

*Full Length Research Paper*

# **Analysis of missing values in simultaneous linear functional relationship model for circular variables**

**S. F. Hassan, Y. Z. Zubairi and A. G. Hussin\***

Centre for Foundation Studies in Science, University of Malaya, 50603 Kuala Lumpur, Malaysia.

Accepted 19 May, 2010

**Missing values arise in many research fields and is a common problem in analysis. Unlike linear data, the methods proposed to handle missing values for circular variables have not been developed. This could be contributed to the closed form of circular variables. An imputation methods based on circular mean by column and sample circular mean are proposed for the simultaneous linear functional relationship model for circular variables. Simulation studies are conducted to assess the performance of the proposed methods. The results suggest that the proposed method provide an adequate approach in handling the presence of missing values for circular variables.**

**Key words:** Circular data, circular mean, missing value, simultaneous linear functional relationship model.

## **INTRODUCTION**

Missing values arise in many research fields and it is a common problem in data analysis. In view of its common occurrence in data collection, many studies have been carried out on how to handle the data set with missing values for linear data. Many approaches have been developed in addressing missing values which begin from the simple procedure including listwise deletion and pairwise deletion. Other approaches are replacement procedure (Tsikrikitis, 2005) which includes mean substitution, hot-deck imputation and regression imputation. By using these methods, all the missing values are replaced with the calculated value based on the chosen method. In particular, by using mean substitution, all the missing values are replaced with the mean of non missing observations. In the other cases, the missing values can be replaced by mean of subgroup of which the observed values are member. Apart from that, another imputation method was used where the estimation of missing values are estimated using model based procedure approach.

By far, the most common way to handle missing values is by deleting those observations with missing values, thus, leading to a complete analysis. However, this

approach decreases the sample size of data and at the same time will reduce the power of statistics which in turn, results in biased estimates when the excluded group is a selective subsample from the study population (Barzi and Woodward, 2004). Therefore, a more pragmatic approach in handling missing values is by using the replacement procedure.

Another aspect that needs to be considered when handling the problem of missing values apart from the types of missing values, is the sample size of data. As mentioned earlier, deletion approach results in a decrease of sample size and statistical power. The imputation approach seems to be a more pragmatic approach. Nevertheless, the issue of biasness should be taken into account in the imputation method.

To date, no work have been done on missing value for circular data. This could be due to the complexity of the circular data itself and the limited statistical software available to analyse such data. In the following section, two methods of data imputation for circular variables are proposed. The imputation methods that are propose are based on the measure of central location where the circular mean substitution is used in this analysis. As an analogue to linear data, the use of mean substitution may be based on the fact that the mean is a reasonable guess of value for a randomly selected observation from a

\*Corresponding author. E-mail: [ghapor@um.edu.my](mailto:ghapor@um.edu.my).

normal distribution (Acock, 2005). In this study, the evaluations of the proposed methods were assessed using simulation studies and illustrated using the wind direction data.

**CIRCULAR DATA**

Circular or directional data (Fisher, 1993; Mardia, 1972; Jammalamadaka and SenGupta, 2001) is rather special as it is unlike the linear data as its values are distributed in a circle. In other words, data measured in the form of angles or two dimensional orientations are unlike the linear data and it cannot be treated in the same way as the linear data.

Because of disparate topologies between a circle and a straight line, difficulties can occur in the statistical analysis for this kind of data. For example, if the angles are recorded in the range  $[0, 2\pi]$  radian or  $(0^\circ, 360^\circ)$ , then the direction close to the opposite end-points are near neighbours in a metric if one refers to the topology of circle, but maximally distant in linear metric. Thus, most of the methods used in statistical analysis of linear data cannot be used in circular data due to the different topology between a circle and a straight line.

Nowadays, for a linear type of data, a wide choice of computer software such as SPSS, Minitab, MATLAB and S-Plus are readily available in the market. However, for circular or directional data, only a few softwares are currently available including Axis and ORIANA software. In this paper, the analysis of missing values in simultaneous linear functional relationship model for circular variables is carried out using S-Plus.

**PROPOSED DATA IMPUTATION OF MISSING VALUES**

In this section, we propose the new method for imputation of missing values for circular variables. Suppose that we have observations from  $q+1$  circular variables, which include missing values. We propose the following two methods for the imputation of missing variables:

**Method 1: Circular mean by column**

The imputation procedures using circular mean by column implies that for each column, the circular mean value for each column is evaluated. The column mean is then used to replace any missing values for the respective columns.

**Method 2: Sample circular mean**

Another imputation procedure that is proposed in the study is to consider sample circular mean to impute into the missing values. The sample circular mean is the mean of the whole dataset excluding the missing values.

For evaluation of the above methods, we propose to apply the simultaneous linear functional relationship model for circular variables. This model is an extension of the linear functional relationship model for circular variables which was first introduced

by Hussin (1997). The details of the model can be defined as follows. Suppose that we have the observation  $x_i$  and  $y_{ji}$  ( $i = 1, \dots, n, j = 1, \dots, q$ ), and there are corresponding unobserved "underlying variables"  $(X_i, Y_{ji})$ . We write;

$$x_i = X_i + \delta_i, \quad y_{ji} = Y_{ji} + \epsilon_{ji}, \quad \delta_i \sim VM(0, \kappa), \quad \epsilon_{ji} \sim VM(0, \nu_j),$$

where  $\delta_i$  and  $\epsilon_{ji}$  are the random variables which follows von Mises distribution with mean direction zero and concentration parameter  $\kappa$  and  $\nu_j$  respectively. Now we assume that the circular variables  $Y_j$  ( $j = 1, \dots, q$ ) are related to  $X$  by the linear relationship:

$$Y_j = \alpha_j + \beta_j X \pmod{2\pi} \tag{1}$$

where  $\alpha_j$  and  $\beta_j$  are unknown parameter.

When  $q = 1$ , the model is known as the linear functional relationship model for circular variables which have been discussed by Hussin (2003) and Caries and Wyatt (2003). These models were difference from the current model that we used in this paper because they only can be used if we have two circular variables. In the case of more than two circular variables, these two models are not longer suitable and the model [1] will be used instead.

Assuming that the ratio of error concentration parameter,  $\lambda_j$  is equal to one, then the log likelihood function is given as below;

$$\log L(\alpha_j, \beta_j, \kappa, X_1, \dots, X_n; x_1, \dots, x_n, y_{11}, \dots, y_{qn}) = -2n \log(2\pi) - (1+q)n \log I_0(\kappa) + \kappa \sum_{i=1}^n \cos(x_i - X_i) + \kappa \sum_{j=1}^q \sum_{i=1}^n \cos(y_{ji} - \alpha_j - \beta_j X_i) \tag{2}$$

By differentiating log likelihood function [2] with respect to parameters  $\alpha_j, \beta_j, X_i$  and  $\kappa$ , the maximum likelihood estimate (MLE) of parameters are given as follows;

**(i) MLE for  $\hat{\alpha}_j$**

$$\hat{\alpha}_j = \begin{cases} \tan^{-1} \left\{ \frac{S}{C} \right\} & S > 0, C > 0 \\ \tan^{-1} \left\{ \frac{S}{C} \right\} + \pi & C < 0 \\ \tan^{-1} \left\{ \frac{S}{C} \right\} + 2\pi & S < 0, C > 0 \end{cases}$$

where  $S = \sum_i \sin(y_{ji} - \hat{\beta}_j \hat{X}_i)$  and  $C = \sum_i \cos(y_{ji} - \hat{\beta}_j \hat{X}_i)$ .

(ii) MLE for  $\hat{\beta}_j$

$$\hat{\beta}_{j1} \approx \hat{\beta}_{j0} + \frac{\sum_i \hat{X}_i \sin(y_{ji} - \hat{\alpha}_j - \hat{\beta}_{j0} \hat{X}_i)}{\sum_i \hat{X}_i^2 \cos(y_{ji} - \hat{\alpha}_j - \hat{\beta}_{j0} \hat{X}_i)},$$

where  $\hat{\beta}_{j1}$  is an improvement of  $\hat{\beta}_{j0}$ .

(iii) MLE for  $\hat{X}_i$

$$\hat{X}_{i1} \approx \hat{X}_{i0} + \frac{\sin(x_i - \hat{X}_{i0}) + \sum_j \hat{\beta}_j \sin(y_{ji} - \hat{\alpha}_j - \hat{\beta}_j \hat{X}_{i0})}{\cos(x_i - \hat{X}_{i0}) + \sum_j \hat{\beta}_j^2 \cos(y_{ji} - \hat{\alpha}_j - \hat{\beta}_j \hat{X}_{i0})},$$

where  $\hat{X}_{i1}$  is an improvement of  $\hat{X}_{i0}$ .

(iv) MLE for  $\hat{\kappa}$

Estimation of  $\hat{\kappa}$  can be obtained by using the approximation given by Fisher (1993),

$$A^{-1}(w) = \begin{cases} 2w + w^3 + \frac{5}{6}w^5 & w < 0.53 \\ -0.4 + 1.39w + \frac{0.43}{(1-w)} & 0.53 \leq w < 0.85 \\ \frac{1}{w^3 - 4w^2 + 3w} & w \geq 0.85 \end{cases}$$

Hence,  $\hat{\kappa} = A^{-1}(w)$  where

$$w = \frac{1}{n(1+q)} \left\{ \sum_i \cos(x_i - \hat{X}_i) + \sum_j \sum_i \cos(y_{ji} - \hat{\alpha}_j - \hat{\beta}_j \hat{X}_i) \right\}$$

### SIMULATION STUDIES

In this section, we investigate the robustness of the two imputation methods described in the previous section by means of simulation studies. Simulation studies were carried out in order to evaluate the performance for each proposed method. For this purpose, programmes are written using S-Plus. The simulation studies are repeated for 5000 times and the values of  $X$  have been drawn from  $X \sim VM\left(\frac{\pi}{4}, 3\right)$  and without loss of generality, the values

of  $\alpha_j = 0$  and  $\beta_j = 1$  for  $j = 1, 2$  are chosen. Hence the proposed model in these simulations are given by  $Y_1 = \alpha_1 + \beta_1 X$  and  $Y_2 = \alpha_2 + \beta_2 X$ .

Two different choices of concentration parameters  $\kappa = 30$  and  $50$  for random error by assuming  $\kappa = v_j$  with sample size  $n = 100$  are considered. The values of

$\kappa$  cover a more realistic range as it is expected the random error of circular variable is less dispersed. For each sample, we randomly assign 5%, 10%, 15%, 20%, 40% and 50% of the missing values, respectively.

In these simulation studies, all parameters,  $\alpha_1, \alpha_2, \beta_1, \beta_2$  and  $\kappa$  are calculated. As for performance indicator purposes, the circular mean and circular distance ( $d$ ) were calculated for  $\alpha_1$  and  $\alpha_2$  since these two parameters are in circular form. For parameters  $\beta_1, \beta_2$  and  $\kappa$ , the mean, estimate bias (EB), and estimate root mean square error (ERMSE) were calculated as follows.

Calculation for  $\alpha_j$  where  $j = 1, 2$

i. Circular Mean,

$$C = \sum \cos(\alpha_j) \quad S = \sum \sin(\alpha_j)$$

$$\bar{\alpha} = \begin{cases} \tan^{-1}\left(\frac{S}{C}\right) & S > 0, C > 0 \\ \tan^{-1}\left(\frac{S}{C}\right) + \pi & C < 0 \\ \tan^{-1}\left(\frac{S}{C}\right) + 2\pi & S < 0, C > 0 \end{cases}$$

ii. Circular Distance,  $d = \pi - |\pi - |\bar{\alpha} - \alpha||$

Calculation for  $\beta_j$  and  $\kappa$  where  $j = 1, 2$

$$\bar{w} = \frac{1}{s} \sum \hat{w}_j$$

iii. Mean,

iv. Estimated Bias,  $EB = |\bar{w} - w|$

v. Estimated Root Mean Square Errors,

$$ERMSE = \sqrt{\frac{1}{s} \sum (\hat{w}_j - w)^2}$$

All biases were calculated based on the corresponding true value that were used in generating the data set and between the new estimated values for the data set with imputed values and labelled as  $\alpha_T$ . The biases were also calculated based on the comparison between the initial parameter which has been estimated by simultaneous linear functional relationship model for circular variables and the new parameters with imputed values and labelled as  $\hat{\alpha}_j$ . The following tables show the results obtained from the simulation studies. Method 1 refers to circular mean by column while Method 2 refers to sample circular mean.

Tables 1 and 2 show the simulation results obtained for  $\kappa = 30$  using both of the proposed methods. The results show that the new means are close to the initial

**Table 1.** Simulation results for  $\alpha_1$  and  $\alpha_2$  using proposed methods for  $\kappa = 30$ .

| Parameter             |            |                        | $\alpha_1$ |          | $\alpha_2$ |          |        |        |
|-----------------------|------------|------------------------|------------|----------|------------|----------|--------|--------|
| True value            |            |                        | 0.0000     |          | 0.0000     |          |        |        |
| Estimated value       |            |                        | 6.2708     |          | 6.2773     |          |        |        |
| Performance indicator | Parameter  | Percentage             | Method 1   | Method 2 | Method 1   | Method 2 |        |        |
| Mean                  | $\alpha_1$ | 5                      | 6.2660     | 6.2661   | 6.2697     | 6.2703   |        |        |
|                       |            | 10                     | 6.2621     | 6.2639   | 6.2667     | 6.2693   |        |        |
|                       |            | 15                     | 6.2589     | 6.2621   | 6.2630     | 6.2659   |        |        |
|                       |            | 20                     | 6.2564     | 6.2602   | 6.2621     | 6.2654   |        |        |
|                       |            | 40                     | 6.2476     | 6.2518   | 6.2548     | 6.2570   |        |        |
|                       |            | 50                     | 6.2407     | 6.2489   | 6.2446     | 6.2520   |        |        |
|                       |            | Circular distance, $d$ | $\alpha_2$ | 5        | 0.0047     | 0.0047   | 0.0075 | 0.0070 |
|                       |            |                        |            | 10       | 0.0087     | 0.0069   | 0.0105 | 0.0079 |
|                       |            |                        |            | 15       | 0.0119     | 0.0087   | 0.0142 | 0.0113 |
|                       |            |                        |            | 20       | 0.0143     | 0.0106   | 0.0151 | 0.0119 |
| 40                    | 0.0231     |                        |            | 0.0190   | 0.0225     | 0.0202   |        |        |
| 50                    | 0.0301     |                        |            | 0.0219   | 0.0327     | 0.0252   |        |        |
| $\alpha_T$            | 5          |                        |            | 0.0171   | 0.0171     | 0.0135   | 0.0129 |        |
|                       | 10         |                        |            | 0.0211   | 0.0193     | 0.0164   | 0.0139 |        |
|                       | 15         |                        |            | 0.0243   | 0.0211     | 0.0202   | 0.0173 |        |
|                       | 20         |                        |            | 0.0267   | 0.0230     | 0.0211   | 0.0178 |        |
|                       | 40         | 0.0356                 | 0.0314     | 0.0284   | 0.0261     |          |        |        |
| 50                    | 0.0425     | 0.0343                 | 0.0386     | 0.0311   |            |          |        |        |

parameters estimated using the simultaneous linear functional relationship model for circular variables as well as the true value if the percentages of missing values are smaller such as 5, 10, 15 and 20%. However, the new means suddenly diverged quite far from the initial parameter once the percentage of missing values increased beyond 20%. In other words, if the percentage of missing values is too high, for example if the percentage of missing values reaches to at least 40%, the estimation seems to diverge from the initial value and produces high value of estimate bias. Thus, it can be inferred that when the percentage of missing values reach more than 40%, the proposed method are no longer suitable in the analysis.

From Table 1, the values of circular distance ( $d$ ) for  $\alpha_1$  and  $\alpha_2$  which correspond to true value ( $\alpha_T$ ) are higher than the values which correspond to the initial parameter estimate ( $\hat{\alpha}_j$ ). It shows that the new mean with imputed values are closer to the initial parameter estimated rather than the true value used in generating the data itself. This is not a surprise as the generated data with mean imputations are quite similar to the generated data itself.

From Table 2, it can be seen that the estimate bias between the new imputed values with initial parameter estimated are smaller for  $\beta_2$  while for  $\beta_1$  and  $\kappa$ , the estimate bias between the new imputed values and true

value are smaller than the bias between and initial one. Therefore, it can be concluded that the new mean for  $\beta_2$  is closer to initial parameter estimated by proposed model, while the new means for  $\beta_1$  and  $\kappa$  are closer to the true value.

Tables 3 and 4 show the simulation results for  $\kappa = 50$  using both of the proposed methods. The results also seems to exhibit the same pattern as for  $\kappa = 30$ , where the mean values are close to the initial parameter estimated as well as the true parameter used in generated the data. The value of estimate bias (EB) and estimate root mean square error (ERMSE) also increases as the percentage of missing values increases to at least 40%. The new mean is closer to the initial parameter estimated as well as to the true parameter, but the increment in the percentage of missing values being imputed using the proposed method has led to high divergence of new mean as well as having large value of estimate bias and estimate root mean square error.

From Tables 1 - 4, it can be seen that the estimate bias of concentration parameter,  $\kappa$  gets larger and larger as the value of  $\kappa$  increases. Hence, it can be said that as the concentration parameter of random error increases, the estimation of  $\kappa$  in analyzing the missing values gives high value of estimate bias as well as their estimate root

**Table 2.** Simulation results for  $\beta_1$ ,  $\beta_2$  and  $\kappa$  using proposed methods for  $\kappa = 30$ .

| Parameter                               |            |            | $\beta_1$ |          | $\beta_2$ |          | $\kappa$ |          |
|---|------------|------------|-----------|----------|-----------|----------|----------|----------|
| True value                              |            |            | 1.0000    |          | 1.0000    |          | 30.0000  |          |
| Estimated value                         |            |            | 0.9989    |          | 1.0016    |          | 28.81327 |          |
| Performance indicator                   | Parameter  | Percentage | Method 1  | Method 2 | Method 1  | Method 2 | Method 1 | Method 2 |
| Mean                                    |            | 5          | 1.0029    | 1.0028   | 1.0072    | 1.0069   | 18.0167  | 18.2924  |
|   |            | 10         | 1.0060    | 1.0048   | 1.0089    | 1.0072   | 13.4841  | 13.5773  |
|   |            | 15         | 1.0080    | 1.0067   | 1.0109    | 1.0096   | 11.0479  | 11.1897  |
|   |            | 20         | 1.0106    | 1.0094   | 1.0112    | 1.0102   | 9.5972   | 9.7367   |
|   |            | 40         | 1.0191    | 1.0193   | 1.0179    | 1.0186   | 7.3679   | 7.4810   |
|   |            | 50         | 1.0247    | 1.0240   | 1.0263    | 1.0235   | 7.0507   | 7.1671   |
| Estimate bias (EB)                      | $\beta_j$  | 5          | 0.0040    | 0.0040   | 0.0056    | 0.0053   | 10.7966  | 10.5209  |
|   |            | 10         | 0.0071    | 0.0059   | 0.0073    | 0.0056   | 15.3292  | 15.2360  |
|   |            | 15         | 0.0091    | 0.0078   | 0.0092    | 0.0080   | 17.7654  | 17.6236  |
|   |            | 20         | 0.0117    | 0.0105   | 0.0096    | 0.0086   | 19.2161  | 19.0766  |
|   |            | 40         | 0.0202    | 0.0204   | 0.0163    | 0.0170   | 21.4454  | 21.3323  |
|   |            | 50         | 0.0259    | 0.0251   | 0.0247    | 0.0219   | 21.7626  | 21.6462  |
|   | $\alpha_T$ | 5          | 0.0029    | 0.0028   | 0.0072    | 0.0069   | 11.9833  | 11.7076  |
|   |            | 10         | 0.0060    | 0.0048   | 0.0089    | 0.0072   | 16.5159  | 16.4227  |
|   |            | 15         | 0.0080    | 0.0067   | 0.0109    | 0.0096   | 18.9521  | 18.8103  |
|   |            | 20         | 0.0106    | 0.0094   | 0.0112    | 0.0102   | 20.4028  | 20.2633  |
|   |            | 40         | 0.0191    | 0.0193   | 0.0179    | 0.0186   | 22.6321  | 22.5190  |
|   |            | 50         | 0.0247    | 0.0240   | 0.0263    | 0.0235   | 22.9493  | 22.8329  |
| Estimate root Mean square error (ERMSE) | $\beta_j$  | 5          | 0.0220    | 0.0225   | 0.0205    | 0.0201   | 15.7752  | 15.3834  |
|   |            | 10         | 0.0314    | 0.0313   | 0.0281    | 0.0285   | 21.1703  | 21.0541  |
|   |            | 15         | 0.0380    | 0.0379   | 0.0364    | 0.0360   | 23.9094  | 23.7457  |
|   |            | 20         | 0.0443    | 0.0442   | 0.0415    | 0.0412   | 25.5140  | 25.3565  |
|   |            | 40         | 0.0708    | 0.0700   | 0.0663    | 0.0710   | 27.9564  | 27.8306  |
|   |            | 50         | 0.0914    | 0.0890   | 0.0895    | 0.0855   | 28.3022  | 28.1738  |
|   | $\alpha_T$ | 5          | 0.0219    | 0.0223   | 0.0210    | 0.0206   | 12.6787  | 12.3765  |
|   |            | 10         | 0.0311    | 0.0311   | 0.0286    | 0.0289   | 16.8452  | 16.7458  |
|   |            | 15         | 0.0378    | 0.0377   | 0.0368    | 0.0364   | 19.1370  | 18.9950  |
|   |            | 20         | 0.0441    | 0.0440   | 0.0419    | 0.0416   | 20.5210  | 20.3824  |
|   |            | 40         | 0.0704    | 0.0697   | 0.0667    | 0.0714   | 22.6753  | 22.5623  |
|   |            | 50         | 0.0911    | 0.0887   | 0.0900    | 0.0859   | 22.9845  | 22.8690  |

mean square error. It is also observed that with the increase in the percentage of missing values and increase in value of concentration parameter  $\kappa$  leads to the increase in biasness for parameter  $\kappa$ .

**ILLUSTRATION USING REAL DATA SET**

As an illustration for the proposed method, the real data set which is the wind direction data collected at three

different levels so that it suits in the prior model, namely, simultaneous linear functional relationship model for circular variables was used. The dataset was recorded at Bayan Lepas airport which is located at Penang Island, north of Malaysia. The measurements was taken on July and August 2005 at time 1200, located at 16.3 m above ground level, latitude 05°18'N and longitude 100°16'E. A total of 62 observations have been recorded at three different pressures with their corresponding height as follows:

**Table 3.** Simulation results for  $\alpha_1$  and  $\alpha_2$  using proposed methods for  $\kappa = 50$ .

| Parameter              |            | $\alpha_1$ | $\alpha_2$ |          |          |          |
|------------------------|------------|------------|------------|----------|----------|----------|
| True value             |            | 0.0000     | 0.0000     |          |          |          |
| Estimated value        |            | 6.2790     | 6.2366     |          |          |          |
| Performance indicator  | Parameter  | Percentage | Method 1   | Method 2 | Method 1 | Method 2 |
| Mean                   |            | 5          | 6.2750     | 6.2757   | 6.2347   | 6.2362   |
|                        |            | 10         | 6.2724     | 6.2713   | 6.2326   | 6.2365   |
|                        |            | 15         | 6.2691     | 6.2696   | 6.2326   | 6.2366   |
|                        |            | 20         | 6.2684     | 6.2673   | 6.2316   | 6.2362   |
|                        |            | 40         | 6.2615     | 6.2616   | 6.2264   | 6.2364   |
|                        |            | 50         | 6.2598     | 6.2576   | 6.2244   | 6.2367   |
| Circular distance, $d$ | $\alpha_1$ | 5          | 0.0040     | 0.0033   | 0.0019   | 0.0005   |
|                        |            | 10         | 0.0066     | 0.0077   | 0.0040   | 0.0002   |
|                        |            | 15         | 0.0099     | 0.0094   | 0.0040   | 0.0000   |
|                        |            | 20         | 0.0106     | 0.0117   | 0.0050   | 0.0004   |
|                        |            | 40         | 0.0175     | 0.0174   | 0.0102   | 0.0003   |
|                        |            | 50         | 0.0192     | 0.0214   | 0.0122   | 0.0001   |
|                        | $\alpha_2$ | 5          | 0.0082     | 0.0075   | 0.0484   | 0.0470   |
|                        |            | 10         | 0.0108     | 0.0119   | 0.0505   | 0.0467   |
|                        |            | 15         | 0.0141     | 0.0136   | 0.0506   | 0.0466   |
|                        |            | 20         | 0.0147     | 0.0159   | 0.0515   | 0.0470   |
|                        |            | 40         | 0.0216     | 0.0216   | 0.0568   | 0.0468   |
|                        |            | 50         | 0.0234     | 0.0256   | 0.0588   | 0.0465   |

- i. at pressure 850 Hpa with 5000 m height as variable  $x$
- ii. at pressure 1000 Hpa with 300 m height as variable  $y_1$
- iii. at pressure 500 Hpa with 19000 m height as variable  $y_2$

Tables 5 and 6 show the results obtained from the analysis for the real data sets using the proposed methods as describe earlier. From the results obtained, it gives a similar trend as in the simulation studies where it can be seen that the estimates are quite good for small percentages of missing values.

Also can be concluded, the increment in percentage of missing values leads to the increment in all biases. In particular, if the percentages of missing values reach to 40% or higher, we can say that analyses give poor estimates and this can be seen from the large value of biases.

Consistent with the findings in the simulation studies, Tables 5 and 6 show that Method 2 gives relatively small value of circular distance,  $d$  in comparison to Method 1. This implies that Method 2 gives better estimation for  $\alpha_1$  and  $\alpha_2$ . The similar results also can be seen for parameter  $\beta_1$  and  $\beta_2$  where Method 2 give the better estimation compared to Method 1 based on the value of estimate bias and their estimate root mean square error for each parameter. The estimation of  $\kappa$  are consistent

as the simulation study where high value of concentration parameter will give a higher value of estimate bias and their estimate root mean square error.

## DISCUSSION

Based on the simulation studies using different concentration parameters namely  $\kappa = 30$  and  $50$ , by imputing values for missing observations in the data, the estimated value of the new mean seems to provide a good estimate.

This can be seen by small values of estimated bias and estimated root mean square error. Therefore, regardless of the value of concentration parameter, the parameter estimation has small bias so long as the percentage of missing values at most 20%. On the other hand, if the percentages of missing values reach at least 40%, the estimates produced from the data set seem inadequate. This can be seen in the high values of biases and can be said to be not acceptable. In short, we can say that if the percentage of missing values in our data is less than or equal to 20%, the analysis can be performed using the proposed methods.

Comparison between both proposed methods also can be made to determine which of the two methods perform better. From the simulation results, it can be seen that the

**Table 4.** Simulation results for  $\beta_1$ ,  $\beta_2$  and  $\kappa$  proposed methods for  $\kappa = 50$ .

| Parameter                               |                  |            | $\beta_1$ |          | $\beta_2$ |          | $\kappa$ |          |
|---|------------------|------------|-----------|----------|-----------|----------|----------|----------|
| True Value                              |                  |            | 1.0000    |          | 1.0000    |          | 50.0000  |          |
| Estimated value                         |                  |            | 0.9954    |          | 1.0167    |          | 47.7866  |          |
| Performance indicator                   | Parameter        | Percentage | Method 1  | Method 2 | Method 1  | Method 2 | Method 1 | Method 2 |
| Mean                                    |                  | 5          | 0.9984    | 0.9979   | 1.0184    | 1.0184   | 27.7332  | 27.9003  |
|   |                  | 10         | 1.0002    | 1.0015   | 1.0199    | 1.0190   | 20.2977  | 20.4992  |
|   |                  | 15         | 1.0039    | 1.0033   | 1.0210    | 1.0203   | 16.5384  | 16.7316  |
|   |                  | 20         | 1.0057    | 1.0058   | 1.0225    | 1.0214   | 14.3439  | 14.4662  |
|   |                  | 40         | 1.0116    | 1.0120   | 1.0284    | 1.0266   | 10.7750  | 10.8732  |
|   |                  | 50         | 1.0172    | 1.0183   | 1.0345    | 1.0311   | 10.2543  | 10.2923  |
| Estimate bias (EB)                      | $\hat{\alpha}_j$ | 5          | 0.0030    | 0.0025   | 0.0017    | 0.0016   | 20.0534  | 19.8863  |
|   |                  | 10         | 0.0048    | 0.0061   | 0.0032    | 0.0023   | 27.4889  | 27.2874  |
|   |                  | 15         | 0.0085    | 0.0079   | 0.0042    | 0.0036   | 31.2482  | 31.0550  |
|   |                  | 20         | 0.0103    | 0.0105   | 0.0058    | 0.0047   | 33.4427  | 33.3204  |
|   |                  | 40         | 0.0162    | 0.0166   | 0.0117    | 0.0099   | 37.0116  | 36.9134  |
|   |                  | 50         | 0.0218    | 0.0229   | 0.0177    | 0.0143   | 37.5323  | 37.4943  |
|   | $\hat{\alpha}_T$ | 5          | -0.0016   | -0.0021  | 0.0184    | 0.0184   | 22.2668  | 22.0997  |
|   |                  | 10         | 0.0002    | 0.0015   | 0.0199    | 0.0190   | 29.7023  | 29.5008  |
|   |                  | 15         | 0.0039    | 0.0033   | 0.0210    | 0.0203   | 33.4616  | 33.2684  |
|   |                  | 20         | 0.0057    | 0.0058   | 0.0225    | 0.0214   | 35.6561  | 35.5338  |
|   |                  | 40         | 0.0116    | 0.0120   | 0.0284    | 0.0266   | 39.2250  | 39.1268  |
|   |                  | 50         | 0.0172    | 0.0183   | 0.0345    | 0.0311   | 39.7457  | 39.7077  |
| Estimate root mean square error (ERMSE) | $\hat{\alpha}_j$ | 5          | 0.0187    | 0.0187   | 0.0176    | 0.0178   | 28.3970  | 28.1876  |
|   |                  | 10         | 0.0269    | 0.0270   | 0.0253    | 0.0251   | 37.2483  | 37.0155  |
|   |                  | 15         | 0.0337    | 0.0329   | 0.0299    | 0.0305   | 41.5185  | 41.3005  |
|   |                  | 20         | 0.0386    | 0.0385   | 0.0353    | 0.0349   | 43.9727  | 43.8359  |
|   |                  | 40         | 0.0561    | 0.0560   | 0.0525    | 0.0525   | 47.9229  | 47.8148  |
|   |                  | 50         | 0.0665    | 0.0681   | 0.0639    | 0.0631   | 48.4952  | 48.4530  |
|   | $\hat{\alpha}_T$ | 5          | 0.0185    | 0.0186   | 0.0255    | 0.0255   | 22.6922  | 22.5315  |
|   |                  | 10         | 0.0265    | 0.0263   | 0.0321    | 0.0314   | 29.8666  | 29.6691  |
|   |                  | 15         | 0.0329    | 0.0321   | 0.0363    | 0.0365   | 33.5447  | 33.3541  |
|   |                  | 20         | 0.0376    | 0.0375   | 0.0415    | 0.0407   | 35.7052  | 35.5838  |
|   |                  | 40         | 0.0550    | 0.0548   | 0.0586    | 0.0581   | 39.2425  | 39.1447  |
|   |                  | 50         | 0.0651    | 0.0667   | 0.0704    | 0.0689   | 39.7602  | 39.7218  |

Method 2 which is sample circular mean is a more superior approach.

Based on the values of estimate bias and estimate root mean square error for each method, it can be seen that the second method, sample circular mean, gives a relatively small bias in comparison to the first method, that is, the circular mean by column. This implies that the second method give better estimate in comparison to the first method. The second method uses the approach where it considers the circular mean for the whole data set which excludes all missing values.

### Conclusion

In this paper, imputation method using mean for simultaneous linear functional relationship model is proposed. Specifically, data sets consisting of three circular variables for three different levels of readings of measurements of wind direction with some missing observations. Two imputation methods are proposed for missing values in the data set which are Method 1 namely circular mean by column and Method 2 known as sample circular mean. Circular mean by column considers mean for each

**Table 5.** Results for  $\alpha_1$  and  $\alpha_2$  for Bayan Lepas.

| Parameter              |            | $\alpha_1$ |          | $\alpha_2$ |          |
|------------------------|------------|------------|----------|------------|----------|
| Estimated value        |            | -0.2108    |          | -0.1740    |          |
| Performance indicator  | Percentage | Method 1   | Method 2 | Method 1   | Method 2 |
| Mean                   | 5          | 0.0867     | 0.0437   | 0.1899     | 0.1580   |
|                        | 10         | 0.1807     | -0.0924  | 0.2008     | 0.1485   |
|                        | 15         | 0.2449     | -0.2886  | 0.1591     | 0.1403   |
|                        | 20         | 0.2148     | -0.4865  | 0.0574     | 0.1102   |
|                        | 40         | -0.2390    | -1.1881  | -0.2521    | 0.1005   |
|                        | 50         | -0.3749    | -1.3684  | -0.3895    | 0.0859   |
| Circular distance, $d$ | 5          | 0.2975     | 0.2545   | 0.3639     | 0.3320   |
|                        | 10         | 0.3915     | 0.1183   | 0.3748     | 0.3225   |
|                        | 15         | 0.4556     | 0.0778   | 0.3331     | 0.3143   |
|                        | 20         | 0.4256     | 0.2758   | 0.2314     | 0.2842   |
|                        | 40         | 0.0283     | 0.9774   | 0.0781     | 0.2745   |
|                        | 50         | 0.1642     | 1.1577   | 0.2155     | 0.2599   |

**Table 6.** Results for  $\beta_1$ ,  $\beta_2$  and  $\kappa$  for Bayan Lepas.

| Parameter                               |            | $\beta_1$ |          | $\beta_2$ |          | $\kappa$ |          |
|---|------------|-----------|----------|-----------|----------|----------|----------|
| Estimated value                         |            | 1.0340    |          | 0.9119    |          | 1.0259   |          |
| Performance indicator                   | Percentage | Method 1  | Method 2 | Method 1  | Method 2 | Method 1 | Method 2 |
| Mean                                    | 5          | 0.9124    | 0.9643   | 0.8425    | 0.8830   | 1.0192   | 1.0205   |
|   | 10         | 0.8316    | 0.8953   | 0.8156    | 0.8986   | 1.0365   | 1.0423   |
|   | 15         | 0.7424    | 0.8401   | 0.8066    | 0.9234   | 1.0679   | 1.0721   |
|   | 20         | 0.6907    | 0.7793   | 0.8201    | 0.9614   | 1.0991   | 1.1074   |
|   | 40         | 0.4902    | 0.7447   | 0.8820    | 1.1525   | 1.2923   | 1.3232   |
|   | 50         | 0.4802    | 0.7799   | 0.9215    | 1.2152   | 1.3894   | 1.4351   |
| Estimate bias (EB)                      | 5          | 0.1216    | 0.0698   | 0.0693    | 0.0288   | 0.0067   | 0.0054   |
|   | 10         | 0.2024    | 0.1387   | 0.0963    | 0.0132   | 0.0106   | 0.0164   |
|   | 15         | 0.2916    | 0.1939   | 0.1053    | 0.0116   | 0.0420   | 0.0462   |
|   | 20         | 0.3433    | 0.2547   | 0.0917    | 0.0495   | 0.0731   | 0.0815   |
|   | 40         | 0.5438    | 0.2893   | 0.0299    | 0.2407   | 0.2663   | 0.2973   |
|   | 50         | 0.5538    | 0.2541   | 0.0096    | 0.3033   | 0.3635   | 0.4092   |
| Estimate root mean square error (ERMSE) | 5          | 0.2722    | 0.1580   | 0.1161    | 0.1244   | 2.0187   | 2.0198   |
|   | 10         | 0.4110    | 0.2781   | 0.1756    | 0.1774   | 2.0337   | 2.0387   |
|   | 15         | 0.5255    | 0.3492   | 0.2141    | 0.2515   | 2.0606   | 2.0643   |
|   | 20         | 0.5941    | 0.4241   | 0.2462    | 0.3218   | 2.0874   | 2.0947   |
|   | 40         | 0.7836    | 0.4901   | 0.3546    | 0.6980   | 2.2566   | 2.2842   |
|   | 50         | 0.7985    | 0.4707   | 0.3729    | 0.8147   | 2.3430   | 2.3845   |

column after excluding all missing values, while sample circular mean treats all observations in number of columns as whole data sets. Finally the circular mean will be evaluated after excluding all missing values.

From the simulation study, it can be shown that Method 2, namely, sample circular mean is more superior in

comparison to Method 1. This is based on the comparison of all performance indicators which are circular distance ( $d$ ), estimate bias (EB) and estimate root mean square error (ERMSE). It can be summarized the estimations are close to the true parameter if the percentage of missing values are smaller, that is, at most



20%. From the study, it also can be said that the method is suitable for the large sample size data set, for instance  $n > 50$ . It can be seen from the good result on simulation study as well as in the real data set. At the same time, it can be seen that all biases also increased as the percentage of missing values increased and this has led to inconsistent estimation. The findings are consistent by varying the values of concentration parameter. Therefore, for simultaneous linear functional relationship model for circular variables, the approach of using mean imputation method seems adequate. Similar to linear data where mean imputation is commonly used, the mean imputation for circular variables for the linear functional relationship model produces small bias which in turn suggest good performance. Furthermore, in the model considered where variables are related to each other, the imputing approach that uses a measure of central tendency performs competitively well in terms of bias.

## REFERENCES

- Acock AC (2005). Working with missing values. *ProQuest Education Journals. J. Marriage Family*, 67(4): 1012-1028.
- Barzi F, Woodward M (2004). Imputations of missing values in practice: Results from imputations of serum cholesterol in 28 cohort studies. *Am. J. Epidemiol.*, 160(1): 34-45.
- Caries S, Wyatt LR (2003). A linear functional relationship model for circular data with an application to the assessment of ocean wave measurements. *J. Agric. Biol. Environ. Stat.*, 8(2): 153-169.
- Fisher NI (1993). *Statistical Analysis of Circular Data*. Cambridge University Press.
- Hussin AG (1997). Pseudo-replication in functional relationship with environmental application. PhD Thesis, School of Mathematics and Statistics, University of Sheffield, England.
- Hussin AG, Chik Z (2003). On estimating error concentration parameter for circular functional model. *Bull. Malaysian Math. Sc. Soc.*, 26: 181-188.
- Jammalamadaka SR, Sen Gupta A (2001). *Topics in Circular Statistics*. World Scientific Publishing Co. Pte. Ltd.
- Mardia KV (1972). *Statistics of Directional Data*. London: Academic Press Inc.
- Tsikriktis N (2005). A review of techniques for treating missing data in OM survey research. *J. Oper. Manage.*, 24: 53-62.