalan Sultan Yahya Petra, 54100
Kuala Lumpur, Malaysia*

*¶University of Hradec Králové
Rokitanskeho 62, 500 03 Hradec Králové, Czech Republic
||kaleabel@gmail.com
***aselamat@utm.my
††roselina@utm.my*

Received 15 October 2018

Accepted 17 January 2020

Published 20 March 2020

One major characteristic of data is completeness. Missing data is a significant problem in medical datasets. It leads to incorrect classification of patients and is dangerous to the health management of patients. Many factors lead to the missingness of values in databases in medical datasets. In this paper, we propose the need to examine the causes of missing data in a medical dataset to ensure that the right imputation method is used in solving the problem. The mechanism of missingness in datasets was studied to know the missing pattern of datasets and determine a suitable imputation technique to generate complete datasets. The pattern shows that the missingness of the dataset used in this study is not a monotone missing pattern. Also, single imputation techniques underestimate variance and ignore relationships among the variables; therefore, we used multiple imputations technique that runs in five iterations for the imputation of each missing value. The whole missing values in the dataset were 100% regenerated. The imputed datasets were validated using an extreme learning machine (ELM) classifier. The results show improvement in the accuracy of the imputed datasets. The work

****Corresponding author.**

This is an Open Access article published by World Scientific Publishing Company. It is distributed under the terms of the Creative Commons Attribution 4.0 (CC BY) License which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

can, however, be extended to compare the accuracy of the imputed datasets with the original dataset with different classifiers like support vector machine (SVM), radial basis function (RBF), and ELMs.

Keywords: Imputation techniques; mechanism of missingness; missing data; missing pattern; multiple imputations.

1. Introduction

Missing data/values describe the absence of important data items in instances of datasets. Data are collected at various points for medical investigation. Lichman¹ observed two possible types of databases in the medical domain. The first type is basically for hospital information systems, which consists of a vast number of attributes. Many of such attributes are not directly required for the diagnosis of the ailments in the patients. The other type of medical database is collected by experts. The databases may contain unique research data on topics based on hypothesis propositions that must be investigated. Missing data is very much pervasive in either of these types of databases, as in many other databases, and most real-world data analysis tasks. García-Laencina *et al.*² and Tran *et al.*³ observed that 45% of datasets in UCI machine learning repository are marred by missing values (MVs), most of which may fall in the category of medical data. The pervasive nature of missing data has been described as one of the most challenging tasks in data science.⁴ The missing observations usually have the potential to be captured but were not captured due to some reasons that may arise from either the patient or the medical personnel. In Ref. 5, two basic patterns of missing data were considered. The patterns are: (a) missing features, a situation when the features exist but values were not taken; therefore the information is lost, or the cost of acquiring the feature is high; and (b) missing label, when the label is inherently missing — that is, a problem that cannot be avoided — which eventually affects the performance of classifiers in the face of ever-growing datasets⁶ in this information age. Missing data affects some operations in medical research. It makes it challenging to extract useful information from datasets; also, feature selection is not always applied to datasets with missing values.³

Missing data is a complex pattern problem that is inherent in equipment malfunction during data extractions, sampling, transcription, transmission, and noise during data preprocessing.⁵ Newgard and Lewis⁷ observed that missingness of data in clinical research could be a result of variables that are complex, time-sensitive, resource-intensive, or the method of collection of longitudinal data. There are other causes of missing values in datasets observed in Ref. 6. They are (i) possibility of the observations being irrelevant, especially in medical data collection, because the primary reason for data collection is for medical investigation and diagnosis, and not really for research purposes. (ii) Inability to record the values when they were collected due to the emergency situation; or patient response avoidance from the respondent for privacy. (iii) The omission of essential features during the data collection plan. (iv) Noncapturing of seemingly available values. Two significant

problems envisaged in Ref. 8 with the presence of a missing value in datasets are (i) reduction in overall statistical power and (ii) statistical bias estimation.

The problem of missing data is being debated for some time.⁹ The best method to fix missing values is to revisit the data collection/extraction process to recollect and correct missing values and noisy values, respectively. This may involve the reinvestigation of patients from various examination units to fix the missing values, but this method of recollection/re-extraction might not be practicable. Therefore, there is a need for fitting techniques with close approximations to the real data. The resolve, however, is that the imputation of missing data has no definite solution.¹⁰

The accuracy of classification from the incomplete dataset is unreliable because some vital information relevant to the analysis might be lost. Besides, some classifiers find it challenging to run in the face of missing values.¹¹ Therefore, there is a necessity for missing value imputation. The focus of imputation is to estimate the possible value of the missing data from observed data and fix the estimated values as a replacement of the lost values; that is, to ensure complete datasets.

Several approaches have been used in research for the treatments of missing data in datasets. Some authors used the percentage of missingness in a dataset as a measure for the choice of imputation technique.⁹ Some implementations treat missing data implicitly. This brings about different results when such treatments are replicated using different applications. Although the difference may not be significant, however, these approaches compromise the scientific soundness of the studies. An explicit approach to handling missing data is a better practice. In Ref. 8, it was observed that the choice of missing data handling imputation technique depends on the research focus, whether it is a pragmatic or an analytical approach. Although missing data is pervasive in data science generally, in this paper, we focus the study on the effect of missing values on the medical datasets as an offshoot to application areas.

From this point on, the paper is organized as follows: Section 2 reviews relevant works on medical datasets with missing data and some techniques of missing data imputation; Sec. 3 describes the characteristics of missing values, i.e. the mechanisms of missingness; various treatments of missing values in datasets are discussed in Sec. 4; Sec. 5 explains the proposed multiple imputations model; Sec. 6 reports the experimental setup; discussion of the results of imputation on Pima Indians Diabetes dataset is in Sec. 7; while Sec. 8 concludes the paper.

2. Review of Literature

Medical datasets have been classified by many researchers. Some of the datasets are complete, while some are incomplete. Several studies had been carried out on missing data and imputation techniques from different perspectives. Some authors explicitly treat missing values using different imputation techniques, while some are passive about it, leading to the assignment of zeros (0), deletions of cases/features, or completely ignoring the missing values. Extreme learning machine (ELM) is widely

used in recent time to solving classification, clustering, compression, forecasting, and regression problems¹² because it tolerates quite a good number of feature mapping functions such as sigmoid, hard-limit, Gaussian, multi-quadratic, wavelet, Fourier series, etc., and it handles large and small datasets efficiently.¹³ Subbulakshmi and Deepa¹⁴ proposed a machine learning paradigm by integrating particle swarm optimization (PSO) technique with extreme learning machines to classify some medical datasets. The hybrid system performed well compared to other classifiers; however, missing values in the datasets were substituted with zeros. This approach is scientifically unfit for accurate results. Zeros do not represent a good imputation of missing values.

Bai *et al.*¹⁵ overviewed the hidden challenges of missing values in medical datasets during preprocessing. They proposed the imputation of the missing values in the medical datasets with categorical attributes, the causes and pattern of missingness in the datasets were, however, not considered. This may result in the wrong choice of imputation technique.

An extensive review was carried out in Ref. 16 on missing value imputations. The authors provide a detailed analysis of various imputation techniques. They grouped imputation techniques into four broad categories: global, local, hybrid, and knowledge-assisted approaches, but there was no experiment conducted to prove any of the imputation techniques discussed in their study. The combination of Gaussian mixture model and extreme learning machine (GMM-ELM) was proposed as a reliable approximation technique for imputing missing data by Sovilj *et al.* in Ref. 17. The results of their work improved the imputation of missing values over the mean imputation technique; however, the execution time was longer. Bai *et al.*¹⁵ addressed categorical attribute missing values in medical datasets using imputation measure. The work, however, used a hypothetical dataset with only nine cases and only two missing values; the characteristics of the missingness was not considered in work. In Ref. 18, Tsai and Chang investigated the effects of filtering outliers from datasets on imputation tasks using instance selection on categorical, numerical, and mixed-type attributes. The effectiveness of the method was tested with *k*-NN and support vector machine (SVM) classifiers. To compare the performances of three Bayesian imputation techniques, Austin and Escobar¹⁹ placed prior distribution on attributes with missing values. Monte Carlo simulation model was used to examine the performance of the sibling models. The result showed that mean and mean square error of logistic and Bayesian models depend on risk factors examined, and the mechanism of missing data that had been used. This gives an insight into the necessity for consideration of the mechanism of missingness for the right choice of imputation technique. Multiple imputations technique with Pohar Perme method was used in Ref. 20 to estimate the net survival for stage-specific colorectal cancer. They concluded that the interpretation of datasets with a high percentage of missing values should be cautious and should be with sensitivity analysis. However, the characteristics of the missingness of data in the dataset were not taken into consideration before the choice of the imputation technique.

In Ref. 21, a hybrid imputation method based on the integration of fuzzy *c*-means (FCM) and the genetic algorithm (GA) for missing traffic volume data was developed. The study based the estimation on inductance loop detector outputs. The result, under prevailing traffic conditions, performed better than conventional methods. All these methods and much more in literature underscore the need for imputation of missing data in a given dataset with missing values.

Although most of the imputation techniques mentioned above attempt to fill the missing values by approximately conforming to the distribution of the datasets, however, the methods of the imputation of values are not explicitly modeled; therefore, further analysis is ignored¹⁷ and thereby lead to bias result. Nguyen *et al.*²² raised some critical points to consider when constructing an imputation model. These are (a) model imputation functional form, (b) feature selection for the model, (c) inclusion of nonlinear relationships in the model, and (d) the best way to handle nonnormal continuous features. They concluded that there is no consensus in literature on how to implement these decisions, these could be evaluated from the nature of the missing values in the datasets. Therefore, there is a need to know the nature of missingness in a dataset for the right choice of imputation technique. In the next section, we attempt to have an overview of the possible nature of missingness in datasets.

3. Mechanisms of Missingness

The focus of this section is to look at the characteristics of missing values in datasets. These characteristics determine the causes of the missingness in the dataset. It is good to know the cause(s) of missing values in a dataset to handle the missingness appropriately.²³ Some literature refers to this as *mechanisms of missingness*. Various mechanisms of missing data values are abound in literature,^{8,20,24–26} but the most popular ones are basically three which shall be considered for the purpose of this study.

3.1. Missing at random

This is a type of missingness that does not occur entirely at random; instead, they occur where there are other variables with complete information that can account for the missingness. It does not necessarily mean that the cases are similar to the complete counterpart. Missing at Random (MAR) is more realistic than missing completely at random (MCAR), and it is mostly applied to missing data imputation in many pieces of literature.⁷ It is based on an ignorable assumption: that is, the available information is sufficient, and the assignment mechanism can be ignored. This case arises when some respondents decide to hide some information that is personal or are unpopular about themselves.²⁷ Logistic regression with the outcome of 1 for the observed and 0 for the missing values is a reasonable option for its treatment. It can be statistically expressed as in (1), thus, for a random attribute X

and a predictor attribute Z , if

$$P(X|x_{\text{miss}}) = P(X|x_{\text{obs}}, Z), \quad (1)$$

then the distribution x is not affected by the values

$$X \in Z.$$

That is, when the missingness is based on the observed factors, then it is independent of the unobserved factors.

Although MAR is popularly accepted in many techniques, the result is still biased or imprecise results are yielded with simple imputation techniques.⁷

3.2. Missing completely at random

MACR occurs if the cause of missing values of observable features and the parameters of unobservable features of interest are independent, and their occurrence is entirely at random. Analyses performed on MCAR datasets are unbiased, although this type of datasets are rare. It is the highest level of randomness. It is expressed as follows:

$$P(X|x_{\text{miss}}) = P(X|x_{\text{obs}}). \quad (2)$$

Any imputation technique can be applied.¹⁸

3.3. Missing not at random

Missing Not at Random (MNAR) is a type of missing data where there is a relationship between the missing data and the reason for the missingness. It occurs when the missingness depends on the probability of the actual value of the missing data⁶ and some/all other observed data.^{8,19} Tsai and Chang¹⁸ observed that this mechanism would be difficult to judge because the missing data are unknown.

The treatments of missing data should be based on the mechanism of missing data, as explained in the next section.

4. Treatments of Missingness

In the treatment of missing data, two broad approaches are conventional in literature. These are (a) omission of missing data and (b) imputing the missing data.¹⁷ Some approaches to the treatment of missing values are outlined in Fig. 1 and later discussed in the following.

4.1. Omission

This approach simply deletes instances with missing data. The approach is common in some regression models, usually referred to as *listwise-deletion (complete case analysis)*.²⁸ It is only valid under the following conditions: (a) the instances with missing values in the sample are negligible; (b) the pattern of the missing data is

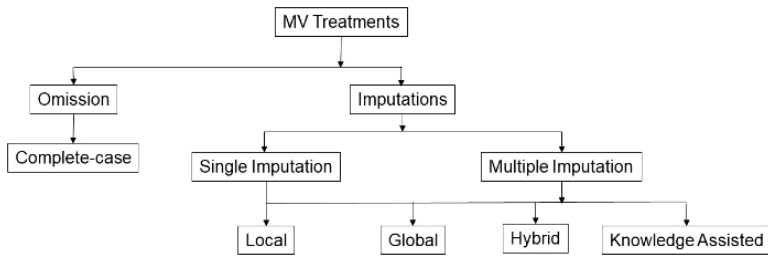


Fig. 1. Treatment of missing values.

MAR or MCAR.¹¹ This approach reduces the sample size, so they often lost vital information because deleted instances may be the essential and deciding factors for predictions and classifications.¹⁷ It limits the study power. Therefore, the imputation technique is preferred to listwise-deletion.

4.2. Imputation techniques

Imputation can be categorized into two: (a) single imputation and (b) multiple imputations.

4.2.1. Single imputation technique

Single imputation uses mean, median, mode, or conditional mean like a predicted value from a regression function evaluation or decision tree^{3,27} to generate the data to be imputed only once. Mean imputation is a standard method used to replace missing values.³ The average or median of the observed feature values is computed and substituted for each missing value for numerical attributes; and mode for simple ones. Unfortunately, this method does not present an actual distribution of features. It underestimates the variance and ignores relationships among the variables in the dataset.^{19,20} These lead to complications of statistical inferences. Imputation of missing data is handled sequentially — one by one. This method work with datasets with a limited number of variables.²⁷

The last observation carried forward: this method replaces every missing value with the last observation. This approach assumes that the result will not change after the last reading.⁸ It is a simple method, so it is accessible. It maintains the actual size of data; however, it might result in bias outcomes. This method is not analytical enough, so there is a need for a more comprehensive imputation technique.

4.2.2. Multiple imputations technique

This method is scientifically plausible to replace values with the modeled results of several imputations that analytically represent the missing data. For example, the regression model will reflect the uncertainty of regression coefficients and the sample variables in the model. Multiple imputations technique analytically creates several values to replace the missing data.²⁹ Different models also predict these replacement

values. It must be known that the aim of multiple imputations technique is not to produce the actual missing value; rather, it attempts to generate scientifically valid results to account for the missing values.³⁰ According to Ref. 31, it is possible that the simulated values may not fall within the expected range. The single idea about multiple imputations, however, is to form N complete datasets from the observed value analytically. N is the number of imputations carried out on the original dataset with missing values, and it produces N different complete datasets.

Armina *et al.*,¹⁶ further detail on local, global, hybrid, and knowledge-assisted imputation techniques, as sketched out in Fig. 1. In the next section, we propose a regression model for multiple imputations technique used in this study.

5. Proposed Multiple Imputations Model

In this section, a regression model is proposed for multiple imputations of missing data, because it draws values randomly from donor instances to predict values that are close to the missing values to be predicted; also regression creates inter-data variability using stochastic elements.³² The results of this data pooling produce correct standard error estimates.

For a dataset with missing data problem of $Y = X_{m+1}$ on m variables, X_1, \dots, X_m are the instances of a random sample when x are incomplete; the incomplete values can be estimated with the regression model shown in (3):

$$E(Y|X_1, \dots, X_m) = \beta_0 + \sum_{m=1}^M \beta_m X_m, \quad (3)$$

where $X = (X_1, \dots, X_m)$ is the correlation coefficient, $\beta = (\beta_1, \dots, \beta_m)$ is the relative size of the coefficient of regression with respect to one another, and β_0 is the intercept.

To make an inference of the regression coefficient, β_m of the missing values and the standard errors e_1, \dots, e_M are obtained for each dataset in M . The mean dataset estimate ($\hat{\beta}$) is given as follows:

$$\hat{\beta} = \frac{1}{m} \sum_{m=1}^M \beta_m. \quad (4)$$

The variation within (V) and between (W) the imputations is shown as follows:

$$V_\beta = W + \left(1 + \frac{1}{m}\right)E, \quad (5)$$

where

$$W = \frac{1}{m} \sum_{m=1}^M e^2, \quad E = \frac{1}{m-1} \sum_{m=1}^M (\hat{\beta}_m - \hat{\beta}).$$

6. Experimental Setup

In order to determine the choice of imputation techniques for a dataset, the Pima Indians Diabetes dataset from UCI was used. Pima Indians Diabetes dataset contains clinical tests and diagnoses of Pima Indian women of 21 years of age and above with diabetes.^{33,34} The dataset is made up of integer and real number data types. It has 768 instances — eight predictive features and a class feature. The features in the dataset are number of times pregnant (V1), a 2-h oral tolerant test for plasma glucose concentration (V2), diastolic blood pressure (V3), triceps skinfold thickness (V4), a 2-h serum insulin (V5), body-mass index (V6), diabetes pedigree (V7), and age (V8). The class (V9) is a binary classification dataset with 1 for positive and 0 for negative. The dataset was indicated to have missing values on the webpage,³³ but the real dataset did not show such signs. However, a critical examination of the dataset shows that variables like plasma glucose concentration, body-mass index, triceps of skinfold thickness, diastolic blood pressure, and 2-h serum insulin cannot be zero for any instance; therefore, it was assumed that all the data values that scored zeros are really missing values and were so treated in this study.

The missing values in the dataset were coded for proper identification by the imputation program. For an analytical presentation of the degree of missingness in the dataset, graphical and numerical summaries²² were used in our report. The percentage of missingness was calculated by instances, features, and values, as shown in Fig. 2.

The percentage of missingness among the instances was considered in this study as common in the literature.^{18,19} The missing values were analyzed in order to know the distributions of missingness among the various features that have missing values in the dataset. The distribution is shown in Table 1.

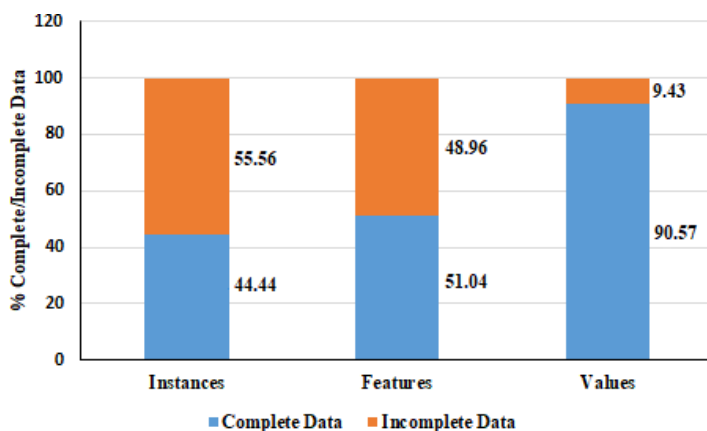


Fig. 2. Percentages of complete/incomplete data in instances, features, and values of Pima Indians Diabetes dataset.

Table 1. Distribution of missing values among the features with incomplete data in the instances of the Pima Indians Diabetes dataset.

	Missing		Valid <i>N</i>	Mean	Std. dev.
	<i>N</i>	Percentage			
V5	74	48.7%	394	155.55	118.776
V4	27	29.6%	541	29.15	10.477
V3	5	4.6%	733	72.41	12.382
V6	1	1.4%	757	32.457	6.9250
V2		0.7%	763	121.69	30.536

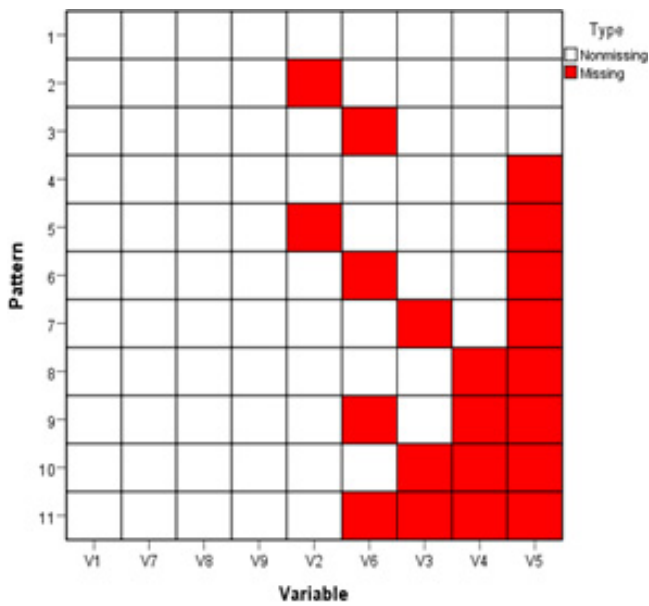


Fig. 3. Pattern of missing values in Pima Indians Diabetes dataset.

A missing pattern describes the set of features in a dataset that has at least an instance in the dataset with missing value(s) for the same feature(s) in the pattern.³

The pattern of missingness of the features in our dataset was analyzed. The result of the analysis is shown in Fig. 3. Each pattern in the figure corresponds to the group of instances with the same patterns of incomplete and complete data.

The variables are sorted in increasing order of missing values (from the least-missing-value feature to the one with the highest missing value) from left to right. This is to enable us to know the type of imputation required to fill up the missingness. Multiple imputations technique is carried out on the original dataset to be able to come up with five sets of complete datasets. The analysis of the imputed datasets for features with missing values is shown in Tables 2–6.

Table 2. Multiple imputations for V2.

Data	Imputation	<i>N</i>	Mean	Std. dev.	Min.	Max.
Original Data		763	121.69	30.536	44.00	199.00
Imputed Values	1	5	125.52	38.578	88.84	170.01
	2	5	119.22	10.672	111.21	137.24
	3	5	145.21	35.676	104.01	191.63
	4	5	111.99	27.614	76.61	140.40
	5	5	135.09	27.640	95.38	172.70
Complete Data After Imputation	1	768	121.71	30.565	44.00	199.00
	2	768	121.67	30.446	44.00	199.00
	3	768	121.84	30.603	44.00	199.00
	4	768	121.62	30.511	44.00	199.00
	5	768	121.77	30.520	44.00	199.00

Table 3. Multiple imputations for V3.

Data	Imputation	<i>N</i>	Mean	Std. dev.	Min.	Max.
Original Data		733	72.41	12.382	24.00	122.00
Imputed Values	1	35	72.55	14.902	43.49	100.49
	2	35	72.25	13.165	44.49	97.34
	3	35	74.56	12.447	43.87	101.37
	4	35	70.85	11.039	45.90	95.96
	5	35	69.06	11.674	50.13	96.20
Complete Data After Imputation	1	768	72.41	12.497	24.00	122.00
	2	768	72.40	12.410	24.00	122.00
	3	768	72.50	12.385	24.00	122.00
	4	768	72.33	12.322	24.00	122.00
	5	768	72.25	12.363	24.00	122.00

Table 4. Multiple imputations for V4.

Data	Imputation	<i>N</i>	Mean	Std. dev.	Min.	Max.
Original Data		541	29.15	10.477	7.00	99.00
Imputed Values	1	227	27.84	10.680	−3.03	63.88
	2	227	28.47	10.895	−1.92	55.99
	3	227	27.45	10.373	−0.62	56.51
	4	227	27.85	10.190	−6.94	58.87
	5	227	27.74	11.160	0.75	57.31
Complete Data After Imputation	1	768	28.77	10.547	−3.03	99.00
	2	768	28.95	10.600	−1.92	99.00
	3	768	28.65	10.469	−0.62	99.00
	4	768	28.77	10.403	−6.94	99.00
	5	768	28.74	10.696	0.75	99.00

Table 5. Multiple imputations for V5.

Data	Imputation	N	Mean	Std. dev.	Min.	Max.
Original Data		394	155.55	118.776	14.00	846.00
Imputed Values	1	374	146.07	111.231	−113.11	489.68
	2	374	153.85	122.138	−227.40	488.04
	3	374	154.91	120.382	−257.32	488.35
	4	374	144.23	126.086	−182.66	540.58
	5	374	134.68	119.561	−180.91	521.16
Complete Data After Imputation	1	768	150.93	115.186	−113.11	846.00
	2	768	154.72	120.349	−227.40	846.00
	3	768	155.24	119.483	−257.32	846.00
	4	768	150.03	122.441	−182.66	846.00
	5	768	145.39	119.538	−180.91	846.00

Table 6. Multiple imputations for V6.

Data	Imputation	N	Mean	Std. dev.	Min.	Max.
Original Data		757	32.457	6.9250	18.200	67.100
Imputed Values	1	11	32.271	7.4383	18.047	46.451
	2	11	34.026	7.7246	22.682	45.062
	3	11	29.469	7.2017	18.687	43.716
	4	11	32.797	6.2029	22.493	42.757
	5	11	28.430	6.8445	13.355	38.423
Complete Data After Imputation	1	768	32.455	6.9274	18.047	67.100
	2	768	32.480	6.9340	18.200	67.100
	3	768	32.415	6.9333	18.200	67.100
	4	768	32.462	6.9117	18.200	67.100
	5	768	32.400	6.9360	13.355	67.100

7. Discussion of Results

As mentioned earlier, the original Pima Indians Diabetes dataset could not be easily noticed as having missing data because every cell in the dataset is completely scored. However, the source (UCI database) categorically stated that the dataset has missing values.

Upon a critical look, it was observed that the missing values in the dataset were scored zero, and it was so treated except the first feature (number of times pregnant — V1), which is assumed can be zero among the selected women.

Figure 2 shows that 44.44% of the features in the dataset are complete. That is, four out of nine variables (class label inclusive) have complete data scored while 55.56% have incomplete (missing) values. Based on the instances (cases) in the dataset, 51.04%, which is 392 cases, have complete data, while 48.96% (366 cases) have incomplete data. For the entire values in the Pima Indians Diabetes dataset

(that is the intersections of the feature and instances values), 90.57% of the cells have values while 9.43% of them are missing.

In Table 1, the distribution of missing values among the cases is presented. The table shows five features with missing values in terms of their quantities and various percentages across the dataset. The features are arranged in descending order of the percentage constituents of their missingness: 2-h serum insulin (V5) has the highest percentage of missing values (48.7%), triceps skinfold thickness (V4) has 29.6%, diastolic blood pressure (V3) has 4.6%, while body-mass index (V6) and plasma glucose concentration (V2) have 1.4% and 0.7%, respectively. Number of times pregnant (V1), diabetes pedigree functions (V7), and age (V8) have no missing values. Also, the numbers of valid values are shown along with their means and standard deviations for all the features with missing values.

Figure 3 depicts the missing patterns of the features. The features are arranged from left to right with features with nonmissing values on the left, through those with the least missing values, to the ones with the highest missing values on the right. The missing pattern chart displays the value pattern for the analysis function. The pattern represents the group of instances that have the same patterns of incomplete and complete data. In Fig. 3, Pattern 1 corresponds to instances with no missing value; Pattern 2 shows instances that have missing values in V2; Pattern 3 represents instances with missing values in V6. Pattern 4 represents instances with missing values in V5; Pattern 5 shows instances with missing values in V2 and V5; Pattern 6 represents the missing values in V6 and V5; Pattern 7 is for missing values in V3 and V5; Pattern 8 represents instances with missing values in V4 and V5; Pattern 9 represents instances with missing values in V6, V4, and V5; Pattern 10 is for missing values in V3, V4, and V5; Pattern 11 is for missing values in V6, V3, V4, and V5. All the patterns show no missing values in V1, V7, V8, and V9. Although the dataset has the potential for 2^9 patterns,¹ only 11 feasible patterns are represented in 768 instances.

The features and patterns are arranged in an orderly manner to reveal the existence of monotonicity in the dataset. From the result of the patterns in Fig. 3, it is clear that the missingness in the dataset is nonmonotone because all missing cells and nonmissing cells are not contiguous; that is, the dataset is MAR as explained in Sec. 3. There are many values to be imputed to achieve monotonicity. Therefore, the use of a monotone (single) method of imputation may not be plausible; the use of multiple imputations technique (in Sec. 4.2.2) becomes the needed option.

Multiple imputations technique was utilized on the dataset features with missing data using (4). The five imputed features were V2, V3, V4, V5, and V6. The order of imputations of the features was V2, V6, V3, V4, and V5 (in the increasing order of the percentage of missingness). The imputation was complete for all missing values in each of the features, and there was no one that was omitted either as a result of “too” many missing values or no missing value. The descriptive analysis of the imputation on each feature is shown in Tables 2–6.

Table 2 shows multiple imputations for the feature V2. Five out of 768 instances in the feature are missing. The missing values were imputed in five iterations, which is a total of 25 runs. The missing values were 100% imputed. The mean, standard deviation, minimum, and maximum values for the original dataset and each imputation run are shown in their respective columns. The same treatment is done for V3, V4, V5, and V6 in Tables 3–6, respectively. The only variation is in the number of missing values in each variable which brings about these numbers of imputed values: V3 with 35 missing values has 175 imputed values, V4 with 227 has 1,135 imputed values, V5 with 374 missing values has 1,870 imputations, and V6 with 11 missing values has 55 imputations.

Observing the characteristics of the datasets before and after the imputations from Tables 2–6, the statistical mean, standard deviation, minimum, and maximum are more stable (not at much variance) after the imputation than the reduced datasets during imputation. For example, during the imputation in Table 2, the mean imputations for the five iterations are 125.52, 119.22, 145.21, 111.99, and 135.09; and the mean imputations for the five iterations of complete datasets are 121.70, 121.67, 121.84, 121.62 and 121.77, which are close to those of the original data. This shows that the imputed datasets are more reliable than the original form, especially in medical diagnosis, which deals with the issue of saving lives.

The results of the simulated datasets are shown in Fig. 4. The results of ELM classification of the complete and incomplete datasets show the validation of multiple imputations technique upon Pima Indians Diabetes dataset. P0, the original incomplete dataset, has the least accuracy of 63.0431%, while all other five (P1–P5) imputed datasets perform better than P0. This proves that multiple imputations technique is a better choice for imputation for a nonmonotone missing dataset for better classification accuracy.

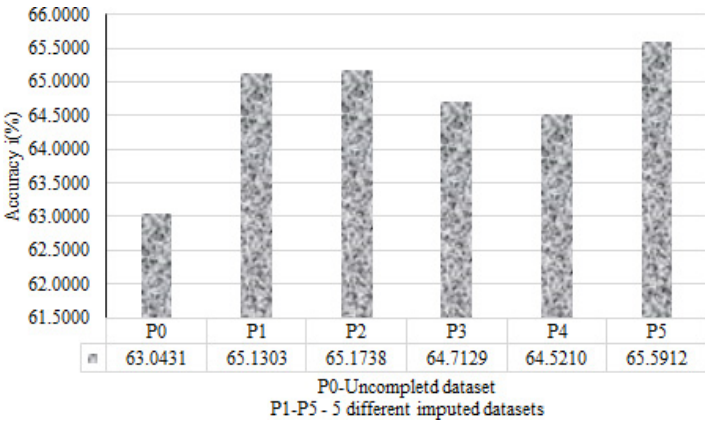


Fig. 4. ELM validation of multiple imputations technique upon Pima Indians Diabetes dataset showing that the results of classification of imputed datasets (P1–P5) are more accurate than the incomplete dataset (P0).

8. Conclusion

This paper proposed the examination of the effect of characteristics of missing data on the choice of imputation technique. The mechanism of missingness and the missing pattern were examined as the bases for the choice of imputation technique for filling missing data in the medical dataset with missing values. In this study, we considered various mechanisms for missing values — MAR, MCAR, and MNAR. Different treatments of the missing data based on the mechanisms of missingness were discussed. We backed up our study with investigations into the pattern of missingness in the Pima Indians Diabetes dataset. The observed pattern of missingness on the dataset showed that multiple imputations technique is more suitable to impute the missing values because it reflects the uncertainty of the undelaying missing data, and the imputations were so treated as shown in Tables 2–6. Our further research work will focus on the performance comparison of different classifiers on the imputed datasets, and suitable optimization technique on a favored classifier in order to improve the accuracy of classification. The work can, however, be extended to compare the accuracy of the imputed datasets with the original dataset with different classifiers like SVM, radial basis function (RBF), and ELMs.

Acknowledgments

This research has been funded by Universiti Teknologi Malaysia (UTM) under Research University Grant Vot-20H04, Malaysia Research University Network (MRUN) Vot 4L876, and the Fundamental Research Grant Scheme (FRGS) Vot 5F073 supported under the Ministry of Education, Malaysia.

References

1. M. Lichman, UCI Machine Learning Repository (2013), School of information and Computer Science, University of California, <https://archive.ics.uci.edu/ml/index.php>.
2. P. J. García-Laencina, J. L. Sancho-Gómez and A. R. Figueiras-Vidal, Pattern classification with missing data: A review, *Neural Comput. Appl.* **19** (2010) 263–282, doi: 10.1007/s00521-009-0295-6.
3. C. T. Tran *et al.*, An ensemble of rule-based classifiers for incomplete data, in *Proc. 2017 21st Asia Pacific Symp. Intelligent and Evolutionary Systems* (IEEE, 2017), pp. 7–12.
4. A. N. Baraldi and C. K. Enders, An introduction to modern missing data analyses, *J. Sch. Psychol.* **48** (2010) 5–37, doi: 10.1016/j.jsp.2009.10.001.
5. H. Gao, X. W. Liu, Y. X. Peng and S. L. Jian, Sample-based extreme learning machine with missing data, *Math. Probl. Eng.* **2015** (2015) 145156:1–145156:11, doi: 10.1155/2015/145156.
6. M. D. R. V. Gimpy, Missing value imputation in multi-attribute data set, *Int. J. Comput. Sci. Inf. Technol.* **5** (2014) 5315–5321.
7. C. D. Newgard and R. J. Lewis, Missing data: How to best account for what is not known, *JAMA* **314** (2015) 940–941, doi: 10.1001/jama.2015.10516.

8. X. Zhu, Comparison of four methods for handling missing data in longitudinal data analysis through a simulation study, *Open J. Stat.* **4** (2014) 933–944, doi: 10.4236/ojs.2014.411088.
9. J. Joseph, How to treat missing values in your data (2016), Data Science Central, <https://www.datasciencecentral.com/profiles/blogs/how-to-treat-missing-values-in-your-data-1>.
10. M. G. Kenward, The handling of missing data in clinical trials, *Clin. Invest.* **3** (2013) 241–250.
11. C. Gautam and V. Ravi, Data imputation via evolutionary computation, clustering and a neural network, *Neurocomputing* **156** (2015) 134–142, doi: 10.1016/j.neucom.2014.12.073.
12. O. A. Alade, A. Selamat and R. Sallehuddin, A review of advances in extreme learning machine techniques and its applications, in *Recent Trends in Information and Communication Technology*, Lecture Notes on Data Engineering and Communications Technologies, Vol. 5 (Springer International Publishing, 2018), pp. 885–895, doi: 10.1007/978-3-319-59427-9.
13. G. Huang, What are extreme learning machines? Filling the gap between Frank Rosenblatt's dream and John von Neumann's puzzle, *Cognit. Comput.* **7** (2015) 263–278, doi: 10.1007/s12559-015-9333-0.
14. C. V. Subbulakshmi and S. N. Deepa, Medical dataset classification: A machine learning paradigm integrating particle swarm optimization with extreme learning machine classifier, *Sci. World J.* **2015** (2015) 418060, doi: 10.1155/2015/418060.
15. B. M. Bai, N. Mangathayaru and B. P. Rani, An approach to find missing values in medical datasets, in *Proc. Int. Conf. Engineering and MIS 2015* (2015), pp. 70:1–70:7.
16. R. Armina, A. M. Zain, N. A. Ali and R. Sallehuddin, A review on missing value estimation using imputation algorithm, *J. Phys., Conf. Ser.* **892** (2017) 012004, doi: 10.1088/1742-6596/892/1/012004.
17. D. Sovilj *et al.*, Extreme learning machine for missing data using multiple imputations, *Neurocomputing* **174** (2015) 220–231, doi: 10.1016/j.neucom.2015.03.108.
18. C. F. Tsai and F. Y. Chang, Combining instance selection for better missing value imputation, *J. Syst. Softw.* **122** (2016) 63–71, doi: 10.1016/j.jss.2016.08.093.
19. P. C. Austin and M. D. Escobar, Bayesian modeling of missing data in clinical research, *Comput. Stat. Data Anal.* **49** (2005) 821–836, doi: 10.1016/j.csda.2004.06.006.
20. M. Falcaro and J. R. Carpenter, Correcting bias due to missing stage data in the non-parametric estimation of stage-specific net survival for colorectal cancer using multiple imputations, *Cancer Epidemiol.* **48** (2017) 16–21, doi: 10.1016/j.canep.2017.02.005.
21. J. Tang *et al.*, A hybrid approach to integrate fuzzy C-means based imputation method with genetic algorithm for missing traffic volume data estimation, *Transp. Res. C, Emerg. Technol.* **51** (2015) 29–40, doi: 10.1016/j.trc.2014.11.003.
22. C. D. Nguyen, J. B. Carlin and K. J. Lee, Model checking in multiple imputation: An overview and case study, *Emerg. Themes Epidemiol.* **14** (2017) 8:1–8:12. doi: 10.1186/s12982-017-0062-6.
23. Y. Liu and V. Gopalakrishnan, An overview and evaluation of recent machine learning imputation methods using cardiac imaging data, *Data (Basel)* **2** (2017) 8, doi: 10.3390/data2010008.
24. P. Diaconis and B. Efron, Computer intensive methods in statistics, *Sci. Am.* **248** (1983) 115–130.
25. G. B. Huang and L. Chen, Convex incremental extreme learning machine, *Neurocomputing* **70** (2007) 3056–3062, doi: 10.1016/j.neucom.2007.02.009.
26. Z. Shang and J. He, Confidence-weighted extreme learning machine for regression problems, *Neurocomputing* **148** (2015) 544–550, doi: 10.1016/j.neucom.2014.07.009.

27. I. Wasito and B. Mirkin, Nearest neighbors in least-squares data imputation algorithms with different missing patterns, *Comput. Stat. Data Anal.* **50** (2006) 926–949, doi: 10.1016/j.csda.2004.11.009.
28. M. Mukaka *et al.*, Is using multiple imputations better than complete case analysis for estimating a prevalence (risk) difference in randomized controlled trials when binary outcome observations are missing? *Trials* **17** (2016) 341:1–341:12, doi: 10.1186/s13063-016-1473-3.
29. A. Gelman and J. Hill, *Data Analysis Using Regression and Multilevel/Hierarchical Models*, 1st edn. (Cambridge University Press, New York, 2007).
30. Z. Zhang, Missing data imputation: Focusing on single imputation, *Ann. Transl. Med.* **4** (2016) 9, doi: 10.3978/j.issn.2305-5839.2015.12.38.
31. L. Rodwell, K. J. Lee, H. Romaniuk and J. B. Carlin, Comparison of methods for imputing limited-range variables: A simulation study, *BMC Med. Res. Methodol.* **14** (2014) 57:1–57:11, doi: 10.1186/1471-2288-14-57.
32. S. M. J. van Kuijk, W. Viechtbauer, L. L. Peeters and L. Smits, Bias in regression coefficient estimates when assumptions for handling missing data are violated: A simulation study, *Epidemiol. Biostat. Publ. Health* **13** (2016) e11598:1–e11598:8, doi: 10.2427/11598.
33. J. W. Smith *et al.*, Using the ADAP learning algorithm to forecast the onset of diabetes mellitus, in *Proc. Annu. Symp. Computer Applications and Medical Care* (1988), pp. 261–265.
34. B. Strack *et al.*, Impact of HbA1c measurement on hospital readmission rates: Analysis of 70,000 clinical database patient records, *Biomed. Res. Int.* **2014** (2014) 781670.