

Determination of the Best Imputation Technique for Estimating Missing Values when Fitting the Weibull Distribution

Ahmad Shukri Yahaya
Nor Azam Ramli
Fauziah Ahmad

Clean Air Research Group
School of Civil Engineering
Universiti Sains Malaysia
Engineering Campus, 14300 Nibong Tebal
Seberang Perai Selatan
PULAU PINANG.

Norlida Mohd Nor
Muhammad Nor Hakim Bahrim
School of Mathematical Sciences
Universiti Sains Malaysia
11800 Minden, PULAU PINANG

Abstract

Missing values are common in any scientific work. Nine single imputation techniques were compared to determine the best technique for estimating missing values when fitting Weibull distributions. Simulation technique was used to obtain random variables for the Weibull distributions. Samples of sizes 30, 50, 200 and 300 were used to simulate the Weibull distribution. To determine the best imputation technique, four error measures were used that are the normalized absolute error, root mean square error, index of agreement and root mean square error. This study shows that no single imputation technique is the best for each sample size and for each percentage of missing values.

Keywords Weibull distribution, Imputation techniques, Performance indicators

1. INTRODUCTION

Missing value is very common in many scientific fields such as statistical, clinical, psychology, environmental research and others. In environmental research especially in air pollution modelling, observations are usually collected by using specialized machines to obtain air pollutants concentrations, wind speed, wind direction, rainfall amount and temperature. These machines sometimes failed to register the observations due to machine malfunction and sometimes when the machine needs servicing. This will result in the occurrence of missing values. Missing data hinder the ability to make exact conclusion or interpretations about the observation. Thus various missing value techniques have been developed to overcome this problem.

Little and Rubin (2002) has discussed about three classifications of missing data that are missing completely at random (MCAR), missing at random (MAR) and non-ignorable missing data (NMAR). Missing completely at random (MCAR) means that the missing value does not depend on any data or variable either observed or missing. The missing data occurs randomly in the datasets. Meanwhile, data are said to be missing at random (MAR) if the failure to observe a value does not depend on the value that would be observed. Non-ignorable missing data (NMAR) occur when missing values are not randomly distributed across observations. For MCAR and MAR, missing values can be deleted. Two common procedures for deletion of missing data are listwise deletion and pairwise deletion. Little and Rubin (2002) indicates that listwise deletion is not an appropriate way to handle missing data because by deleting the data, it will decrease the sample size and resulting in the decreasing of statistical power. Meanwhile the pairwise deletion may preserve more information. Thus imputation of missing values should be considered. The use of single imputation techniques for estimating missing values when fitting the Weibull distribution for carbon monoxide data was discussed by Yahaya et al. (2005).

They used the mean-based imputation technique and found that the mean estimation method using one datum above and one datum below is the best method. The percentage of missing values is about seven percent. The mean-based imputation technique was used by Yahaya et al. (2008) to estimate the missing values. They used simulation techniques to obtain random observations from five different Weibull distributions that represent left-skewed distributions, normal distribution and right-skewed distributions. The percentage of missing values used was 5%. This simulation study shows that the mean all method is the best imputation method to be used. Mohamed Noor Norazian (2008) describes the use of single imputation techniques to estimate missing values for a one year hourly PM₁₀ data. They used the linear quadratic and cubic interpolation techniques, nearest neighbour, mean-before-after and mean-before as the imputation techniques. Simulated percentage missing values of 5%, 10%, 15%, 25% and 40% were used. They concluded that the best imputation technique is mean-before-after. Junninen et al. (2004) describes in detail various single and multiple imputation techniques that can be used for six air quality datasets. They found that the more complex imputation techniques known as hybrid models and multiple imputations are better than the other techniques.

Plaia and Bondi (2006) analyzed the space-time variability of PM₁₀ concentration using meteorological variables. This paper uses space-time information on PM₁₀ level concentrations in eight monitoring sites that are called Site-Dependent Effect method (SDEM). Four simulated incomplete data were generated and five methods including SDEM have been applied. They found that the best imputation techniques are those that consider both space and time information that is the SEDEM method. Weibull distribution have been used to fit distributions in air pollution studies and to determine return periods (Yusoff et al., 2009; Seinfeld and Pandis, 1997; Maffei, 1998). This distribution have also been used successfully in fitting distributions for wind speed (Shoji, 2005; Jaramillo and Borja, 2004 and Yahaya et al. (2007). Weibull distribution has also been used in life testing and reliability theory. This paper compares nine single imputation methods for estimating missing values for the Weibull distributions of sizes 30, 50, 200 and 300. Five percentages of missing values were chosen at random that are 5%, 10%, 15%, 20% and 25%. These imputation techniques were chosen because it can be easily computed.

2. MATERIALS AND METHODS

2.1 The Weibull distribution

This distribution was originally derived by Fisher and Tippet in 1928 as an asymptotic extreme value distribution. In 1939 the Swedish physicist Weibull derived the same distribution on the basis of practical requirements in the analysis of material breaking strength. It was not until 1951, however, when one of Weibull's articles received wide circulation among engineers concerned with modelling the statistical variation of their data, that this distribution became prominent in the engineering community. Weibull's name has since been associated with this distribution (Bury, 1999). The Weibull density function contains two parameters, sigma (σ) and mu (μ). The σ value acts as a scale parameter and μ value acts as the location parameter that determines the form and 'skew' of the distribution (Piegorisch and Bailer, 1997).

$$f(x) = \left(\frac{\lambda}{\sigma}\right) \left(\frac{x}{\sigma}\right)^{\lambda-1} \exp\left\{-\left(\frac{x}{\sigma}\right)^{\lambda}\right\} \quad x \geq 0, \quad \sigma, \lambda \geq 0 \quad (1)$$

The cumulative distribution function (cdf) takes the form

$$F(x) = 1 - \exp\left\{-\left(\frac{x}{\sigma}\right)^{\lambda}\right\} \quad x \geq 0, \quad \sigma, \lambda \geq 0 \quad (2)$$

where σ is the scale parameter and λ is the shape parameter. The scale parameter controls the spread of the distribution, and the shape parameter controls the form of the distribution.

2.1.2 Single imputation methods

Nine imputation methods were used to deal with missing values that are mean above, mean below, mean above below, median above and median below, nearest neighbour, linear interpolation, spline interpolation and regression (Yahaya et al., 2005). The description of method used in this study is provided in Table 1.

Let y_1, y_2, \dots, y_n be the n observations of a time series data and k are the missing values that are denoted by

$y_1^*, y_2^*, \dots, y_k^*$ with $k < n$. Then the observed data with missing values are given by

$$y_1, y_2, \dots, y_n, y_1^*, y_{n+1}, y_{n+2}, \dots, y_{n_2}, y_2^*, y_{n_2+1}, y_{n_2+2}, \dots, y_k^*, y_n$$

The formulae for the first estimated values using the imputation methods are given in Table 1 according to the notations given above.

Table 1 Single imputation methods

2.2 Performance Measures

Four types of indicators were used to verify the best plotting formula to estimate the parameters of the Frechet distribution. Two error measures that is the normalized absolute error (*NAE*) and the root mean square error (*RMSE*) and two accuracy measures that is index of agreement (*IA*) and coefficient of determination (R^2) were used. To produce a good estimator, the error measures must approach zero and the accuracy measure should approach one. Table 2 gives the formula for the performance indicators.

Table 2 Performance Indicators

3. DATA

Simulation of random variables for Weibull distribution was used to obtain the observed values. The values of the shape and scale parameters were chosen to represent various forms of the Weibull distribution. The range of values for the shape parameter is (0,5) and for the scale parameter it is from 20 to 90. These values were chosen because it represents the Weibull distribution that is often used in air pollution studies in Malaysia (Sedek et al., 2006; Yusoff et al., 2009; Nurulilyana Sansuddin, 2010).

Sample of sizes 30, 50, 100, 200 and 300 were used. Each of the sample sizes were replicated ten times for each Weibull distribution. For each sample sizes, 100000 random variables were generated. Multiplicative congruential generator (Law and Kelton, 2000) of the form

$$X_i = (397204094 X_{i-1}) \bmod (2^{31} - 1)$$

was used to simulate the random variables. The above multiplicative congruential generator was found to have good statistical properties.

From these simulated Weibull random variables 5%, 10%, 15%, 20% and 30% missing values were created at random. Air pollution data are usually obtained using automated machines and it was found that usually about less 20% of data are missing (Sedek et al., 2006; Yusoff et al., 2009; Nurulilyana Sansuddin, 2009). These missing values were then estimated using the eleven imputation methods. The errors between observed and estimated missing values were then obtained and the best imputation method was obtained.

4. RESULTS AND DISCUSSIONS

The results of the simulation study using mean-based imputation, median-based imputation, nearest neighbour, linear interpolation, cubic spline and linear regression are given in this section. Table 3 gives the result of the performance indicators when there are 5% missing values. When sample size is 30, the best technique is nearest neighbour and followed by median above and median below respectively. For this case, only the normalized absolute error was obtained because the number of missing value is very small. When $n = 50$ median below is the best and followed by median above and nearest neighbour. The cubic spline technique dominates for $n = 200$. The median based technique is the next best imputation technique. However, nearest neighbour is the top and followed by linear interpolation and median above technique. Overall when the missing value is 5%, the best imputation techniques are nearest neighbour and the two median-based techniques.

Table 3 Performance Indicators with 5% missing values

Table 4 shows the performance of the imputation techniques when the percentage of missing value is 10%. When the sample size is 30, linear regression technique is found to be the best imputation technique. However, for sample sizes 50, 200 and 300, the nearest neighbour is the best. Thus when there are 10% of missing values, the technique to be used is the nearest neighbour and this is followed by linear interpolation and mean above.

Table 4 Performance Indicators with 10% missing values

The results of the imputation techniques for 15% missing values are given in Table 5. For sample of size 30, median below gives the best result and for sample size 200, the mean above is the best imputation technique. The best imputation technique for sample sizes 50 and 300 is the nearest neighbour. The best techniques to be considered when there are 15% missing values are linear interpolation, nearest neighbour and mean above.

Table 5 Performance Indicators with 15% missing values

Table 6 shows the performance measures at 20% missing values for all studied sample sizes. The best imputation technique for $n = 30$ and $n = 50$ is linear interpolation. For $n = 200$, the best imputation technique is the mean below whereas for $n = 300$ median below is the best. This result show that the best imputation technique varies according to sample size. However without considering the percentage of missing values, the best techniques are linear interpolation, mean above and mean below. Median below is at ranking number four.

Table 6 Performance Indicators with 20% missing values

At 25% missing values, the results of the performance indicators for the imputation techniques are shown in Table 7. Again the best imputation techniques depend on the sample sizes. Nearest neighbour technique is shown to be the best for $n = 30$ and $n = 50$ and is second best for $n = 200$. Median above is best for $n = 200$ and second best for $n = 30$ and $n = 50$. Linear interpolation technique is shown to be the best for $n = 300$. At 25% missing values, the best technique is nearest neighbour, median below and mean below.

Table 7 Performance Indicators with 25% missing values

We also look the best imputation techniques for each sample size. When the sample sizes are small ($n = 30$ and $n = 50$), the best techniques are median below and nearest neighbour. For $n = 200$, the best techniques are median below and median above. For $n = 300$, the best techniques are linear interpolation and nearest neighbour with the median above being the fourth best technique.

5. CONCLUSION

Estimating missing values are very important in order to obtain an accurate representation of the data. This paper compares nine simple imputation techniques for estimating missing values when fitting the Weibull distribution. Weibull distribution was chosen because it is widely used in air pollution studies, hydrologic studies and reliability studies. The data for the study was simulated so that it represents air pollutants concentration levels. Sample of sizes 30 and 50 were used to represent small sample and with 200 and 300 to represent large sample. The percentage of missing values used for this study are 5% representing small percentage of missing values, 10% and 15% representing medium percentage of missing values, 20% and 25% representing large percentage of missing values. The results of the study show that no single imputation technique is the best for each sample size and for each percentage of missing values. For small sample sizes, the best techniques are median above, median below and nearest neighbour and for large sample sizes, the best techniques are medium above, linear interpolation and nearest neighbour. This shows that resistant estimators give favourable results because the simulated Weibull distributions are skewed to the right and thus contains extreme values. For small percentage of missing values, the imputation technique that should be considered are nearest neighbour, median above and median below. In contrast, for medium and large percentage of missing values the techniques that should be used are linear interpolation and nearest neighbour.

6. REFERENCES

- Bury, K. (1999). *Statistical Distribution in Engineering*. New York: Cambridge University Press
- Jaramillo, O.A. and Borja, M.A.(2004). Wind speed analysis in La Ventosa, Mexico: a bimodal probability distribution case. *Renewable Energy*, 29, 1613–1630
- Junninen, H., Niska, N., Tuppurainen, K., Ruuskanen, J. & Kolehmainen, M. (2004). Methods for Imputation for Missing Values in Air Quality Data Sets. *Atmospheric Environment*, 38, 2895-2907
- Little, R.J.A. & Rubin, D.B. (2002) *Statistical Analysis with Missing Data*. 2nd edition. Canada : John Wiley and Sons
- Maffei, G. (1998) .Prediction of carbon monoxide acute air pollution episodes. Model formulation and first application in Lombardy. *Atmospheric Environment*, 33 (23), 3859 – 3872
- Mohamed Noor Norazian, Yahaya Ahmad Shukri, Ramli Nor Azam, Abdullah Mohd Mustafa Al Bakri (2008). Estimation of missing values in air pollution data using single imputation techniques, *ScienceAsia Journal*, 34, 341-345

Nurulilyana Sansuddin (2009) *Modeling Locational Differences And Prediction Of Temporal Concentration Of PM₁₀ Using Time Series Analysis*, Ph.D Thesis, University Sains Malaysi

Piegorsch, W. W. and Bailer, A. J. (1997) *Statistic for Environment Biology and Toxicology*. London: Chapman & Hall Publication

Plaia, A. and Bondi, A.L. (2006) Single Imputation Method of Missing Values in Environmental Pollution Data Sets. *Journal of Atmospheric Environment*, 40, 7316-7330

Sedek, J. N. M., Ramli, N. A. and Yahaya, A. S. (2006) Air quality predictions using lognormal distribution functions of particulate matter in Kuala Lumpur, Malaysia. *Journal of Environmental Management*, 7, 33 – 41

Seinfeld, J. H. and Pandis, S. N. (1997). *Atmospheric Chemistry and Physics, from Air Pollution to Climate Change*, New York: Wiley

Shoji, T. (2005). Statistical and geo-statistical analysis of wind: A case study of direction statistics, *Computers & Geosciences*, 32, 1025-1039

Yahaya, A.S., Ramli, N.A, and Yusof, N.F.F.M. (2005) Effect of estimating missing values on fitting distributions, *Proceeding of The International Conference on Qualitative Sciences and Its Applications (ICOQSA)*, 1-5. 6-8 December, Penang, Malaysia

Yahaya, A.S., Ramli, N.A., Abdul Wahab, A.A.(2007). Finding The Best Wind Speed Distribution: A Case Study, *World Engineering Congress 2007 (WEC2007)*, 5-9 August 2007, Penang, Malaysia, 14-19.(CD Proceedings

Yahaya, A.S., Abdul Hamid, H., Mohd Nor, N., Wong, L. S. (2008) Comparison Of Mean Imputation Techniques In Estimating Missing Values For The Weibull Distribution, *The Poceedings of the 16th National Mathematical Sciences Symposium, Symposium*, 2-5 June 2008, Kota Bharu, Kelantan, 183-188

Yusoff, N.F.F., Ramli, N.A., Yahaya, A.S., Sansuddin, N., Ghazali, N.A., AlMadhoun, W.A. (2009). Monsoonal Differences and Probability Distribution of PM10 Concentration. *Environmental Monitoring & Assessment* 163(1-4), 655-667

6. TABLES AND FIGURES

Table 1 Single imputation methods

Methods	Description and Formula
Mean Above - The missing values will be replaced by the mean value of the above missing data	$y_1^* = \frac{1}{n_1} \sum_{i=1}^{n_1} y_i$
Mean Below- The missing values will be replaced by the mean value of the below missing data	$y_1^* = \frac{1}{(n_2 - n_1 + 1)} \sum_{i=n_1+1}^{n_2} y_i$
Mean Above Below - The missing values will be replaced by the average of one existing data above and below the missing values	$y_1^* = \frac{1}{2}(y_{n_1} + y_{n_1+1})$ (for data ranked from smallest to largest)
Median Above - The missing value will be replaced by the median value of the above missing data	$y_1^* = \begin{cases} y\left(\frac{n_1 + 1}{2}\right) & n \text{ odd} \\ \frac{y\left(\frac{n_1}{2}\right) + y\left(\frac{n_1}{2} + 1\right)}{2} & n \text{ even} \end{cases}$ (for data ranked from smallest to largest)
Median Below - The missing value will be replaced by the median value of the below missing data	$y_1^* = \begin{cases} y\left(\frac{n_2 - n_1 + 1}{2}\right) & n \text{ odd} \\ \frac{y\left(\frac{n_2 - n_1 + 1}{2}\right) + y\left(\frac{n_2}{2} + 1\right)}{2} & n \text{ even} \end{cases}$ (for data ranked from smallest to largest)
Linear Interpolation - The missing value will be replaced by drawing a straight line between two neighbouring data	$y_i^* = y_0 + (y_1 - y_0) \left(\frac{x_i - x_0}{x_1 - x_0} \right)$ where (x_0, y_0) and (x_1, y_1) are the initial and last values
Spline Interpolation - The missing value will be replaced by using low degree polynomials in each of the interval point	$y_i^* = a_0 + a_1x + a_2x^2 + a_3x^3$ where $x_k < x < x_{k+1}$
Nearest Neighbour - The missing value will be replaced by the nearest value	$y_i^* = y_{i-1} \text{ or } y_{i+1}$
Linear Regression - The missing value will be replaced by regression of the unobserved variables against observed ones for that datasets	$y_i^* = \beta_0 + \beta_1x_i + \varepsilon_i$ where β_0 and β_1 are the intercept and slope parameters of the regression model and are estimated by the least squares method, ε_i is the error term

Table 2 Performance Indicators

Measures	Formula
Normalized absolute error (NAE)	$NAE = \frac{\sum_{i=1}^N P_i - O_i }{\sum_{i=1}^N O_i}$
Root Mean Square Error (RMSE)	$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (P_i - O_i)^2}$
Index of Agreement (IA)	$1 - \frac{\sum_{i=1}^N (P_i - O_i)^2}{\sum_{i=1}^N (P_i - \bar{O} + O_i - \bar{O})^2}$
Coefficient of Determination (R^2)	$R^2 = \frac{1}{N} \frac{\sum_{i=1}^N [(P_i - \bar{P}) + (O_i - \bar{O})]^2}{\sigma_p \sigma_o}$

*N is the number of observations, O_i is the observed values, P_i is the predicted values, \bar{O} is the average of the observed values, \bar{P} is average of predicted values, σ_p is standard deviation of observed values and σ_o is standard deviation of observed values.

Table 3 Performance Indicators with 5% missing values

Imputation Method	NAE				RMSE				R^2				IA			
	30	50	200	300	30	50	200	300	30	50	200	300	30	50	200	300
Sample size →																
Nearest	28.93 5	3.30 5	2.35 7	0.38 2	N A	65.13 1	106.8 1	43.31 0	N A	0.2 5	0.69 0	0.76 7	N A	0.50 9	0.67 5	0.90 3
Mean Abv	199.7 4	4.89 6	1.54 4	0.59 6	N A	68.39 6	72.17 4	55.61 6	N A	0.2 5	0.59 7	0.65 8	N A	0.38 6	0.64 1	0.74 6
Mean AB	113.9 5	4.16 2	1.45 6	0.68 1	N A	59.06 3	67.75 1	61.51 5	N A	0.2 5	0.71 0	0.66 4	N A	0.37 7	0.60 6	0.64 3
Mean Bel	47.13 3	3.48 1	1.30 7	0.58 0	N A	50.71 5	65.69 7	55.32 3	N A	0.2 5	0.57 8	0.66 5	N A	0.43 9	0.65 4	0.74 6
Med Abv	100.3 9	2.67 4	1.03 8	0.55 7	N A	53.13 1	66.44 5	61.29 0	N A	0.2 5	0.63 6	0.76 2	N A	0.47 5	0.67 9	0.72 3
Med Bel	36.50 7	1.83 4	0.77 8	0.55 9	N A	45.99 7	58.26 1	61.59 7	N A	0.2 5	0.61 5	0.76 0	N A	0.44 2	0.69 3	0.72 1
Linear	226.0 7	5.31 1	1.78 3	0.43 4	N A	69.72 1	80.78 4	45.98 0	N A	0.2 5	0.66 8	0.72 6	N A	0.42 3	0.68 8	0.86 9
Spline	252.8 2	6.98 2	0.81 3	0.63 5	N A	91.47 2	58.04 7	53.15 8	N A	0.2 5	0.62 1	0.66 6	N A	0.40 1	0.72 7	0.89 7
Regression	176.5 2	5.78 1	1.33 9	0.98 7	N A	67.95 2	70.74 0	81.70 0	N A	0.2 5	0.61 4	0.48 5	N A	0.25 5	0.40 5	0.23 5

Table 4 Performance Indicators with 10% missing values

Imputation Method	NAE				RMSE				R ²				IA			
	30	50	200	300	30	50	200	300	30	50	200	300	30	50	200	300
Nearest	3.13 2	0.65 4	0.64 8	0.55 1	236.0 69	67.2 40	68.3 68	64.5 61	0.33 6	0.54 8	0.78 9	0.80 8	0.38 3	0.72 5	0.80 7	0.84 2
Mean Abv	1.46 9	0.77 4	0.59 0	0.59 1	167.2 83	76.3 54	67.1 95	71.0 37	0.32 5	0.42 2	0.76 8	0.70 3	0.52 9	0.57 8	0.75 0	0.72 4
Mean AB	1.51 2	0.86 9	0.65 6	0.60 1	165.7 24	81.4 62	73.7 01	61.9 45	0.39 9	0.40 9	0.72 3	0.74 2	0.50 0	0.49 6	0.64 5	0.71 1
Mean Bel	1.61 2	0.84 1	0.59 1	0.59 4	167.9 96	79.2 46	67.0 44	71.0 86	0.33 4	0.39 3	0.76 2	0.70 4	0.49 4	0.54 4	0.74 9	0.72 3
Med Abv	1.23 9	0.65 9	0.66 4	0.55 2	167.3 28	77.1 09	80.2 10	76.1 92	0.34 7	0.46 8	0.73 8	0.71 3	0.58 2	0.57 0	0.63 2	0.68 9
Med Bel	1.14 9	0.69 7	0.65 2	0.55 9	163.1 36	79.2 63	79.1 80	76.3 92	0.36 7	0.50 8	0.75 1	0.71 4	0.57 6	0.56 4	0.63 4	0.68 8
Linear	2.23 9	0.53 9	0.65 1	0.50 4	179.7 13	76.1 27	74.4 43	53.7 89	0.43 2	0.40 7	0.81 2	0.81 4	0.47 9	0.55 7	0.79 1	0.81 5
Spline	4.45 8	0.87 9	1.11 2	0.73 4	255.6 31	73.5 99	102. 07	66.5 14	0.41 3	0.46 8	0.78 0	0.77 0	0.44 7	0.72 4	0.78 7	0.82 7
Regression	1.22 1	1.04 0	0.88 1	0.87 3	151.6 49	89.1 10	87.7 61	73.5 26	0.38 1	0.48 2	0.60 1	0.59 8	0.58 1	0.33 9	0.34 1	0.29 6

Table 5 Performance Indicators with 15% missing values

Imputation Method	NAE				RMSE				R ²				IA			
	30	50	200	300	30	50	200	300	30	50	200	300	30	50	200	300
Nearest	1.16 9	0.64 7	0.43 1	0.29 1	49.2 42	66.2 73	56.3 08	36.7 81	0.45 7	0.56 1	0.78 9	0.84 5	0.70 8	0.70 3	0.84 3	0.93 3
Mean Abv	0.84 6	0.64 7	0.40 6	0.37 4	45.4 19	66.5 41	37.8 88	39.9 75	0.34 0	0.57 6	0.79 4	0.82 1	0.68 5	0.69 1	0.83 8	0.90 2
Mean AB	0.82 4	0.62 9	0.49 5	0.46 4	42.7 56	66.5 59	41.7 85	45.5 91	0.40 5	0.57 2	0.79 7	0.85 0	0.58 8	0.64 7	0.77 1	0.85 3
Mean Bel	0.80 3	0.68 6	0.41 2	0.38 1	41.2 84	69.4 69	37.9 37	40.1 83	0.35 5	0.58 7	0.79 0	0.81 9	0.67 2	0.63 0	0.83 6	0.90 1
Med Abv	0.70 8	0.59 7	0.40 7	0.32 4	44.9 89	76.3 48	40.5 42	41.5 70	0.40 3	0.47 0	0.83 8	0.80 9	0.67 7	0.67 9	0.83 3	0.87 9
Med Bel	0.64 4	0.71 4	0.40 3	0.32 8	40.3 82	81.0 46	40.2 17	41.8 81	0.40 5	0.52 5	0.83 5	0.80 4	0.71 3	0.60 3	0.83 2	0.87 7
Linear	0.78 1	0.67 0	0.40 4	0.32 0	44.1 11	68.1 59	46.9 80	37.6 93	0.42 0	0.63 9	0.77 8	0.86 8	0.77 6	0.69 1	0.82 0	0.92 6
Spline	2.97 6	1.28 6	0.84 2	0.60 0	133. 46	118. 86	85.2 63	53.2 34	0.40 8	0.58 0	0.56 1	0.79 1	0.56 3	0.64 6	0.69 5	0.91 8
Regression	1.00 3	0.81 8	0.82 1	0.86 9	49.3 67	75.8 26	61.3 50	82.2 99	0.42 5	0.55 7	0.71 6	0.55 6	0.55 9	0.48 6	0.29 7	0.24 0

Table 6 Performance Indicators with 20% missing values

Imputation Method	NAE				RMSE				R ²				IA			
	30	50	200	300	30	50	200	300	30	50	200	300	30	50	200	300
Nearest	0.84 2	0.66 9	0.76 2	0.28 2	93.3 36	86.5 72	54.2 62	34.7 45	0.58 9	0.66 5	0.89 8	0.85 5	0.76 4	0.64 6	0.83 6	0.93 9
Mean Abv	0.83 6	0.62 6	0.67 5	0.33 5	92.6 15	78.4 30	42.3 65	40.4 89	0.57 7	0.68 8	0.84 7	0.86 7	0.74 6	0.72 1	0.84 8	0.89 9
Mean AB	0.84 0	0.60 6	0.72 5	0.44 4	95.7 91	83.8 02	41.5 51	48.1 44	0.48 7	0.63 1	0.84 7	0.86 0	0.67 8	0.64 0	0.81 4	0.81 2
Mean Bel	0.79 0	0.52 3	0.67 4	0.33 3	94.3 90	74.8 23	42.2 43	40.9 19	0.51 7	0.63 2	0.84 6	0.86 7	0.74 9	0.72 9	0.84 8	0.89 4
Med Abv	0.76 5	0.62 5	0.55 4	0.35 1	91.3 55	84.3 67	41.8 85	42.5 72	0.54 1	0.63 8	0.86 8	0.86 7	0.73 0	0.62 9	0.84 0	0.88 4
Med Bel	0.67 2	0.58 4	0.54 3	0.35 8	89.1 58	82.8 11	41.6 99	43.6 21	0.53 0	0.63 4	0.86 7	0.86 6	0.76 3	0.63 6	0.84 0	0.87 4
Linear	0.65 8	0.54 0	0.64 0	0.32 7	86.6 12	79.9 23	43.6 78	40.9 15	0.56 6	0.66 0	0.86 8	0.87 0	0.76 8	0.69 1	0.86 5	0.89 9
Spline	1.33 1	1.17 0	1.08 1	0.70 2	107. 85	131. 02	78.1 72	67.9 49	0.54 7	0.60 8	0.85 4	0.57 6	0.69 6	0.57 1	0.78 2	0.70 8
Regression	1.09 2	0.87 1	0.95 8	0.85 5	114. 10	93.1 30	55.7 57	80.2 84	0.48 7	0.48 4	0.66 1	0.53 8	0.44 4	0.51 3	0.34 3	0.16 4

Table 7 Performance Indicators with 25% missing values

Imputation Method	NAE				RMSE				R ²				IA			
	30	50	200	300	30	50	200	300	30	50	200	300	30	50	200	300
Nearest	0.52 9	0.36 2	0.41 1	0.40 6	38.1 19	32.0 15	33.5 82	42.3 51	0.60 8	0.76 5	0.88 3	0.86 5	0.84 4	0.93 9	0.93 1	0.90 7
Mean Abv	0.69 8	0.37 9	0.42 0	0.36 9	46.9 86	36.5 25	33.5 55	36.8 48	0.44 5	0.70 5	0.88 1	0.87 3	0.65 7	0.89 9	0.91 6	0.93 2
Mean AB	0.62 0	0.44 1	0.52 9	0.44 9	43.5 40	37.9 73	38.3 60	37.2 69	0.63 3	0.70 1	0.83 9	0.87 6	0.79 4	0.83 8	0.86 9	0.89 9
Mean Bel	0.59 4	0.38 7	0.41 7	0.36 3	42.7 08	37.2 64	33.5 51	36.6 58	0.55 1	0.68 6	0.87 7	0.87 1	0.81 2	0.88 1	0.91 7	0.93 3
Med Abv	0.52 3	0.30 5	0.38 5	0.33 8	44.1 87	32.9 27	33.4 03	38.6 49	0.61 2	0.74 0	0.88 6	0.85 7	0.81 6	0.91 0	0.90 6	0.92 4
Med Bel	0.49 8	0.36 6	0.39 2	0.33 3	41.5 79	40.5 54	33.6 98	38.6 46	0.60 0	0.68 7	0.88 5	0.85 6	0.82 9	0.86 2	0.90 6	0.92 4
Linear	0.53 3	0.41 8	0.46 8	0.40 2	42.7 89	37.3 60	33.8 54	34.8 30	0.62 1	0.76 3	0.86 5	0.89 6	0.76 6	0.88 9	0.90 8	0.92 9
Spline	1.56 2	0.43 6	0.94 9	0.76 7	110. 17	32.9 96	74.9 02	78.0 48	0.55 3	0.76 8	0.73 0	0.79 1	0.60 5	0.93 4	0.80 2	0.83 5
Regression	0.99 3	0.89 8	0.90 6	0.97 0	56.8 81	71.8 95	64.1 17	72.5 89	0.54 4	0.56 0	0.68 1	0.64 8	0.56 7	0.38 3	0.35 5	0.19 7