



Close-Knit-Regression: An Efficient Technique in Estimating Missing Completely at Random Data

Ahmed Abdulkadir^a and Bannister Jerry Zachary^{a*},
Nafisat Yusuf^a and Kabiru Musa^{a*}

^aAbubakar Tafawa Balewa University, Bauchi, Nigeria.

Authors' contributions

This work was carried out in collaboration among all authors. Author AA designed the study, wrote the first draft of the manuscript, managed the literature searches and the analyses of the study. Author BJZ, NY and KB wrote the protocol and performed the statistical analysis. All authors read and approved the final manuscript.

Article Information

DOI: 10.9734/AJPAS/2023/v24i3528

Open Peer Review History:

This journal follows the Advanced Open Peer Review policy. Identity of the Reviewers, Editor(s) and additional Reviewers, peer review comments, different versions of the manuscript, comments of the editors, etc are available here:

<https://www.sdiarticle5.com/review-history/53035>

Original Research Article

Received: 04/01/2020
Accepted: 09/03/2020
Published: 05/09/2023

Abstract

The study aimed at using the Close-Knit Regression (CKR) technique to approximate values absent because of the missing completely at random mechanism. Bivariate datasets were generated and simulated for MCAR mechanism at low (10%) and high (60%) rates. The CKR method was used and compared alongside other single imputation techniques like mean imputation, simple regression and K- Nearest Neighbors (K-NN). The difference between parameter estimates like mean, correlation coefficient (r), maximum, minimum and standard deviation which were gotten using predicted data and those using the original data as well as assessment of error rates like mean absolute error (MAE) and root mean square error (RMSE) were used as metrics in deciding the efficiency of the techniques. Results showed that the CKR technique was the best from those considered, with its estimated data having parameter estimates closer to that of the original whilst

*Corresponding author: Email: batesthommie@gmail.com;

having the least error rates at 10% (MAE of 0.01 and RMSE of 0.047) and 60% (MAE of 0.021 and RMSE of 0.073) in comparison to other methods, CKR technique is a suitable single imputation technique which produces estimates close to the original data and parameters with low error rates when data are MCAR.

Keywords: *Missing completely at random; close knit regression; mechanism; parameter estimates; mean absolute error; root mean square error.*

1 Introduction

The possession of high quality data is primarily important in research studies, a statistician, no matter his level of expertise can do from little to nothing without access to reliable information on the phenomena he wishes to assess.

It is in fact safe to say, that one can not depend on the results of any investigation if the data source is not verifiable. In the real world however, data collection is affected by so many factors, ranging from human error or apparatus failure to voluntary or involuntary non response or invalid answers by some participants and even loss of life [1].

While some of the aforementioned dynamics are mitigatable, most are not within the complete control of the researcher which makes avoiding them almost hopelessly inevitable [2], leading to unwanted errors, lack of consistency alongside redundancy and inadequacy in data sets. This in turn can heavily compromise the process and outcome of data analysis if not making it impossible to proceed in some cases.

When there are no values recorded in required information fields during research, it is referred to as *missing data* [3]. It is the lack of input, where input is needed. It can also be referred to as information that should have been present but isn't, for peculiar reasons [4]. According to McKnight et al. [5] the causes of missing data can be usually traced to:

- (a) The study participants, which entails errors on the part of subjects or their refusal to provide information for personal reasons (participant characteristics).
- (b) The study design, having to do with the structure of the data collection methods and how its tedious and overbearing nature could discourage participants from providing complete data (design characteristics).
- (c) The interaction of (a) and (b) above that has to do with the repercussions from the contact of study participants with design, an example of this is when some subjects in clinical trials are too sick to continue. There have also been cases of missing values due to the aforesaid reasons occurring in non-indigenous forms, they camouflage among genuine data making the task of spotting them a strenuous one [6].

Prevention as they say is better than cure and this applies greatly to missing data, some of the ways researchers can curb the effect of missing data is by adopting a well organised study and being meticulous with data collection [2], if however missing values occurs as it is a near definite possibility [7], there are a plethora of methods for handling them. These methods depend largely on the underlying structures and reasons for occurrences of missing values.

The course of action that led to missing values existing in a data set is referred to as the *mechanism of missing data* [8]. Little and Rubin [9] gave a deft classifying system of missing values basing mainly on their probabilities. When the probability of a variable being missing is independent of all other variables (observed and unobserved) in the data set, the mechanism in place is Missing Completely at Random (MCAR), a good example is skipping of certain items on a questionnaire by respondents due to oversight. Sometimes, the probability of a variable being missing is dependent on other observed variables in the data, this defines as Missing at Random (MAR), for example, women might exclude their age response in the demographic section of a questionnaire for sociological reasons. The last mechanism is Not Missing at Random (NMAR) and this happens when the probability of missing value occurrence is dependent on both observed and unobserved variables, take for example data on the IQ scores with data missing for subjects with low IQ values. The lack or presence of constancy in the way data are missing is referred to as its *pattern*. A univariate pattern happens

when values are absent for only one variable. When missing values are dependent on each other it is termed to have occurred in a monotonic pattern, arbitrary patterns occur in random fashion [5].

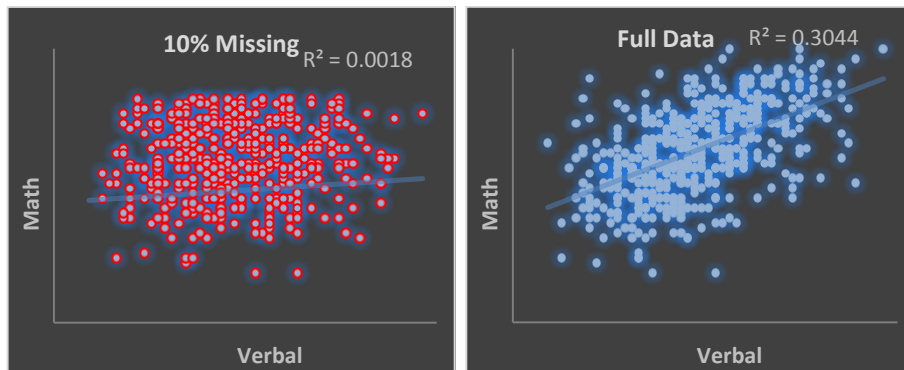


Fig. 1. Showing the effect of missing data on a scatter plot randomly generated

Notwithstanding the advent of super computers with high end estimating powers in the 21st century, the problem of missing value estimation has continued to trouble researchers and scientists alike [10]. Its predominance in datasets if not addressed, being one of the many causes of bias when estimating parameters [11], hence weakening the statistical and empirical powers of estimators. There are a plethora of techniques for handling missing data ranging from complete/available case analysis to single imputation methods, likelihood based approaches and multiple imputation techniques [12]. Single imputation being one of the most flexible and general methods is easier and more direct than other techniques this in turn makes it more popular. Single imputation techniques however, tend to ignore uncertainty and almost always underestimates variance, like it was evident in the research of Paniagua et al. [13].

This study aimed to develop and apply the close-knit-regression (CKR) approach as a single imputation method, methods, investigate its advantages and disadvantages (if any) alongside three (3) other selected single imputation techniques in widespread use, which are mean imputation, simple linear regression and K-Nearest Neighbour (K-NN). Which for a wide scope, will make use of these methods tried on generated data which will be simulated for MCAR mechanism at low and high rates of 10% and 60% respectively under a univariate pattern.

1.1 Assessment on data techniques from surveyed literature

We could classify methods for handling missing data can into four [9,5], and these are:

- a) Data Deletion Methods: Which includes List wise deletion, Pairwise deletion, Available item analysis, Individual growth curve analysis, Multisample analysis etc.
- b) Data Augumentation Methods: Comprising of Maximum Likelihood based methods, Expectation Maximization, Markov Chain Monte Carlo Method, Weighing and Dummy Code Adjustments etc
- c) Single Imputation Methods: With the following methods
 - i- Constant Replacement (like the Mean, ML mean, Median Substiution and Zero Imputation)
 - ii- Random replacement (Hot Deck, Cold Deck), Model Based (Bayesian/Monte Carlo, ML)
 - iii- Not Random replacement like the One Condition approach (Group Mean, Group Median, Last Observation carried Forward, Next Observation carried backward) and Multiple Conditions (Mean Previous Observations, Mean Subsequent Observations, Last and Next Average, Regression, Regression with Error) etc.
- d) Multiple Imputation Methods.

We will now look at past methods used by researchers in handling missing data, which we cataloged with respect to the aforementioned methods.

1.1.1 Data deletion methods

Nirelli et al. [14] on handling missing data, found the largest differences in standard errors between the original data and two simulated missing data mechanisms, MCAR and MAR to have occurred while using the complete case method. The further effect of bias appearance when using Complete Case Analysis (CCA) was seen in the work of Nakai [15] on a simulated longitudinal dataset, where even though estimated means were close to the original values showing no disadvantage, the Mean Square Error (MSE) of the estimate was doubled. They concluded that CCA method was best for low and fair missing rates (under 15%). Guan and Yusoff [11] on the other hand whilst also working a longitudinal study reported not just bias in standard error but a significant compromise in parameter estimate using CCA as opposed to the original data set and other imputation methods.

Persisting problems of sample size and test efficiency reduction was also reported in Nakai [16] study when complete case analysis was used to predict missing values in a longitudinal analysis of 1000 MCAR datasets with constant variance and AR(1) for varying correlation structures, but still positives were taken from the results when computations were made at low missing and correlation rates of 5% and $\rho=0.1$ respectively.

1.1.2 Data augmentation methods

Dong and Peng [17] in their work on demonstrating three principled data methods, which were, Multiple imputation (MI), Full information maximum likelihood (FIML) and EM. Standard errors from EM were closer to those based on the complete data. Susianto et al. [18] assessed the EM algorithm alongside four imputation methods in a comparative study. Performances were compared using mean square error (MSE) and mean absolute error (MAE). EM algorithm performed better than Markov Chain Monte Carlo Method (MCMC) but came inferior to the Yates Method. MCMC was also outperformed by EM method in terms of accuracy in the work of Takahashi [19].

Dong and Peng [17] incorporated the FIML technique as mentioned earlier, similar results were obtained under the three missing data conditions of 20%, 40% and 60%. They prosit using FIML when parameters are to be estimated cause they don't introduce "too much" randomness in data sets. The Missing Indicator Method/Missing Indicator has provided researchers with a sufficient alternative to listwise deletion, its main advantage is the sample size is usually not compromised. ML Methods have also being reviewed to be are more advantageous than MI methods [11,20,21].

1.1.3 Single imputation methods

Literatures surveyed on single imputation methods indicate that The CN2 and C4.5 algorithms are perhaps the two most simplest of all imputation methods, they in general replace missing values with the mode from entries of the variable considered. A study by Grzymala-Busse and Hu [20] categorized them both as not very good estimators of missing values. These findings were also supported by Batista and Monard [1] in later studies which compared the two aforesaid algorithms with more precise procedures like the K-NNI method. The mean imputation has been found to recurrently underestimate standard error of parameters [11,22,23]. Simple regression and using conditional means were both deemed more effective method than mean imputation [14].

1.1.4 Multiple imputation methods

Nakai et al. [16] conducted a simulation study to investigate the efficiency of four typical imputation methods with longitudinal data setting under MCAR. Tests were done at varying missing rates (5%, 30%, 50%), MI method had the least bias and best coverage probability, it was concluded to be the most effective of all imputation methods (others tried were LOCF, CCA and mean imputation). Results from the work of Schmitt et al. [24] which compared 6 different types of imputation methods used multiple imputations by chained equation (MICE) to compare the performance of four real data sets under MCAR assumptions showed that MI's performance was termed "not consistent" with its best results gotten in small data sets and the worst in large ones. It also took the longest time for estimation (about half an hour). From all six (6) techniques assessed.

Given the number of repetitious cases of missing values post data collection, a good portion of statisticians have since proposed a variety of single imputation methods that handle such inconveniences The CKR was developed to make up for some of the shortcomings of other popular single imputation methods. The proposed CKR

method is expected to not overly underestimate variance while providing more accurate estimates since the imputations are conditionally random, systematic and likely to be different for each missing point.

2 Materials and Methods

Data simulations will be performed in R using the `ampute` function as proposed by Schouten et al. [25] which works with `mice`, `vim` and `MASS` packages. Continuous defined datasets of one thousand observations (N=1000) will be generated, which will be composed of two fairly correlated variables (V1, V2) as most real world variables are, be aware that the covariance matrix should be semi definite. Summary of variables and conditions used in this study will be specified in Table 1 below.

Table 1. Summary of variables and conditions used in this study

Variables	Correlation (r)	Missing mechanism	Missing pattern	Distribution	Missing rate	Techniques
V1(Independent variable) and V2(Dependent Variable)	Fairly correlated (0.4)	MCAR	Univariate on dependent variable (V2)	Both are standard normal V1, V2 ~ N (0,1)	10% & 60%	Single Imputation: 1) Mean Imputation 2) K-NN 3) Regression Imputation 4) CKR (Proposed Method)

2.1 Data set simulation

After the data set generation is complete, `ampute` function has several other arguments which specify the nature of your missing data. First is the proportion, which in our study will vary from 10% to 60%. Next is the specification of missing mechanism which for our study will be of the univariate kind acting on the dependent variable.

Another important argument in the `ampute` function is the one that lets you select the frequency of missing-ness across the data, `ampute` divides original data into multiple subsets, where the number of subsets which has values in proportions that sum must equal one using a single value will suit the univariate pattern assigned earlier.

Specification of the mechanism to be MCAR is the next step after which assigning the weights which determines the relative missing-ness in the data set with respect to the variables, A weighted sum score uses a linear regression with coefficients assigned, it is of the form

$$wss_i = W_1 \cdot V1_i + W_2 \cdot V2_i \tag{1}$$

where wss_i is the weighted sum score of case i , $V1_i$ and $V2_i$ are the variable values of case i and W_1 and W_2 are the specified weights.

Keeping in mind that MCAR is completely random and the variables don't influence its being missing, a zero weight is assigned to both variables. The last argument in `ampute` is not applicable to the MCAR mechanism.

2.2 Techniques considered in the study

A total of four single imputation techniques were considered in this study, three (3) are already commonly used and the -last is the method proposed, they are:

- a) Mean Imputation
- b) K-NN
- c) Simple Regression
- d) CKR (Proposed Method)

A brief description of these methods will be in focus.

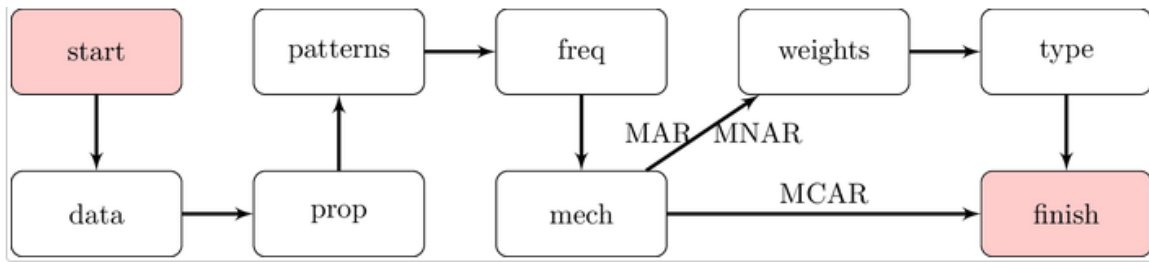


Fig. 2. Flowchart showing steps in the ‘Ampute’ process (Adapted from Schouten et al. [25])

2.2.1 Mean imputation

The mean imputation is one of the most popularly known methods. It replace the missing values in a variable with the mean of all present values for continous data, while it replaces the missing values with the mode in discrete data. The disadvantages of the mean method is mainly on how it tends to underestimate variance by repeating values since the mean is a constant, correlation coefficient values are also stunted cause of the repititive nature of its outcome. Mathematically If x_{ij} of the k -th class C_k is missing, then it is replaced by

$$V2_i = \sum_{i: v2_i \in C_k} \frac{v2_i}{n_k} \quad (2)$$

2.2.2 K-Nearest Neighborhood (K-NN)

The K-NN method replaces the missing values by considering the given number of occurences that are most similar to the value of interest. It has numerous advantages, as it can be used for both qualitative and quantitative features in a data set, it doesn’t make use of a predictive model too, the K-NN method also considers the correlation structure of the data. The first set back of this method is in the consideration of what distance function to use, it also requires a lot of time which is based on the choice of k . The procedure is as follows:

- Given a data set $V2$, Divide $V2$ into two parts. Let $V2_{mis}$ be the set containing the instances in which at least one of the features is missing. The remaining instances with complete feature information form a set called $V2_{pres}$.
- For each vector $V2$ in $V2_{mis}$: Divide the instance vector into observed and missing parts as $V2 = [V2_{obs}; V2_{miss}]$. Calculate the distance between $V2$ and all the instance vectors from the set $V2_{pres}$. Use only those features in the instance vectors from the complete set $V2_{pres}$, which are observed in the vector $V2$.
- Use the K closest instances vectors (K -nearest neighbors) and perform a majority voting estimate of the missing values for categorical attributes. For continuous attributes replace the missing value using the mean value of the attribute in the k -nearest neighborhood. The median could be used instead of the mean in cases of categorical data.

The K-NN takes into consideration the correlation structure of the data set and is so an improvement on using the mean.

2.2.3 Regression method

This is usually used for univariate or monotone missing data pattern. The first step involves building a model from the observed data. Predictions for the incomplete cases are then calculated under the fitted model, and serve as replacements for the missing data. The demerits of this method is usually the model estimated values are usually more artificial than the true values, also the technique could suffer from a lack of precision especially if there are no relationships among the values in the data set and the attribute with missing data, it is sometimes a tedious process too, since depending on the number of variables with missing data, so many models could be created.

Suppose that there are 2 variables $V1, V2$ in a data set and missing data are uniformly or to impute the missing values for a variable, one first constructs a regression model using observed data on $V1$ through $V2$.

$$V2 = \beta_0 + \beta_1 V1 \tag{3}$$

The regression model in above yields the estimated regression coefficients β_0, β_1 and the corresponding covariance matrix. Based on these results, one can impute one set of regression coefficient. from the sampling distributions of β . Next, the missing values in $V2$ can be imputed by plugging β_0, β_1 into the above equation and adding a random error ε resulting in one complete data set.

2.2.4 Close-Knit Regression (CKR): Proposed method

The proposed method combines certain aspects of the K-NN regression with simple linear regression, The close-knit-regression has two stages, first the close-knit sample-selection-stage where numerical values present in the incomplete data set that we think would give us the best estimate of missing data points are selected. Then the estimation stage where linear regression is applied to the selected sample and a model is built to use in interpolating (preferably) or extrapolating missing data points. It was built to handle univariate missing patterns.

Given two fairly correlated variables ($V2, V1$). Let $V1$ (v_{1i} 's) be the complete data set of the predictor variable, and $V2$ (v_{2i} 's) the outcome variable with some missing values, for a univariate missing pattern in ($V1, V2$). To use the close-knit regression algorithm of $V2$ on $V1$ to estimate missing values in $V2$, we follow the steps below:

- a) Sort the entire data set, by re-arranging the complete predictor variable $V1$ in ascending or descending order.
- b) For a value say $V2_n$ missing in the outcome variable $V2$, compute all $|V1_n - V1_i|$'s, a set of absolute differences.
- c) Say the smallest absolute difference is obtained at $V1_i = V1_a$

$$\implies |V1_n - V1_a| < \text{all } |V1_n - V1_i| \text{ for all values of } i \text{ not equal to } a.$$

And it is so that $V1_a$ has a corresponding non-missing value in $V2$ say $V2_a$. form a set of closely knitted samples, C and add $(V1_a, V2_a)$ as the first set of element, that is $C = \{(V1_a, V2_a)\}$,

- d) i) If $V1_n - V1_a > 0$ i.e $V1_n > V1_a$ then for the next entry $V1_b$ with a corresponding $V2_b$ value, search for values closest to $V1_n$ i.e the smallest $|V1_n - V1_b|$ where $V1_n - V1_b < 0$ i.e $V1_n < V1_b$.
 ii) If on the other hand, $V1_n - V1_a < 0$ i.e $V1_n < V1_a$ then for the next entry $V1_b$ with a corresponding $V1_b$ value, search for values closest to $V1_n$. i.e the smallest $|V1_n - V1_b|$ where $V1_n - V1_b > 0$ i.e $V1_n > V1_b$.
 iii) If no such values exists as in i or ii, then for the next entry $V1_b$ with a corresponding $V2_b$ value, only search for values closest to $V1_n$ i.e the smallest $|V1_n - V1_b|$.
- e) In similar fashion, sets of bivariate entries ($V1, V2$) are added to the set C till a chosen number of elements which is the close knitted sample size (n) is reached.

$$n\{C\} = n$$

- f) Simple-linear regression involving the elements of C is then performed to obtain coefficients, these are then used to estimate the missing data point $V2_n$.
- g) The procedure is repeated till there are no missing data points in $V2$.

The logic behind this method is straightforward, once a missing data point is located in our outcome variable $V2$, find data points in the predictor variable $V1$ that are nearest to value that was supposed to have generated the missing point in $V2$. Then use a selected number of those points in $V1$ to build a model involving non missing points in $V1$ and $V2$ which will be used to give the best predictor equation of the missing point in $V2$. The method is expected to produce good parameter estimates while not inflating their standard errors.

2.3 Performance measures

The indicators used to asses the precision of the missing data techniques relative to the complete data are the correlation coefficients (r), means, minimums, maximums, ranges, mean absolute errors and root mean square errors, they are described briefly below.

2.3.1 Comparison of parameters

Firstly, the arithmetic mean of the complete data and imputed data will both be calculated and compared using the basic formula given by:

$$\bar{V} = \frac{1}{n} \sum_{i=1}^n \hat{V}_i \tag{4}$$

Where \bar{V} is the mean of the data in focus, n is the size, and \hat{V}_i the data points. The mean will tell us about the comparative centrality of our datasets. Next, the correlation and standard deviation of the complete data and imputed data will also both be estimated and assessed comparatively using the Pearson correlation coefficient formula given by the two formulas respectively

$$r_{V2V1} = \frac{\sum_{i=1}^n (V1_i - \bar{V1})(V2_i - \bar{V2})}{n s_{v1} s_{v2}} \tag{5}$$

$$s_{v2} = \sqrt{\frac{1}{n} \sum_{i=1}^n (V2_i - \bar{V2})^2} \tag{6}$$

Where for n data points, $V1_i$ and $V2_i$ are the values of both complete and estimated data points of $V1$ and $V2$, with means and standard deviations $\bar{V1}/\bar{V2}$ and s_{v1}/s_{v2} respectively. The values of each of the correlation gotten from MDTs will be compared with that of the complete data. Contrasting the correlation coefficients and standard deviation will tell us about the spread and the strength as well as direction of the bivariate linear relationships and existing in the full and imputed data sets respectively.

The maximum and minimum values from the complete and estimated data points of the outcome variable Y in focus will gotten and there on, used to calculate the range to give us a quick sense of the spread.

Maximum value of $\hat{V}_i = \max(\hat{V}_i)$, Minimum value of $\hat{V}_i = \min(\hat{V}_i)$,

Range $\hat{V}_i = \max(\hat{V}_i) - \min(\hat{V}_i)$.

2.3.2 Comparison of errors

We here will be comparing the error arising from the differences in values between the complete simulated data and that estimated we will be using the Mean absolute error (MAE) and the Root mean square error (RMSE). To compute the MAE and RMSE, the difference between the estimated dataset points (D_{est}) and complete data set points (D_{com}) will be used to get the MAE and RMSE which represents the sample standard deviation of the MAE [26].

$$MAE = \frac{\sum_{i=1}^n |D_{com} - D_{est}|}{n} \tag{7}$$

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (D_{com} - D_{est})^2}{n}} \tag{8}$$

3 Results and Discussion

Results of data analysis after simulations are presented in this chapter, the techniques were applied to the datasets altered to suit the conditions given in Table 1.

3.1 Presentation of Results

Results are shown in terms of the proximity of the parameters estimated using techniques to that from the original dataset (Table 2 and Table 3) and then consideration was given to the error rates the parameters generated (Table 4 and Table 5).

3.1.1 Comparison of parameter estimates

The summary of statistics of the originally generated data before missing conditions were implemented showed that fair correlation between the variables V1 and V2 with, $r_{(V1,V2)} = 0.4$, Our variable of concern was V2 where the minimum and maximum values were -0.98 and 8.67 resulting in a range of 9.65. V2 also had a mean and standard deviation of -0.0001 and 1.001 respectively.

Results of Table 2, at 10% missing-ness for MCAR mechanism, CKR (Our proposed method) produced estimates with the best proximity to the full data with correlation coefficient of 0.39, mean of -0.003, minimum of -0.71, maximum of 8.3 and a range of 9.01. Results of simple linear regression were closely related to those of k-NN and Mean. The simple linear regression technique produced results with a correlation coefficient of 0.37, a mean of 0.001, minimum of -0.5, as well as a maximum and range of 8.2 and 8.7 respectively. The mean imputation technique produced data points with a correlation of 0.36, mean and range of 0.01 and 7.74 respectively while having a minimum of -0.44 and a maximum of 7.30. For K-NN imputation, the generated data points had a correlation coefficient of 0.36 with mean and a range of 0.03 and 8.03. The least value was -0.63 and the highest was 7.4.

Table 2. The four parameters before and after estimation from MCAR with techniques at 10% rate

Parameter	MCAR @ 10%				
	Full data	Single imputation technique			Proposed method CKR
		Simple Reg	Mean	K-NN	
<i>r</i>	0.40	0.37	0.36	0.36	0.39
Mean	-0.0001	0.001	0.01	0.003	-0.003
Min	-0.98	-0.5	-0.44	-0.63	-0.71
Max	8.67	8.2	7.30	7.4	8.3
Range	9.65	8.7	7.74	8.03	9.01
Std. Dev.	1.001	1.3	0.74	1.51	1.4

After the missing rate was increased to 60%, results as seen in Table 3 showed that CKR estimated data sets produced results with the best correlation estimate of 0.34. The coefficients of correlation produced by using Simple regression, mean and K-NN were 0.31, 0.3 and 0.29 respectively. CKR had the best mean estimate from the single imputation methods with a value of -0.12, values from KNN, Simple Regression and mean were the next in line with 0.14, 0.15 and 0.11 respectively. The proposed CKR produced a data set with a range of 8.51. K-NN, mean and simple linear regression produced data sets with ranges of 8.22, 8.7 and 7.81 respectively. Our proposed method also produced data points with a minimum of -0.12. Other single imputation techniques like K-NN, Mean and simple regression had their least figures as -0.23,-0.1 and -0.4 respectively. Simple linear regression, CKR and KNN methods produced maximum estimates of 7.99, 7.9 and 7.71.

Table 3. The four parameters before and after estimation from MCAR with techniques at 60% rate

Parameter	MCAR @ 60%				
	Full data	Single imputation technique			Proposed method CKR
		Simple Reg	Mean	K-NN	
<i>r</i>	0.40	0.31	0.30	0.29	0.34
Mean	-0.0001	0.11	0.15	0.14	-0.12
Min	-0.98	-0.4	-0.1	-0.23	-0.61
Max	8.67	7.99	8.1	7.71	7.9
Range	9.65	8.5	7.81	8.22	8.51
Std. Dev.	1.001	1.28	0.66	1.39	1.3

3.1.2 Comparison of parameter estimates

The MAE and RMSE values are shown in Tables 4 and 5. Small values are in general preferable as they imply better accuracy of missing data techniques. We earmarked (in boldface) small MAE values, with those less than or equal to (\leq) 0.01 being indicative of methods with good precision.

Results of Table 4 show error rates from estimations of MCAR simulated datasets at 10 %, In general low MAE values were from CKR and K-NN techniques which were each 0.01. Simple regression and mean imputation techniques had MAE values of 0.02 and 0.04 respectively. Using the CKR method gave us an RMSE of 0.047. K-NN, Simple regression and the mean imputation gave us RMSE values of 0.048, 0.064 and 0.074 respectively.

Table 4. Errors between original and predicted data from MCAR at 10% rate

Error	MCAR @ 10%			
	Mean	Single imputation		Proposed method
		Simple Reg.	K-NN	CKR
MAE	0.04	0.02	0.01	0.01
RMSE	0.074	0.064	0.048	0.047

MAE values in boldface are less than or equal to (\leq) the 0.01 threshold

After Intensifying the missing-ness to 60% as seen in Table 5, The proposed CKR method gave us an MAE of 0.02 and an RMSE of 0.073. For the K-NN method MAE value was 0.07 while simple regression and mean imputation had values of 0.05 and 0.09 respectively. The RMSE value from using the K-NN method was 0.101. The mean imputation technique had the highest RMSE with a value of 0.117 and that for simple regression was a value of 0.078 which was higher than that of our proposed CKR method.

Table 5. Errors between original and predicted data from MCAR at 60% rate

Error	MCAR @ 60%			
	Mean	Single imputation		Proposed method
		Simple Reg.	K-NN	CKR
MAE	0.09	0.05	0.07	0.02
RMSE	0.117	0.078	0.101	0.073

4 Discussion

The missing mechanism considered in this study was MCAR at two (2) missing rates (low or 10% - high or 60%) which was simulated on a bivariate dataset with a univariate missing pattern on the outcome variable V2 after which the techniques were applied and data analysis on estimated data took place. The performance of the methods were compared regarding parameter estimates such as correlation coefficients, means, standard deviation/error, minimum, maximum and range alongside MAE and RMSE error metrics.

Results show that all single imputation techniques tended to produce consistent parameter estimates in MCAR simulated data sets at all missing rates considered which was expected since the methods didn't have to deal with problems of non-normality [15]. While this is so, it is important to consider that the precision of all techniques reduced slightly as missing-ness increased from 10% to 60%. The mean imputation also produced reasonable estimates, which is largely due to the complete randomness of our missing values making the reduced sample a random subset of our original data as suggested by Nakai et al. [16]. Our proposed CKR regression performed as the best among all methods considered as it gave closer estimates to the original and didn't grossly underestimate our standard error as presumed. MAE rates for the CKR and MICE technique fell on and below the postulated threshold of 0.01 respectively. The simple regression and K-NN techniques in general faired better than mean imputation which had the highest MAE and RMSE rates from 10% to 60% missing-ness. The findings are consistent with those in reviewed literature and confirm their recommendations [1,22,26,27].

5 Conclusion and Recommendations

In accordance to our aim of developing and investigating the efficacy of CKR, it was found to be well suited for MCAR mechanism as it outperformed other single imputation techniques, this was evident in the nearness of its parameter estimates to that of the original data and its relatively low MAE and RMSE rates, the performance of K-NN and Simple regression were very nearly at par. The slight superiority of the CKR over the two previously mentioned techniques was attributed to the idea that the concept of CKR is mainly the amalgamation of them both with only nuances in execution the proposed CKR also proved to be the most robust among all single imputation techniques as changes in its error rates while increasing missing proportions were low. The CKR technique was concluded to be an effective single imputation technique in comparison to its counterparts considered in this study, it was seen to perform its very best in MCAR conditions having low missing rates of about 10%.

Competing Interests

Authors have declared that no competing interests exist.

References

- [1] Batista GEAPA, Monard MC. An analysis of four missing data treatment methods for supervised learning. *Applied Artificial Intelligence*. 2003;17(5–6):519–33.
- [2] Kang H. The prevention and handling of the missing data. *Korean Journal of Anesthesiology*. 2013;64(5):402-6.
- [3] Vaishnav RL, Patel KM. Analysis of various techniques to handling missing value in dataset. *International Journal of Innovative and Emerging Research in Engineering*. 2015;2(2):191–95.
- [4] Nwakuya MT, Nwabeuze JC. Relative efficiency of estimates based on percentages of missingness using three imputation numbers in multiple imputation analysis. *European Journal of Physical & Agricultural Sciences*. 2016;4(1):63–69.
- [5] McKnight PE, McKnight KM, Sidani S, Figueredo AJ. *Missing data: A gentle introduction*. The Guilford Press. New York London; 2007.
- [6] Buchman S. *Overview of approaches for missing data*; 2018.
- [7] Nakai M, Ke W. Review of the methods for handling missing data in longitudinal data analysis. *Int. Journal of Math. Analysis*. 2011;5(1):1–13.
- [8] Siddique J, Harel O, Crespi CM. Addressing missing data mechanism uncertainty using multiple-model multiple imputation: Application to a Longitudinal Clinical Trial. *Annals of Applied Statistics*. 2012;6(4):1814–37.
- [9] Little RJA, Rubin DB. *Statistical analysis with missing data*. A John Wiley & Sons, Inc., Publication; 1997.
- [10] Schafer JL. *Analysis of incomplete multivariate data*. Library of Congress Cataloging in Publication Data; 1997.
- [11] Guan NC, Yusoff MSB. Missing values in data analysis: ignore or impute? *Education in Medicine Journal*. 2011;3(1):6–11.
- [12] Little RJA, Rubin DB [Ed]. *Statistical analysis with missing data*. A John Wiley & Sons, Inc., Publication; 2002.

- [13] Paniagua D, Amor PJ, Echeburua E, Abad FJ. Comparison of methods for dealing with missing values in the EPV-R. *Psicothema*. 2017;29(3):384–89.
- [14] Nirelli LM, Larsen MD, Croghan IT, Schroeder DR, Offord KP, Hurt RD. Comparison of methods for handling missing data in a collegiate survey of tobacco use. Working Paper: 3439–46; 2005.
- [15] Nakai M. Simulation study: Introduction of imputation methods for missing data in longitudinal analysis. *Applied Mathematical Sciences*. 2011;5(57):2807–18.
- [16] Nakai M, Chen DC, Nishimura K, Miyamoto Y. Comparative study of four methods in missing value imputations under missing completely at random mechanism. *Open Journal of Statistics*. 2014;4:27-37. Available:<https://pdfs.semanticscholar.org/620b/c6d3e15b2e5b252ef3b1f53c9148b6989148.pdf>
- [17] Dong Y, Peng CYJ. *Principled Missing Data Methods for Researchers*”. SpringerPlus; 2013. DOI; 10.1186/2193-1801-2-22. PubMed
- [18] Susianto Y, Notodiputro KA, Kurnia A, Wijayanto H. A comparative study of imputation methods for estimation of missing values of per capita expenditure in central java a comparative study of imputation methods for estimation of missing values of per capita expenditure in central java. *IOP Conf. Series: Earth and Environmental Science* 58. 012017; 2017.
- [19] Takahashi M. Statistical inference in missing data by MCMC and Non-MCMC multiple imputation algorithms: Assessing the effects of between-imputation iterations. *Data Science Journal*. 2017;16(37):1–17.
- [20] Grzymala-Busse JW, Hu M. A Comparison of several approaches to missing attribute values in data mining” *Rough sets and current trends in computing*. 2001;2005(Chapter 46):378–85.
- [21] Allison PD. Handling missing data by maximum likelihood. *SAS Global Forum 2012 Statistics and Data Analysis*. 2012;1–21.
- [22] Truxillo C. A comparison of missing data handling methods. SAS® Institute Inc, Cary, NC; 2002.
- [23] Acuña E, Rodriguez C. The treatment of missing values and its effect on classifier accuracy. *Classification, Clustering, and Data Mining Applications*. 2004;639–47.
- [24] Schmitt P, Mandel J, Guedj M. A comparison of six methods for missing data imputation. *Journal of Biometrics and Biostatistics*. 2015;6(1):1–6.
- [25] Schouten RM, Lugtig P, Vink G. Generating missing values for simulation purposes: A multivariate amputation procedure. *Journal of Statistical Computation & Simulation*. 2018;88(15):2909–2930.
- [26] Nookhong J, Kaewrattanapat N. Efficiency comparison of data mining techniques for missing-value imputation. *Journal of Industrial and Intelligent Information*. 2015;3(4):305–9.
- [27] Zhang Z. Missing data imputation: Focusing on single imputation. *Annals of Translational Medicine*. 2016;4(1):9.

© 2023 Abdulkadir et al.; This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Peer-review history:

The peer review history for this paper can be accessed here (Please copy paste the total link in your browser address bar)

<https://www.sdiarticle5.com/review-history/53035>