

# A Study of EM Algorithm as an Imputation Method: A Model-Based Simulation Study with Application to a Synthetic Compositional Data

Yisa Adeniyi Abolade, Yichuan Zhao

Department of Mathematics and Statistics, Georgia State University, Atlanta, Georgia, USA

Email: yabolade1@gsu.edu

**How to cite this paper:** Abolade, Y.A. and Zhao, Y.C. (2024) A Study of EM Algorithm as an Imputation Method: A Model-Based Simulation Study with Application to a Synthetic Compositional Data. *Open Journal of Modelling and Simulation*, 12, 33-42.

<https://doi.org/10.4236/ojmsi.2024.122002>

**Received:** January 4, 2024

**Accepted:** March 5, 2024

**Published:** March 8, 2024

Copyright © 2024 by author(s) and Scientific Research Publishing Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

## Abstract

Compositional data, such as relative information, is a crucial aspect of machine learning and other related fields. It is typically recorded as closed data or sums to a constant, like 100%. The statistical linear model is the most used technique for identifying hidden relationships between underlying random variables of interest. However, data quality is a significant challenge in machine learning, especially when missing data is present. The linear regression model is a commonly used statistical modeling technique used in various applications to find relationships between variables of interest. When estimating linear regression parameters which are useful for things like future prediction and partial effects analysis of independent variables, maximum likelihood estimation (MLE) is the method of choice. However, many datasets contain missing observations, which can lead to costly and time-consuming data recovery. To address this issue, the expectation-maximization (EM) algorithm has been suggested as a solution for situations including missing data. The EM algorithm repeatedly finds the best estimates of parameters in statistical models that depend on variables or data that have not been observed. This is called maximum likelihood or maximum a posteriori (MAP). Using the present estimate as input, the expectation (E) step constructs a log-likelihood function. Finding the parameters that maximize the anticipated log-likelihood, as determined in the E step, is the job of the maximization (M) phase. This study looked at how well the EM algorithm worked on a made-up compositional dataset with missing observations. It used both the robust least square version and ordinary least square regression techniques. The efficacy of the EM algorithm was compared with two alternative imputation techniques, k-Nearest Neighbor (k-NN) and mean imputation ( $\bar{x}$ ), in terms of Aitchison distances and covariance.

## Keywords

Compositional Data, Linear Regression Model, Least Square Method, Robust Least Square Method, Synthetic Data, Aitchison Distance, Maximum Likelihood Estimation, Expectation-Maximization Algorithm, k-Nearest Neighbor, and Mean imputation

---

## 1. Introduction

Compositional data exclusively consists of relative information. These entities are part of a broader entity. Typically, closed data or data that aggregates to a constant value, such as 100%, is commonly documented. An illustrative instance within the field of medicine involves the examination of the constituent elements present in bodily fluids such as blood and urine. The statistical linear model is frequently employed to uncover latent associations among relevant random variables due to its user-friendly nature and interpretability. In the domain of machine learning and its associated disciplines, ensuring the quality of data is a significant difficulty. The quality of the underlying data plays a crucial role in determining the quality of information obtained through Machine Learning algorithms, as these algorithms rely solely on data for their functioning. One significant concern pertaining to data quality involves the presence of missing data, particularly in compositional datasets. The linear regression model is a widely employed statistical modeling technique that is utilized across various applications to ascertain correlations between variables of interest. The method of maximum likelihood estimation (MLE) is commonly employed to estimate the parameters of linear regression by determining the values that maximize the likelihood function given the observed data. The obtained model can be utilized for doing partial effects analysis on the independent variables as well as for making predictions about future outcomes.

However, it is important to note that many datasets often exhibit missing observations. In the context of research, it is possible for participants to opt out of providing a response to a survey query, for files to finally undergo destruction, or for data to be inadequately preserved. The process of recommencing data collecting and recovery in this case will incur financial expenses and necessitate a significant amount of time. The matter pertaining to the evaluation of inadequate data necessitates attention. The expectation-maximization (EM) technique has been proposed as a potential solution for scenarios with missing data due to its robust convergence properties. The estimation of parameters in statistical models that involve unobserved variables or unobserved data is commonly achieved by the repetitive use of the Expectation-Maximization (EM) technique, which seeks to find the maximum likelihood or maximum a posteriori (MAP) estimates. The M phase of the EM algorithm involves the estimation of parameters that maximize the expected log-likelihood obtained during the E stage. The

E-step, in contrast, constructs a function that calculates the expected value of the log-likelihood, using the current estimate as the parameters. In the subsequent E phase, these estimates are subsequently employed to determine the distribution of the latent variables or missing data.

Despite the wide investigation of the EM algorithm as an imputation tool, there is a lack of knowledge regarding its effectiveness when used for compositional data. This study investigates the performance of the EM method on a synthetic compositional dataset with missing observations. Two imputation techniques, namely the robust least square version and least square, are utilized and evaluated. The EM technique is applied to simulated studies by making iterative assumptions about a compositional dataset with random missing data and outliers, assuming a normal distribution. The effectiveness of the EM method was evaluated by comparing its results with two commonly employed imputation techniques, namely k-Nearest Neighbor (k-NN) and mean imputation ( $\bar{x}$ ), in terms of Aitchison distance and covariance [1]. Based on the conducted trials, it was shown that the robust variant of the EM algorithm exhibited superior performance compared to alternative imputation strategies.

## 2. Methodology

### 2.1. Linear Regression Model

We take into consideration a one-dimensional estimator and response linear regression model. Let's say we have  $n$  observations in our dataset. We define the predictor  $X = (x_1, x_2, \dots, x_n)$ , and the response  $Y = (y_1, y_2, \dots, y_n)$ . For the  $i^{\text{th}}$  observation, we assume that  $y_i$  and  $x_i$  are related by the linear regression model and in Equation (1):

$$y_i = \beta_0 + \beta_1 x_i + \epsilon, \quad \epsilon \sim NID(0, \sigma^2) \quad (1)$$

We assume that  $x_i \sim N(\alpha, \delta^2)$ , *i.i.d.* Under such assumptions, the conditional distribution of  $Y$  given  $X$  is  $[Y | X] \sim N(\beta_0 + \beta_1 X, \sigma^2)$ . Then we can write down the joint probability density of  $X$  and  $Y$  given by

$$\begin{aligned} f(y_i, x_i) &= f(y_i | x_i) f(x_i) \\ &= \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2\sigma^2}(y_i - \beta_0 - \beta_1 x_i)^2} \times \frac{1}{\delta\sqrt{2\pi}} e^{-\frac{1}{2\delta^2}(x_i - \alpha)^2} \end{aligned} \quad (2)$$

### 2.2. Missing Values

The data is not completely observed in a lot of real-world scenarios. We expand our model to have response values  $Y$  fully observed and only  $m$  of the predictor values observed (*i.e.*,  $n - m$  predictor values missing), and the response values  $Y$  are fully observed. We can arrange the dataset so that the first  $m$  observation is fully observed.

$$X_{comp} = (x_1, \dots, x_m, x_{m+1}, \dots, x_n) = (X_{obs}, X_{miss}) \quad (3)$$

Equation (4) thus allows for the decomposition of the entire data log-likelihood for the model into the observable and missing parts.

$$\begin{aligned}
 L(\theta; X, Y) &= s_{i=1}^n L(\theta; x_i, y_i) \\
 &= -2n \log \sqrt{\sigma^2 2\pi} - 2n \log \sqrt{\delta^2 2\pi} - \frac{1}{2\sigma^2} \sum_{i=1}^m (y_i - \beta_0 - \beta_1 x_i)^2 \\
 &\quad - \frac{1}{2\sigma^2} \sum_{i=m+1}^n (y_i - \beta_0 - \beta_1 x_i)^2 - \frac{1}{2\delta^2} \sum_{i=1}^m (x_i - \alpha)^2 \\
 &\quad - \frac{1}{2\delta^2} \sum_{i=m+1}^n (x_i - \alpha)^2
 \end{aligned} \tag{4}$$

where  $\theta = (\beta_0, \beta_1, \sigma^2, \alpha, \delta^2) \in \mathbb{R}^5$

### 2.3. The EM Algorithm Formulation

The issue with the above calculation is that  $X_{mis}$  with not observed and needs to be estimated. One reasonable approach is that we simply require each  $x_{m+1}, \dots, x_n$  to be replaced by its conditional expectation given the observed data,  $X_{obs}$  and  $Y$ .

1) *E-step*:

$$\begin{aligned}
 &E \left[ \sum_{i=m+1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \right] \\
 &= \sum_{i=m+1}^n \left( (y_i - \beta_0)^2 + \beta_1^2 E(X_i^2 | y_i, \theta^*) - 2\beta_1 (y_i - \beta_0) E(X_i | y_i, \theta^*) \right) \\
 &\quad E \left[ \sum_{i=m+1}^n (x_i - \alpha)^2 \right] \\
 &= \sum_{i=m+1}^n \left( E(X_i^2 | y_i, \theta^*) - 2\alpha E(X_i | y_i, \theta^*) + \alpha^2 \right)
 \end{aligned} \tag{5}$$

where  $E[X_i | y_i, \theta^*]$  and  $E[X_i^2 | y_i, \theta^*]$  are the first and second conditional moments, respectively. Since  $X$  and  $Y$  have a bivariate normal distribution, we can derive the conditional of  $X$  given  $Y$  and  $\theta^*$

$$E[X_i | y_i, \theta^*] \sim N \left( \alpha + \frac{\beta_1 \delta^2}{\sigma^2 + \beta_1^2 \delta^2} (y_i - \beta_0 - \beta_1 \alpha), \frac{\sigma^2 \delta^2}{\sigma^2 + \beta_1^2 \delta^2} \right) \tag{6}$$

Then we can easily find the conditional first and second moment of  $X_{mis}$  given  $Y$  and  $\theta^*$ , denoted as  $M^1$  and  $M^2$  respectively.

$$M_i^1 = \alpha + \frac{\beta_1 \delta^2}{\sigma^2 + \beta_1^2 \delta^2} (y_i - \beta_0 - \beta_1 \alpha) \tag{7}$$

$$M_i^2 = \left( \alpha + \frac{\beta_1 \delta^2}{\sigma^2 + \beta_1^2 \delta^2} (y_i - \beta_0 - \beta_1 \alpha) \right)^2 + \frac{\sigma^2 \delta^2}{\sigma^2 + \beta_1^2 \delta^2} \tag{8}$$

With these terms computed above, the E-step formulation is shown in Equation (9) below.

$$\begin{aligned}
 Q(\theta, \theta^*) &= L(\theta; X, Y) \\
 &= -2n \log \sqrt{2\pi} - n \log \delta - n \log \sigma \\
 &\quad - \frac{1}{2\sigma^2} \sum_{i=1}^m \left( (y_i - \beta_0)^2 + \beta_1^2 M^2 - 2\beta_1 (y_i - \beta_0) M^1 \right) \quad (9) \\
 &\quad - \frac{1}{2\sigma^2} \sum_{i=m+1}^n (y_i - \beta_0 - \beta_1 x_i)^2 - \frac{1}{2\delta^2} \sum_{i=m+1}^n (x_i - \alpha)^2 \\
 &\quad - \frac{1}{2\delta^2} \sum_{i=1}^m (M^2 - 2\alpha M^1 - \alpha^2)
 \end{aligned}$$

where  $M^1, M^2$  are given in Equations (7) and (8).

2) *M-step*.

The M-step maximizes  $Q(\theta, \theta^*)$  calculated in the E-step. Solving  $\frac{\partial Q(\theta, \theta^*)}{\partial \theta} = 0$ , we get the following results. The updated estimates of  $\beta'$  is just the OLS solution to the model, *i.e.*,  $\beta' = (X^T X)^{-1} (X^T Y)$

$$\begin{bmatrix} \beta'_0 \\ \beta'_1 \end{bmatrix} = \begin{bmatrix} n & \sum_{i=1}^n \widetilde{x}_i^* \\ \sum_{i=1}^n \widetilde{x}_i^* & \sum_{i=1}^n \widetilde{x}_i^{2*} \end{bmatrix}^{-1} \begin{bmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n \widetilde{x}_i^* * y_i \end{bmatrix}$$

where  $\widetilde{X}^* = (X_{obs}, M^{1*}) \in \mathbb{R}^n$  and  $\widetilde{X}^{2*} = (X_{obs}^2, M^{2*}) \in \mathbb{R}^n$  are estimated completed predictor under current estimated parameter  $\theta^*$ . Similarly, the other updated parameters will be:

$$\begin{aligned}
 \sigma^{2'} &= \frac{1}{n} \sum_{i=1}^n \left( (y_i - \beta'_0) - 2\beta'_1 y_i - \beta'_0 \right) \widetilde{x}_i^* + (\beta'_1)^2 \widetilde{x}_i^{2*} \\
 \alpha' &= \frac{1}{n} \sum_{i=1}^n \widetilde{x}_i^* \\
 \delta^{2'} &= \frac{1}{n} \sum_{i=1}^n \left( \widetilde{x}_i^{2*} - 2\alpha' \widetilde{x}_i^* + \alpha'^2 \right)
 \end{aligned}$$

### 2.4. Convergence of EM Algorithm

We will discuss the convergence of EM algorithm in a more general setting. Suppose we have dataset of  $m$  independent examples and want to fit a parametric model  $p(x, z)$  to the dataset, then the log likelihood function is:

$$\begin{aligned}
 l(\theta) &= \sum_{i=1}^m \log p(x^{(i)}; \theta) \\
 &= \sum_{i=1}^m \log p(x^{(i)}, z; \theta)
 \end{aligned}$$

where  $z$  are the latent random variables. Explicitly finding the maximum likelihood estimate of the parameter  $\theta$  is quite hard, but if  $z^{(i)}$  is observed, the estimation would be easy. Let  $Q_i(z) \geq 0, \sum_z Q_i(z) = 1$ . Since  $f(x) = \log(x)$  is a concave function and by Jensen's Inequality, we get the lower-bound of  $l(\theta)$ :

$$\sum_i \log p(x^{(i)}; \theta) = \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log p(x^{(i)}, z^{(i)}; \theta) \quad (10)$$

$$= \sum_i \log \sum_{z^{(i)}} Q_i(z^{(i)}) \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \tag{11}$$

$$\geq \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \tag{12}$$

Note that this inequality holds for any distribution  $Q_i$ , it gives the lower bound on  $l(\theta)$ . Later we will show that the  $l(\theta)$  increases monotonically with successive iterations of EM if the lower-bound is tight at  $\theta$ . We know that the Jensen's Inequality holds with equality if the random variables are constant. So, it suffices to satisfy:

$$\frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} = c$$

i.e.  $Q_i(z^{(i)}) \propto p(x^{(i)}, z^{(i)}; \theta)$

where  $c$  is constant and does not depend on  $z^{(i)}$ . Under this assumption, since  $\sum_z Q_i(z) = 1$ , we get:

$$Q_i(z^{(i)}) = \frac{p(x^{(i)}, z^{(i)}; \theta)}{\sum_z p(x^{(i)}, z; \theta)} = \frac{p(x^{(i)}, z^{(i)}; \theta)}{p(x^{(i)}, z; \theta)} = p(z^{(i)} | x^{(i)}; \theta)$$

So  $Q_i$ 's is just the posterior distribution of the  $z^{(i)}$ 's given  $x^{(i)}$  and the setting of  $\theta$ . Now we introduce the iteration of EM algorithm:

While  $\|\theta^{(t)} - \theta^{(t-1)}\| > \epsilon$  do

*E-step*

Compute  $Q_i^{(t)}(z^{(i)}) := p(z^{(i)} | x^{(i)}; \theta^{(t)})$

*M-step*

Compute

$$\theta^{(t+1)} := \arg \max_{\theta} \sum_i \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i^{(t)}(z^{(i)})}$$

*End*

Consider that:

$$l(\theta^{(t+1)}) \geq \sum_i \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta^{(t+1)})}{Q_i^{(t)}(z^{(i)})} \tag{13}$$

$$\geq \sum_i \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta^{(t)})}{Q_i^{(t)}(z^{(i)})} \tag{14}$$

$$= l(\theta^{(t)}) \tag{15}$$

The first inequality holds because of (13). By the definition of  $\theta^{(t+1)}$ , (15) is obvious. The last equality holds because lower-bound in (13) is tight at  $\theta = \theta^{(t)}$  under our previous assumption. So, the sequence  $\{l(\theta^{(t)})\}_t$  is both upper

bounded (by 0) and increasing. Hence, in EM algorithm the log likelihood converges monotonically.

### 3. Application

This section covers the model-based simulation research application of the EM algorithm with least square and resilient least square regression on compositional data. We also analyze the resilience and efficiency of the EM approach and compare its output to two other commonly used imputation techniques, k-Nearest Neighbor (k-NN) and mean imputation ( $\bar{x}$ ), to address missing data.

#### 3.1. Data Description

Compositional data is a unique kind of non-negative data that contains the pertinent information not in the actual data values but rather in the ratios between the variables. An observation  $x = (x_1, \dots, x_D)$  is a D-part composition if, and only if,  $x_i > 0$ ,  $i = 1, \dots, D$ , and according to Aitchison [2], the ratios between the components include all the important information.

$$S^D = \{[x_1, \dots, x_D] : x_i > 0 (i = 1, \dots, D), x_1 + \dots + x_D = 1\} \text{ and}$$

$$(x_1, \dots, x_D) = \frac{(w_1, \dots, w_D)}{w_1 + \dots + w_D} \text{ where } (w) \text{ denote } = \text{ total weight and } (w_1, \dots, w_D)$$

are the component weights. According to Aitchison [2], compositional data is not directly represented in Euclidean space. The Aitchison distance  $d_A$  is a suitable way to measure the D-part composition known as (simplex) [3]. According to Egozcue *et al.* [4], the isometric log-ratio (*ilr*) is used to convert the D-dimensional simplex into the real space  $\mathbb{R}^{D-1}$ . With this transformation, the Aitchison distance can be expressed as

$d_A(x, y) = d_E(ilr(x), ilr(y))$ , where  $d_E$  denotes the Euclidean distance. The data in this simulation is generated by a normal distribution on the simplex, denoted by  $\mathbb{N}_s^D(\mu, \Sigma)$  (Mateu-Figueras, Pawlowsky-Glahn, and Egozcue) [5]. We generated 10,000 realizations of a random variable [6]  $X \sim \mathbb{N}_s^4(\mu, \Sigma)$  with  $\mu = (0, 2, 3)^T$  and  $\Sigma = \left( (1, -0.5, 1.4)^T, (-0.5, 1, -0.6)^T, (1.4, -0.6, 2)^T \right)$  [7].

#### 3.2. Experimental Design

- To assess the effectiveness of the EM algorithm, we look at the results using least squares (LS) and its robust version (RLS) across a range of missing data rates (contamination levels between 5% and 10%) and outlier rates (1%, 3%, 5%, and 10%) expressed in terms of their Aitchison distance ( $d_A$ ) [8].
- We also look at the EM algorithm's output in terms of the covariances of various rates of outliers (1%, 3%, 5%, and 10%) and missing data (contamination levels ranging from 5% to 10%).

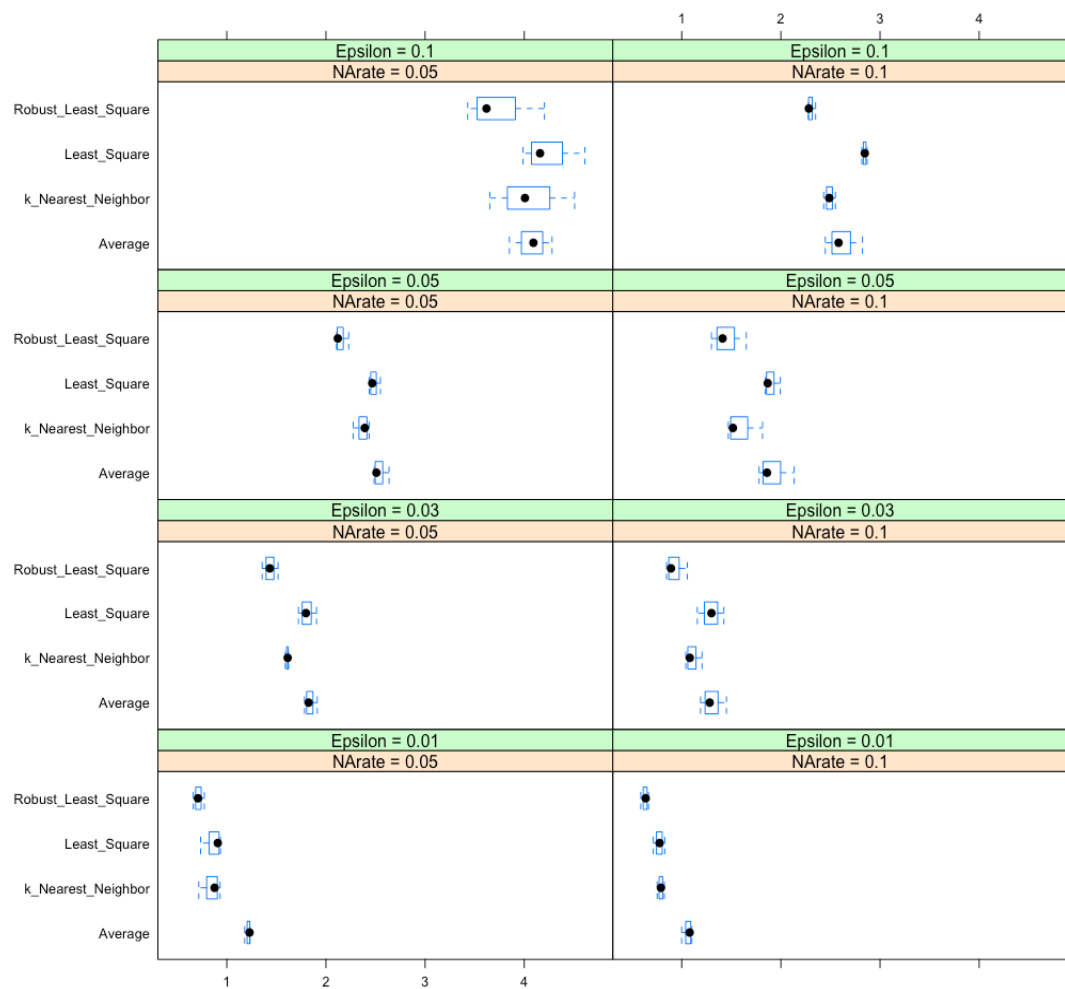
#### 3.3. Results and Analysis

The detailed results for all experiments are discussed in this section (Table 1, Table 2, Figure 1, Figure 2).

**Table 1.** Performance metrics for different imputation methods in terms of distance.

S/N	Epsilon <sup>1</sup>	NArate <sup>2</sup>	xMean <sup>3</sup>	kNN <sup>4</sup>	LS <sup>5</sup>	RLS <sup>6</sup>
1	0.01	0.05	1.208117	0.8815717	0.9277800	<b>0.7544389</b>
2	0.03	0.05	1.869092	1.6451996	1.9300590	<b>1.5242680</b>
3	0.05	0.05	2.613898	2.2722186	2.5338586	<b>2.1552349</b>
4	0.10	0.05	4.044281	3.8628814	4.1707580	<b>3.6228353</b>
5	0.01	0.10	1.087026	0.7727159	0.8874163	<b>0.6802822</b>
6	0.03	0.10	1.419570	1.2479738	1.3840970	<b>1.0872886</b>
7	0.05	0.10	1.798377	1.5482930	1.7860411	<b>1.4121292</b>
8	0.10	0.10	2.693001	2.4922940	2.8084783	<b>2.2551779</b>

<sup>1</sup>Epsilon denote rate of outlier at (1%, 3%, 5% and 10%); <sup>2</sup>NArate denote missing rate (contamination level at 5% and 10%); <sup>3</sup>xMean denote arithmetic mean imputation method; <sup>4</sup>kNN denote k-Nearest Neighbor imputation method; <sup>5</sup>LS denote least square regression version of EM algorithm; <sup>6</sup>RLS denote robust least square regression version of EM algorithm; \*Bold numbers indicate the best performing imputation method for a given epsilon and missing rate.



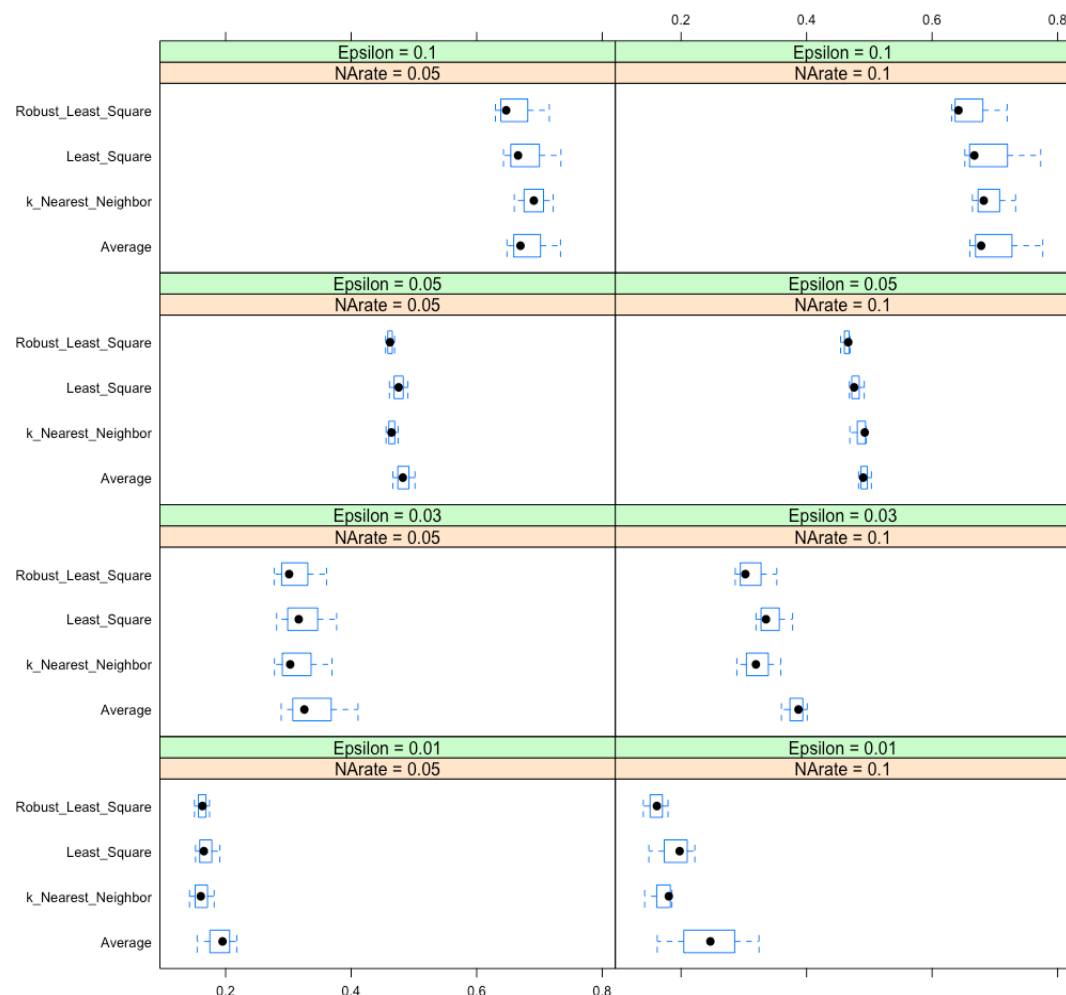
**Figure 1.** Performance for different imputation methods in terms of distance.



**Table 2.** Performance metrics for different imputation methods in terms of Covariance.

S/N	Epsilon <sup>1</sup>	NArate <sup>2</sup>	xMean <sup>3</sup>	kNN <sup>4</sup>	LS <sup>5</sup>	RLS <sup>6</sup>
1	0.01	0.05	0.1982435	0.1686442	0.1756952	<b>0.1668150</b>
2	0.03	0.05	0.4107212	0.4087216	0.3966203	<b>0.3836261</b>
3	0.05	0.05	0.5342290	0.5234262	0.5219798	<b>0.5126078</b>
4	0.10	0.05	0.7064511	0.7005467	0.7037723	<b>0.6945066</b>
5	0.01	0.10	0.2274753	0.1765070	0.1968090	<b>0.1695136</b>
6	0.03	0.10	0.4497979	0.3934609	0.4346386	<b>0.3849345</b>
7	0.05	0.10	0.5563971	0.5472681	0.5541094	<b>0.5150961</b>
8	0.10	0.10	0.7163915	0.6997279	0.7142401	<b>0.6907288</b>

<sup>1</sup>Epsilon denote rate of outlier at (1%, 3%, 5% and 10%); <sup>2</sup>NArate denote missing rate (contamination level at 5% and 10%); <sup>3</sup>xMean denote arithmetic mean imputation method; <sup>4</sup>kNN denote k-Nearest Neighbor imputation method; <sup>5</sup>LS denote least square regression version of EM algorithm; <sup>6</sup>RLS denote robust least square regression version of EM algorithm; \*Bold numbers indicate the best performing imputation method for a given epsilon and missing rate.



**Figure 2.** Performance for different imputation methods in terms of Covariance.

## 4. Conclusions

In this study, it is postulated that the missing rate follows a random pattern. Consequently, both the least square and robust least square EM algorithm approaches are employed to address the issue of missing compositional data. The simulated compositional synthetic data was utilized as the practical dataset in our study. It has been shown that when the occurrence of missing and outlier data decreases, the Aitchison distance approaches zero. Furthermore, the robust least square variant of the EM method yields the smallest values in this regard making it the best performing algorithm. Distance values close to zero represent the most efficacious imputation strategies.

We also compare the Expectation-Maximization (EM) method to arithmetic mean imputation and k-Nearest Neighbor and find that the robust EM version is the best of the bunch when it comes to co-variances. In this study, we provided evidence to support the effectiveness of all strategies when applied under conditions of low missing rates, specifically those below 5%. The EM algorithm demonstrates improved performance as the rate of missing data increases. However, when the incidence of missing data exceeds 10%, most imputation methods become inadequate in generating a dependable approximation, hence exacerbating the difficulty in recovering the data.

## Conflicts of Interest

The authors declare no conflicts of interest regarding the publication of this paper.

## References

- [1] Aitchison, J. (2002) Simplicial Inference. In: Viana, M.A.G. and Richards, D.S.P., Eds., *Contemporary Mathematics Series, Vol. 287: Algebraic Methods in Statistics and Probability*; American Mathematical Society, Providence, 1-22.
- [2] Aitchison, J. (1986) *The Statistical Analysis of Compositional Data*. Chapman & Hall, London.
- [3] Aitchison, J. (1989) Measures of Location of Compositional Data Sets. *Mathematical Geology*, **21**, 787-790. <https://doi.org/10.1007/BF00893322>
- [4] Martín-Fernández, J.A., Egozcue, J.J., Olea, R.A. *et al.* (2021) Units Recovery Methods in Compositional Data Analysis. *Natural Resources Research*, **30**, 3045-3058. <https://doi.org/10.1007/s11053-020-09659-7>
- [5] Pawlowsky, V., Olea, R.A. and Davis, J.C. (1995) Estimation of Regionalized Compositions: A Comparison of Three Methods. *Mathematical Geosciences*, **27**, 105-127. <https://doi.org/10.1007/BF02083570>
- [6] Weltje, G.J. (1997) End-Member Modeling of Compositional Data: Numerical-Statistical Algorithms for Solving the Explicit Mixing Problem. *Mathematical Geosciences*, **29**, 503-549. <https://doi.org/10.1007/BF02775085>
- [7] Rehder, U. and Zier, S. (2001) Comment on "Logratio Analysis and Compositional distance by Aitchison *et al.* (2000)". *Journal of Mathematical Geology*, **32**, 741-763.
- [8] Dempster, A.P., Laird, N.M. and Rubin, D.B. (1977) Maximum Likelihood from Incomplete Data via the EM Algorithm. *Journal of the Royal Statistical Society*, **39**, 1-22. <https://doi.org/10.1111/j.2517-6161.1977.tb01600.x>