

**Handling Missing Values in Relational Event History Data using Multiple Imputation:
A Framework in Social Network Research**

Myrthe Prins (6753566)

Mahdi Shafiee Kamalabad

Gerko Vink

Master Applied Data Science, Utrecht University

July 7, 2024

Abstract

Background - Missing data is a problem that is common. It affects the accuracy and introduces biases in social network analysis, which can have a significant effect on the interpretation of findings. Relational event history (REH) data, a type of social network data, is becoming increasingly available due to new technological developments and can enhance the understanding of dynamic social networks. However, research on handling missing values in social network data is limited and statistical tools for incomplete REH data are underdeveloped. This paper focuses on using multiple imputation to handle missing values within REH data.

Methods – Relational event history model analysis is first performed on the fully observed dataset to produce true estimates. Next, a simulation study is conducted to introduce missingness to this fully observed dataset, assuming missing completely at random (MCAR) and right-tailed missing at random (MAR). After multiple imputation, the relational event model is applied on the simulations and the results are compared to the analysis of the fully observed dataset.

Results – The results of the relational event model of the simulations and the true estimates show inconsistency in the significance of the results. The simulations generally have a low bias, good coverage rate and low average width. A higher proportion of missingness resulted in a decrease in the performance. Multiple imputation thus produces unbiased inferences under the MCAR and MAR mechanism, however unexpected significant results are found.

Conclusion – This study provides insights into the use of multiple imputation for producing valid inferences when applied on REH data. It shows that under the assumption of MCAR and MAR, multiple imputation can be a valid method for missing data in REH data when the percentage of missingness is not too high. Further research is needed to confirm and expand upon the results obtained in this study.

	3
1. Introduction	4
1.1 Multiple imputation	5
1.2 Missingness mechanisms	6
1.3 Social Networks	8
1.4 Relational Event History data	9
1.5 Relational Event Model	10
1.6 Current research	11
2. Data	12
3. Method	13
3.1 Missing data generation	13
3.2 Multiple imputation	15
3.3 Data analysis	16
4. Results and analysis	18
5. Discussion and conclusion	23
References	25
Appendix A	29
Appendix B	33

1. Introduction

Social networks are defined by a relation among a collection of individuals. Social network analysis seeks to understand social networks by examining the patterns and structures of these relationships. A specific type of social network data is relational event history (REH) data, which can be defined as time-ordered sequences of social interactions between a set of individuals or entities (Butts & Marcum, 2017, Meijerink-Bosman et al., 2022). This data has, unlike panel data, a high resolution precision and records all relational events at an exact moment in time. This makes REH data ideal for social researchers to deeply investigate social phenomena (Back, 2021).

Missing data is a problem that is common across various research domains, affecting the accuracy of estimates and potentially introducing biases in parameter estimates. Thus, it can have a significant effect on the interpretation of findings. Most statistical analyses require complete data (Schouten et al., 2018). As a result, the presence of missing data may result in potentially wrong conclusions drawn from the data, which can lead to far-reaching consequences. Complete case analysis, a simple approach to handling missing data, involves removing the observations with missing values completely. However, this method causes loss of information, a decrease in statistical power and may introduce bias (Schafer & Graham, 2002). Multiple imputation (MI) offers a solution to this problem. MI is a method that creates multiple complete versions of the data by replacing missing values by plausible data values based on the patterns and relationships found in the observed data. By creating multiple complete versions, MI quantifies uncertainty in estimating missing values and is focused on preserving information, rather than throwing it away. It therefore minimizes the risk of drawing incorrect conclusions (Vink & Van Buuren, 2014). Understanding the appropriate techniques for addressing missing data according to specific circumstances is crucial.

The impact of missing data is larger when the data has a complex structure. Network data, which is highly structured, is particularly affected by missingness as the network structure should be preserved (Borgatti et al., 2006; Smith et al., 2017). Missing values in REH data poses significant challenges for social researchers to investigate complex social and behavioral phenomena as it can lead to invalid statistical inferences. However, research on the influence and how to handle missing values in social network data is limited (Huisman & Krause, 2017). Statistical tools for REH data are also currently almost underdeveloped. This paper focuses on using MI to address missing values within relational event history (REH) data and thereby enhances the understanding of the evolution of social relations in continuous time.

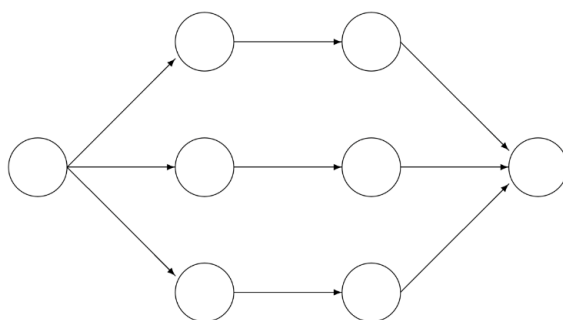
1.1 Multiple imputation

Multiple imputation is a method used to handle missing data by estimating and replacing each missing value multiple times. Each missing value is imputed $m > 2$ times, leading to m completed datasets. The m completed datasets are then analyzed independently and pooled using Rubin's rules for combining estimates and standard errors (Rubin, 1987, pp. 76; Schafer & Graham, 2002). Rubin's rules take into account both within- and between-imputation variability (Van Buuren, 2018). MI differs from single-value imputation methods by imputing missing values multiple times, which creates multiple complete datasets. Single-value imputation methods estimate what each missing value might have been and replace it once with a single value. The multiple complete datasets generated with multiple imputation are analyzed separately. The results of the analyses are combined and multiple imputation thus accounts for the uncertainty in the imputation process, avoiding false precision that can occur with a single imputation method. Replacing a missing value with a single value in single-imputation may suggest an unrealistic accuracy and certainty. MI provides accurate estimates for the metrics of interest. It also minimizes the risk of drawing false-positive or false-negative conclusions (Li et al., 2015).

Multiple imputation consists of two stages. Firstly, imputations (replacement values) for missing values are generated. This results in many datasets with different replaced missing values, as is shown in the imputed data part in Figure 1. The imputations are generated based on statistical characteristics of the data. Secondly, the imputed datasets are analyzed (Analysis results in Figure 1) and the results of these analyses are combined (Li et al., 2015). This is the pooled result (Figure 1).

Figure 1

Steps in multiple imputation: imputation, analysis and pooling



Incomplete data Imputed data Analysis results Pooled result

Note: Scheme of main steps in multiple imputation. From *Flexible Imputation of Missing Data*. (2nd edition., p. 19), by S. van Buuren, 2018, CRC Press. Copyright by 2018 by Taylor & Francis Group, LLC.

It is often assumed that data when using MI, is sampled from infinite populations. However, if sampling from a finite population, the standard pooling rules may overestimate the variance of the estimates. This leads to a loss of statistical efficiency as confidence intervals are wider than necessary. In the case of pooling multiple imputations for finite populations, simplified pooling rules that exclude sampling variance and only account for the variation caused by the mechanism that created the missing data, need to be used (Vink & Van Buuren, 2014). In accordance with the guidelines by Vink & Van Buuren (2014), these simplified pooling rules were applied in this study to combine multiple imputations for finite populations.

1.2 Missingness mechanisms

Solving a missing data problem is challenging. Although there are techniques such as multiple imputation (Rubin, 1987; Little & Rubin, 2002) which are proven to be effective and intuitive, it is important to think about the characteristics of the missingness. The degree to which the observed and unobserved data are connected, may be of great influence on the validity of the imputation method. Inclusion of a variable that correlates either with the incomplete variable or with the missing values improves parameter estimates (Collins et al., 2001). This is the reason why predictor variables are often included in imputation methods.

Formulating beliefs about the extent to which the observed data also applies for the missing parts of the data is essential for handling missing data. We distinguish three different missingness mechanisms: 1) Missing Completely At Random (MCAR); the probability of being missing is the same for all cases. 2) Missing At Random (MAR); the probability of being missing is related to the observed data. 3) Missing Not At Random (MNAR); the probability of being missing varies for reasons unrelated to the observed data (Rubin, 1987; Van Buuren, 2018). In the current literature, MCAR is also referred to as Not Data Dependent (NDD), MAR as Seen Data Dependent and MNAR as Unseen Data Dependent (UDD) (Hand, 2020). The mechanisms can be further explained using a data matrix Y with y_{ij} either observed or missing, where Y_{obs} is the observed data and Y_{mis} is the missing data. Matrix R is considered the response indicator with $R_{ij} = 1$ if y_{ij} is missing and $R_{ij} = 0$ if y_{ij} is observed. Ψ are fixed parameters of the probability model (Van Buuren, 2018).

Missing Completely At Random (MCAR)

When the data is missing completely at random (MCAR), the probability of a variable being missing is independent from the observed and unobserved data. The missingness is thus not related to the data. This can be represented by the formula:

$$Pr(R=1|Y_{obs}, Y_{mis}, \psi) = Pr(R=1|\psi)$$

The missing values are solely induced by ψ and independent from the observed and unobserved data. The observed data and the missing data are thus exact representations of the true data model. When the data are MCAR, the remaining data can be considered a simple random sample of the full dataset (Mack et al., 2018). MCAR missingness is also called *ignorable* missingness, because bias is not introduced and power can be restored with modern treatment (Rubin, 1967; Little et al., 2013).

Missing At Random (MAR)

With MAR mechanism, the probability of a value being missing depends on the values of the observed variable but not the unobserved data. This can be represented by the formula:

$$Pr(R=1|Y_{obs}, Y_{mis}, \psi) = Pr(R=1|\psi, Y_{obs})$$

The observed and missing data represent different parts of the population. MAR missingness, like MCAR, is also referred to as *ignorable* missingness, as the bias is recoverable and power can be restored with modern treatment, such as multiple imputation (Rubin, 1976; Little et al., 2013).

Missing Not At Random (MNAR)

With MNAR mechanism, the probability of a value being missing depends on unobserved information. The missingness is related to events or factors which are unknown. This can be noted as:

$$Pr(R=1|Y_{obs}, Y_{mis}, \psi) = Pr(R=1|\psi, Y_{obs}, Y_{mis})$$

The observed data alone is not enough to infer about the population. The missing data are called *nonignorable* (Rubin, 1967). The observed and unobserved data represent different and unique parts of the true data.

When the variables in a dataset show low correlations, the identification between MAR and MCAR missingness for the observed data may become difficult. This is because low correlations provide less information about the relationships between variables, which makes it difficult to determine whether the missingness is completely at random (MCAR) or related to other observed variables (MAR). MI with MAR mechanisms would then primarily limit statistical power and increase variance without necessarily reducing the bias. This is also true for assuming MNAR missingness when data is highly correlated, as strong correlations can make it difficult to distinguish whether the missingness is related to the unobserved variables or only related to the observed variables. Therefore, it is important to consider which mechanism to assume based on

the observed data structure (Schouten & Vink, 2018). Imputation methods can yield valid results provided that the missingness mechanism is not MNAR. The percentage of missingness should also not be too high as this can lead to biased imputation models and inaccurate results (Li et al., 2015; Sterne et al., 2009).

1.3 Networks

Network data is highly structured, missing data can therefore have a large impact on this type of data. Networks are a collection of nodes (points) joined together in pairs by edges (lines). There are many examples of systems which can be represented by networks in the physical, biological and social sciences. These sciences all have different types of networks, which can be divided into four broad categories: technological networks, information networks, biological networks and social networks (Newman, 2018). Here, we will focus on social networks.

Social networks

Social networks can be defined as any network in which the nodes (*actors*) represent individuals, such as friends, family members or classmates, and the edges (*ties*) represent the relationships/connection between them, such as friendships or interactions (Newman, 2018). Social networks are based on the representation of social structure in terms of a set of social entities, such as people and organizations, that are connected via relationships (Wasserman & Faust, 1994; Carrington et al., 2005). Social networks are not always static, the relationships can be dynamic, they can change over time. However, some social networks exhibit stability over time. Dynamic social networks require different models than static networks (e.g. Snijders, 2001; Krivitsky & Handcock, 2014). Social networks have a high flexibility, many different definitions of an edge are possible, and they can thus serve as a good representation of different social phenomena (Butts, 2009; Newman, 2018). Most current models view relationships as evolving over time, discrete or continuous. The changes in the relationships are driven by mechanisms, such as reciprocity (the tendency for directed ties to be reciprocated) or transitivity (the tendency for nodes connected to the same node to also be connected to each other). The presence and strength of these mechanisms can be estimated from the intertemporal network data. These models thus allow for better understanding of social networks. Many researchers are interested in the network perspective because it provides a framework to understand patterns in relationships among interacting units within a social environment (Wasserman & Faust, 1994).

Social network analysis can provide insight into the underlying relationships between individuals, which can reveal patterns (in e.g. interactions, communication, relationships) that cannot be detected from the individual observations alone. Social network analysis is based on the assumption that relationships among interacting individuals are important (Wasserman & Faust, 1994; Serrat, 2017). It shows the formal and informal relationships between individuals and can be used to understand what facilitates or impedes the existence of ties between edges. Social network analysis has gotten much more interest in social and behavioral sciences as the availability and technical tools of social network data are increasing (Carrington et al., 2005).

1.4 Relational Event History data

Relational Event History (REH) data is a type of social network data, which describes a time-ordered series of interactions between actors in a network. These interactions are also known as relational events. Minimally, the relational events contains information about the actors that are involved in the event and the time of the event (Meijerink-Bosman et al., 2023). It captures an action initiated by one entity and directed toward another entity within its environment at a specific point in time (Butts & Marcum, 2017).

Table 1 provides an example of a REH dataset, showing the relational events per row. The relational event consists of a sender, receiver and time, in which the sender is the sender of the action, the receiver the receiver of the action, and the time shows the time at when the interaction took place (see Table 1). Both the sender and receiver can consist of humans, animals, objects or a combination of multiple type of actors. Actions can consist of a variety of relationships between the actors. Multiple of these events combined and ordered by time given a time window, result in REH data (Marcum & Butts, 2015).

Table 1

Example of Relational Event History (REH) dataset

Time	Sender ID	Receiver ID
11849.2	18	2
11854.2	2	18
11885.2	18	2
11890.2	2	18
12232.2	2	17

Note. Adapted from *Apollo 13* dataset. Time is in seconds from onset of the Apollo 13 mission. Sender ID and Receiver ID represents respectively the sender and receiver of the message, indicated with a number.

The interactions between social entities in REH data are discrete instances. This is in contrast to the conventional social network setting in which the ties are temporally extensive such as friendships or family members. REH data focuses on individual interactions that occur at specific time moments, as in conventional social network data the relationships are presented as ongoing over time (Butts, 2008).

REH data is becoming increasingly available due to the development of technology and has the potential to greatly contribute to the understanding of dynamic social networks. REH data also has a high precision, which makes it particularly useful for analyzing social networks. REH data is distinct from panel data in the sense that the ties in REH data are short lasting, there are no unobserved tie changes and relational events occur in exact moments in time. These differences in combination with the rapid increase of available REH data, its high precision and the potential to greatly contribute to the understanding of dynamic social networks, make REH data valuable for understanding social networks.

1.5 Relational Event Model

Analysis of REH data can help researchers answer complicated research questions, such as at the most basic level “what drives what happens next?” in a complex sequence of interdependent events (Marcum & Butts, 2015). In more detail it can give insight into who is interacting with whom at what specific time, predict future interactions, reveal how interaction dynamics change over time, assess the impact of past events on future interactions and identify what drives interactions. Therefore, this type of data is becoming more and more popular for the analysis of relational dynamics. Relational event dynamics are fundamentally about sequential relational structures, which differs from the conventional social network analysis as the primary interest thereof is the simultaneous relational structure.

REH data is characterized by the inherent dependency between nodes and edges (Meijerink-Bosman et al., 2023). Consequently, this type of data is not handled very well by traditional statistical methods such as linear regression because this does not take into account the temporal sequencing and interdependence of events. Specialized tools are therefore needed for analysis. Relational Event Model (REM) is a gold standard statistical model that can take this into account (Butts, 2008). The REM is built to analyze continuous, detailed, social interaction data, such as is found in REH data (Meijerink et al., 2023). It is used to understand

communication/interaction structures based on observed social interactions in real-time (Butts, 2008).

REM can be used to examine the frequency and time to activation among relational events. The probability of future interactions can be determined by the event rate (λ). This can be represented as follows:

$$\log \lambda(s, r, t) = \sum \beta_p x_p(s, r, t)$$

where the event rate between sender “ s ” and receiver “ r ” at time point “ t ”, $\lambda(s, r, t)$, is modelled as a loglinear function of statistics (both endogenous and exogenous statistics can be included), which shows the propensity of an event to occur. In the loglinear function, $X_p(s, r, t)$ refers to the p -th statistic for the actor pair (s, r) at time “ t ” and β_p refers to the model parameter related to the statistic X_p . The endogenous statistics contain the internal information until a given time point (e.g. past interactions), while the exogenous statistics contain external information (e.g. attributes, age, sex). With the REM, endogenous and exogenous factors can be investigated that predict the probability of subsequent events happening (Butts, 2008). The risk set consists of all possible interactions at time “ t ” and usually consist of $N(N-1)$ events, with N being the total number of actors in the network. It is constructed to calculate the event rate for all possible events at specific time point and for predicting the future events. The event rate, $\lambda(s, r, t)$, thus represents the rate of occurrence at time “ t ” for sender-receiver pair (s, r). The outcome of the REM shows the extent to which the specified statistics affect social interaction behavior in the network (Meijerink-Bosman et al., 2023).

1.6 Current research

REH data is important as it is becoming increasingly available, offers a high resolution and captures detailed histories of events. This makes this type of data valuable for studying complex social systems and addressing crucial research questions. REM is used to analyze REH data. We already established that missing data can cause serious problems. This is especially the case for complex structures, such as in REH data, because the network structure should be preserved (Borgatti et al., 2006; Smith et al., 2017). REM is very sensitive to missingness. Failing to address missing data appropriately or using naïve approaches can lead to potentially wrong conclusions in the analysis. Currently, research on the impact and how to handle missing values in social network data is limited (Huisman, 2009; Huisman & Krause, 2017). Statistical tools for addressing missingness in REH data are also currently underdeveloped, although REH data is

everywhere. Therefore, developing valid methodologies to address missingness in REH data is essential.

This paper gives more insight into the use of MI as an approach to missingness in REH data. It aims to assess the effectiveness of MI to produce accurate estimates of missing data and evaluates its impact on the REM estimates. The research question of this paper is: How effective is multiple imputation in handling missing data in Relational Event History data to produce valid inferences? The remainder of this paper has the following structure. First, I will provide information about the data used in the study. I will then outline the methods used in the study in the third section and present the results and analysis in the fourth section. In the fifth section, the conclusion and limitations of these findings will be discussed.

2. Data

The REH data used in this study is part of the Apollo 13 dataset. Apollo was a program by NASA in which people aimed to travel to the moon for the first time, Apollo 13 was the seventh mission in this program. The mission ended early due to an explosion in the oxygen tank. The communication between the flight and ground crew ensured a safe return on the ground. This unusual occurrence resulted in a well-documented dataset capturing the crisis communication between the flight and ground crew. The real-time playback of the events following the incident is available on the Apollo 13 Real-time website (*Apollo 13 Real-time*, n.d.). The full Apollo 13 dataset after the incident is also publicly accessible on GitHub (Tseng, 2017).

The part of the Apollo 13 dataset used, consists in total of 38982 rows and three columns: time, sender and receiver. Each row represents a single, directed communication event, which is the relational event. The sender column represents the actors that initiated the communication at a corresponding time point. The receiver column represents the actors that were the target of this communication initiated by the sender. The time column consists of unique exact time points of the communication initiated. The actors in the sender and receiver column are indicated with numbers, with numbers 1-16 representing the ground crew and 17-19 representing the flight crew. The subset of the Apollo 13 dataset used, includes 16 unique actors. Numbers 14 to 16 were not present in this particular subset. The dataset is fully observed, no data preparation is needed. The Apollo 13 dataset gives an ideal illustration for the use of MI on handling missing values in REH data.

3. Method

To assess the effectiveness of MI on handling missing values in REH data, the criteria for evaluation must be determined. First, missing values are induced in the complete dataset (amputation) and the missingness mechanism and proportion of missingness are defined. Subsequently MI is performed, generating multiple complete datasets, each of which is analyzed using REM. The results of these analyses are combined (pooled). Additionally, the REM is applied to the fully complete dataset. The accuracy of estimates of MI and its impact on the REM estimates is examined. Below, each step will be discussed in more detail (see Appendix B for full R code).

3.1 Missing data generation

The generation of missing values in a complete dataset is called amputation. The Apollo 13 dataset is a fully observed dataset. To evaluate MI on this dataset, missing values have to be induced. The package MICE (Multivariate Imputation via Chained Equations) (Van Buuren & Groothuis-Oudshoorn, 2011) is used in R (R Core Team, 2023). This implements a method to handle missing data. The method is based on Fully Conditional Specification (FCS), which imputes each incomplete variable using separate models, such as regression models and predictive mean matching. MICE works by iteratively imputing missing values using predictive models depending on the observed data (Van Buuren & Groothuis-Oudshoorn, 2011).

Based on “Strategies for simulated missingness” (Vink, 2022), missingness was simulated with the “*ampute()*” function (Schouten et al., 2018) from the package MICE. Model-based finite populations was used, where a single finite observed set is taken as the comparative truth and missingness is induced via simulation. Monte Carlo simulations refer to simulations in which many random values are sampled from a posterior distribution, which allows for uncertainty, variability and distribution in the estimates. The induced missingness provides the necessary variation for a Monte Carlo simulation (Hammersley & Handscomb, 1964). The sampling variance from the evaluations of the imputation performance can be eliminated as the sampling variance becomes irrelevant when dealing with a finite population (Vink and Buuren, 2014). Once every element in the population has been included in the sample, no variance is left to estimate. The noise induced by the sampling mechanism will not be taken into account, as this is not the topic of interest of this research.

In this study the proportions of missing is set to 10%, 30% and 50% to simulate real-world scenarios in which data might be incomplete with different percentages of missingness. The percentage of missingness indicates the percentage of data rows that will have missing

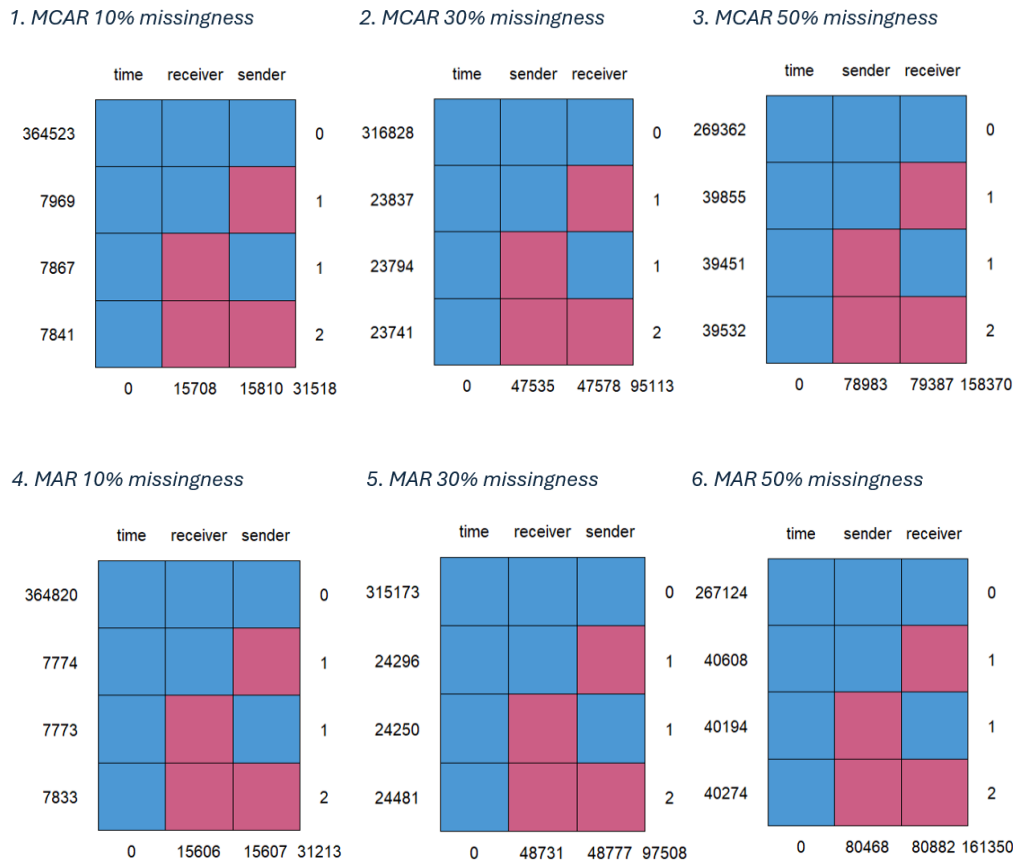
values in either the sender or receiver column or both. Missingness was simulated a hundred times to ensure the validity of the results and to minimize the Monte Carlo error. The missing data generation process was conducted two times: once assuming MCAR as the missingness mechanism and once assuming right-tailed MAR. In real-world datasets, MCAR is rarely the case, because it is often unrealistic to assume that the probability of missing values is completely independent of other observed variables (Van Buuren, 2018). The likelihood of data being missing often depends on the values of the observed dataset, therefore this study also observes multiple imputation assuming MAR. A right-tailed MAR mechanism indicates that the probability of the data being missing increases with higher values. The Apollo dataset contains categorical variables for sender and receiver. Because sender and receiver are not continuous variables, right-tailed MAR could be interpreted in the context of having implicit ordered IDs within the categorical variables. The categorical variables can be thought of as having an implicit order, higher IDs could correlate with more influential or active actors which makes them having a higher chance of being missing. By assuming a right-tailed MAR mechanism, the performance of multiple imputation under conditions that are more complex and realistic than MCAR can be tested.

The `ampute()` function in MICE creates pattern-based missingness, where the patterns govern the relation between missingness and observed data. The weights for `ampute()` patterns are not specified, resulting in all patterns occurring with equal probability. To prevent a combination of actors from going completely missing in the dataset, 1500 observations were conserved from the dataset and not amputed. This way, all dyads (pairs of variables) are still covered, thereby preserving necessary dyads for comparison with the true data after imputation.

The missing data pattern considered for amputation are all possible combinations of missingness involving sender (actor 1) and receiver (actor 2) per missingness mechanism and percentage (Figure 2). Each row in the missing data pattern represents a unique missingness pattern. The numbers on the left indicate the frequency of each pattern for all simulations per missingness mechanism and proportion combined. Cells highlighted in red indicate missing values, while blue cells indicate observed values. The numbers at the bottom illustrate the count of missing values in each column.

Figure 2

Missing data pattern of all simulations combined per missingness percentage and missingness mechanism



Note. The rows represents the missingness patterns. The numbers on the left indicate the frequency of each missing data pattern. Red cells indicate missing values in those specific rows. Blue cells indicate the observed values. The numbers at the bottom indicate the count of missing values in each column.

3.2 Multiple imputation

The variables sender, receiver and time are used as predictors for the missing values. Other variables (exogenous statistics) that can be derived from the variables present in the dataset, such as whether the sender and receiver are in the same location, are not included as predictors as it is not realistic to assume that it is known whether sender and receiver were in the same location if one of them was missing. However, using this information could improve the performance. In the past, it was thought that using a lower number of imputations, around 3-5 imputations, was sufficient to obtain excellent results (Schafer & Olsen, 1998). The number of imputations refers to the number of datasets generated, each containing missing values generated based on the number of iterations. More recent research by Van Buuren (2018) and

Graham et al. (2007) show that a higher number of imputations than previously recommended, would lead to better outcomes. However, considering the computational efficiency, the number of imputations is set to five. The number of iterations is also set to five, as inferential validity is often achieved after five to ten iterations. Proper convergence is generally assumed to be achieved in 5 to 20 iterations (Oberman et al., 2021). The number of iterations refers to the number of iterations through each variable to estimate new values for missing data (Van Buuren, 2018).

A custom version of the predictive mean matching (pmm) method (Van Buuren, 2018), “pmm.conditional” (Vink, 2023) is used as the method for imputation. This method works by first selecting a small set of candidate donors that are the closest to the missing value from the other variables. One randomly drawn donor from the set of candidate donors is taken to replace the missing value. Predictive mean matching avoids creating loops in the data, which makes it possible to use time as a predictor and prevents actors initiating interactions with themselves. The imputed value for the corresponding receiver or sender cannot be the same, this will be ensured using the “pmm.conditional” method. The generation of MCAR missingness can be seen as the situation in which all weight values are zero because there is no relationship between the data values and the missing data. Since the probability of values being missing under the MAR mechanism by definition depends on the value of the observed variables, only the weights of the variables that will be imputed is set to zero (Schouten et al., 2018).

The convergence and plausibility of the imputation is checked. The convergence is checked by the convergence plot for MAR and MCAR simulations. The convergence plot shows the change in imputation estimates as the number of imputations increases and can thus tell whether the algorithm has stabilized, or further iterations are needed for reliable imputations. Plausibility is checked by the density plots for MAR and MCAR simulations. The density plots show the distribution of the imputed values and the skewness of the imputed data distribution. The quantile-quantile also shows the distribution of the imputed values (Nguyen et al., 2017). The Kolmogorov-Smirnov test (Kolmogorov & Smirnov, 1933) is used to compare the distribution of the imputed values to a theoretical distribution in the quantile-quantile plot.

3.3 Data analysis

The statistics (effects) to be included in the REM are defined. The remify() function (Arena et al, 2023) and remstats() (Meijerink-Bosman et al., 2023) are used to calculate the statistics on the dataset. The model includes three endogenous statistics, as used by Shafiee Kamalabad

et al. (2023): reciprocity, indegree sender and outdegree receiver. Additionally, one exogenous statistic is included: same location. The definitions of the statistics are:

Reciprocity

Reciprocity measures the tendency for directed ties to be reciprocated, which reflects the propensity of person A to initiate an event towards person B as a function of the volume of past events A received from B. Thus, the rate of a relational event occurring from A to B is positively affected by the volume of prior instances of a relational event from B to A (Leenders et al., 2016, Coleman et al., 1990).

Indegree sender

Indegree sender reflects the tendency for senders with a high number of incoming ties to initiate ties. This indicates whether individuals who frequently receive interactions are more or less likely to initiate interactions themselves (Carrington et al., 2005).

Outdegree receiver

Outdegree receiver reflects the tendency for receivers with a high number of outgoing ties to receive ties from others. This indicates whether individuals who frequently initiate interactions are more or less likely to receive interactions themselves (Carrington et al., 2005).

Same location

Same location reflects the tendency of more interactions when the sender and receiver are both at mission control/both in space.

The data is processed to create a dataset for use with the Cox proportional hazard (coxph) function from the survival package (Therneau, 2023). A risk set is created, consisting of every possible interaction at a given time combined with the status of whether that interaction took place in the observed data and whether the actors are at the same location. The status column indicates 1 if interaction took place and 0 if it did not. The same location column indicates 1 if the actors were in the same place and 0 if not. Then, the REM is fitted on all simulations and the results are pooled. The analysis is run over a hundred simulations for the imputed datasets, assuming MCAR and MAR. These results are pooled and averaged for both MCAR and MAR. The imputation method is evaluated using the following measures, as proposed by Van Buuren (2018):

- Raw Bias (RB) and Percent Bias (PB): measures difference between the expected value of the estimate and the true value of the estimate. RB should be close to zero and is defined

as: $RB = E(\hat{\beta}) - \beta$. PB is calculated by dividing the RB by the true estimate and multiplying by hundred. PB should not exceed 5%.

- Coverage Rate (CR): proportion of confidence intervals that contain the true parameter value. CR is affected by the estimate and the confidence interval. The CR should not fall below the nominal rate. The nominal rate is a pre-specified rate at which the true estimate is expected to fall in the confidence interval. The CR should ideally be around the nominal rate. In this study, the nominal is set at 90%. If the CR falls below the nominal rate, the method is too optimistic, leading to false positives. Conversely, a too high CR, above 95%, indicates a too wide CI and an inefficient method that leads to too conservative inferences.
- Average width (AW): difference between lower and upper end of the confidence interval (CI), average width for the CI. AW should be as small as possible, but not lower than nominal rate. Small AW indicates statistical efficiency, but should not be too small to cause the CR to fall below 90%. AW is also an indication for how well the standard deviations are estimated.

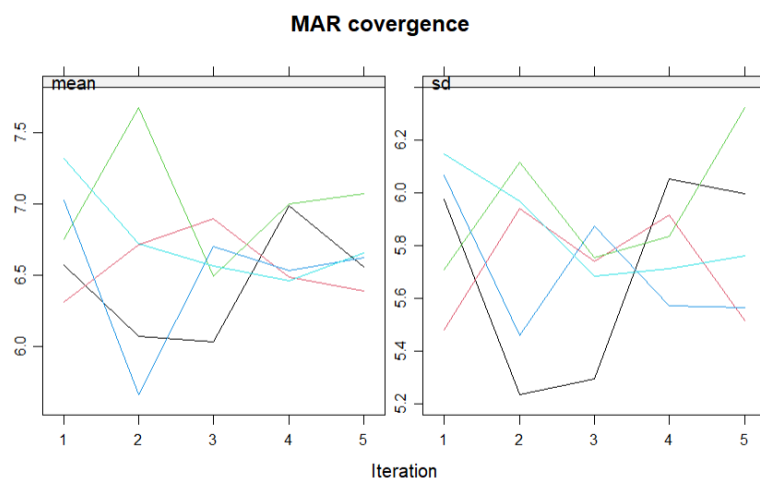
An optimal method has a raw bias close to zero and a coverage rate near 0.95. Methods that fulfill these assumptions are called randomization valid (Rubin, 1987). A shorter confidence interval is preferred over a longer confidence interval (Van Buuren, 2018). The estimates generated from the imputations are compared to the true estimates. These measures are compared for the three proportions of missingness in combination with the mechanisms of MCAR and MAR.

4. Results and analysis

The convergence and plausibility of the imputation are checked through convergence plots and density plots (see Appendix A). The convergence plots derived from a single imputation illustrate that complete convergence is not fully achieved, however after approximately 5 iterations most plots show some convergence. Figure 3 shows one of the convergence plots under the MAR assumption with 30% missingness.

Figure 3

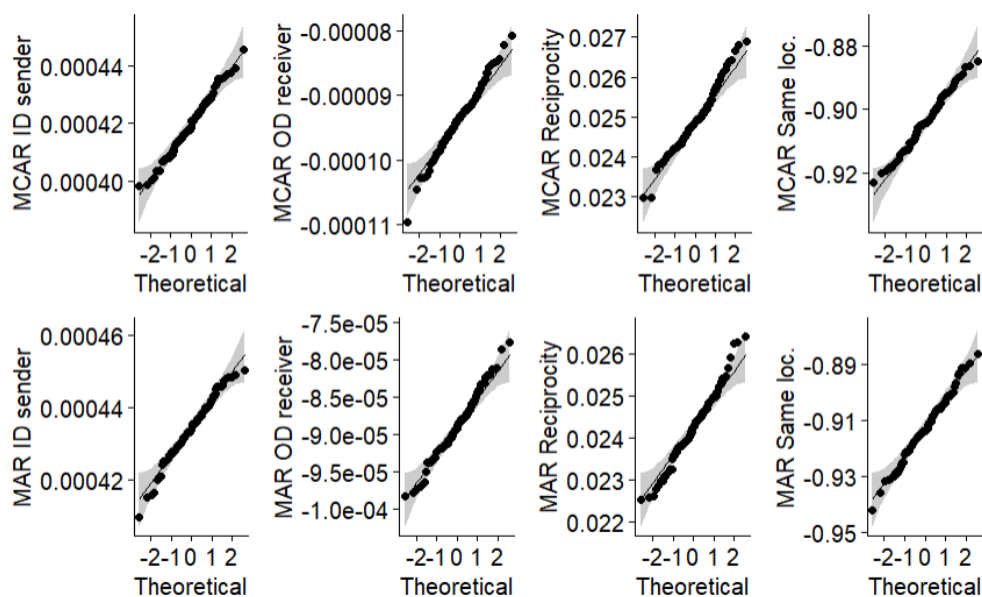
Convergence plot under MAR assumptions with 30% missingness, derived from a single imputation



The convergence of MI could probably improve with more iterations. Considering the computational efficiency, the number of imputations in this study was set to five. The distribution of the imputed values is shown by the quantile-quantile plots (Nguyen et al., 2017). The quantile-quantile plots, see Figure 4 for an example under MCAR and MAR 30% missingness, show that the estimates are not normally distributed for all statistics after all, with some showing more extreme deviations than others.

Figure 4

Quantile-quantile plots for distribution of estimates for each statistic under MCAR and MAR assumptions with 30% missingness



The density plots for both MCAR and MAR show a (slight) deviation from the theoretical distribution assumption (see Appendix A). The distribution for MAR missingness is less skewed towards the right than the distribution for MCAR missingness. This is not what we would expect, as a right-tailed MAR mechanism is used. Even though the missingness is higher on the right side with right-tailed MAR mechanism, the distribution of the missingness is less skewed towards the right compared to MCAR.

Results true analysis

The REM applied on the fully observed data gives the true model of the subset of Apollo 13 data, enabling comparison with the REM results of the simulations under MCAR and MAR. Table 1 shows a positive but not significant effect for *reciprocity* ($\beta = 0.0233$, $p = .209$), which indicates a very small tendency for individuals to reciprocate interactions. However, this effect is not statistically significant. There is no strong evidence to support that individuals tend to reciprocate interactions. *Outdegree receiver* has a negative but also not significant effect ($\beta = -9.023e-5$, $p = .225$). The effect of *outdegree receiver* is not statistically significant, which implies no evidence for the negative fact that individuals who initiate more interactions, receive less interactions. The small positive significant *indegree sender* estimate indicates that individuals with a higher indegree (i.e., who are more likely to receive ties from others) have a higher tendency to initiate interactions ($\beta = 0.0004$, $p < .001$). The negative significant effect of *same location* suggests that individuals being in the same location (both in the flight or ground crew) leads to fewer interactions ($\beta = -0.8629$, $p < .001$).

Table 1

REM results on a subset of the Apollo 13 dataset

Statistic	Estimate	Standard Error	p-value
Reciprocity	2.332e-2	1.856e-2	0.209
Indegree sender	4.314e-4	7.398e-5	< .001
Outdegree receiver	-9.023e-5	7.437e-5	0.225
Same location	-8.629e-1	3.217e-2	< .001

Note: The subset consists of the first 38982 rows of the Apollo 13 dataset

Results simulations assuming MCAR

Table 2 shows the REM results of the simulations averaged assuming MCAR. The simulations are conducted with different proportion of missingness: 10%, 30% and 50%. The raw bias of the estimates are for all missingness percentages close to zero, which indicates that the bias is negligible. The percent bias is below 5% for all statistics with a missingness proportion of 10%, which indicates that there is very little bias. The missingness being assumed

MCAR could explain this (Oberman & Vink, 2023). However, increasing the proportion of missingness to 30% and 50% created a percent bias above 5% for *outdegree receiver*, *reciprocity* and *same location*, which can be due to the nature of the statistics to rely on the structure of the network. These statistics are influenced significantly by the network itself and if a higher percentage of the data is missing, there is less information available for accurate imputation, possibly resulting in a higher bias after imputation. The coverage rate of *reciprocity* and *indegree receiver* with missingness proportion of 30% and 50% are below 90%, which is a suboptimal coverage. The coverage rate of *same location* for all missingness percentages is extremely low, which can potentially be due to the fact that it is dependent on the same location column in the risk set, which is not based on the imputed network like the other statistics. This can lead to an identical status in the same location column in each risk set across simulations and thus result in very similar estimates. In case of low nominal coverage, there is a likelihood of drawing false positive conclusions, which is also known as a type I error. The coverage rate of *reciprocity* with 10% missingness, *indegree sender* with 10% missingness and *outdegree receiver* for all missingness proportions fall within the acceptable range. The confidence interval is small under all missingness proportions for all the statistics, *reciprocity*, *indegree sender*, *outdegree receiver* and *same location*. The average width is close to zero for all effects. The estimates of the imputed model are all significant, which differs from the true model in which only the estimate of *indegree sender* and *same location* is significant.

Table 2

Averaged REM results on the simulations under the MCAR assumption for each missingness proportion

Statistic	Prop	Estimate	Std. error	P-value	CI	CV	RB	PB	AW
Reciprocity	0.1	2.38e-2	5.34e-4	< .001	[2.23e-2, 2.53e-2]	0.93	5.11e-4	2.55	2.97e-3
	0.3	2.49e-2	9.75e-4	< .001	[2.22e-2, 4.47e-2]	0.84	1.60e-3	6.92	5.41e-3
	0.5	2.58e-2	1.37e-3	< .001	[2.27e-2, 2.82e-2]	0.77	2.51e-3	10.76	7.59e-3
Indegree sender	0.1	4.28e-4	5.34e-6	< .001	[4.13e-4, 4.42e-4]	0.93	-3.89e-6	1.24	2.97e-5
	0.3	4.20e-4	9.70e-6	< .001	[3.93e-4, 4.47e-4]	0.89	-1.16e-5	2.99	5.39e-5
	0.5	4.11e-4	1.31e-5	< .001	[3.78e-4, 4.46e-4]	0.76	-2.01e-5	4.79	7.26e-5
Outdegree receiver	0.1	-9.11e-5	2.01e-6	< .001	[-9.67e-5, -8.56e-5]	0.92	-9.09e-7	2.53	1.12e-5
	0.3	-9.36e-5	3.92e-6	< .001	[-1.04e-4, -8.27e-5]	0.93	-3.32e-6	5.30	2.18e-5
	0.5	-9.48e-5	5.39e-6	< .001	[-1.02e-4, -8.84e-5]	0.91	-4.55e-6	6.61	3.00e-5
Same location	0.1	-8.78e-1	4.77e-3	< .001	[-8.91e-1, -8.65e-1]	0.39	-1.53e-2	1.78	2.65e-2
	0.3	-9.04e-1	1.10e-2	< .001	[-9.34e-1, -8.73e-1]	0.24	-4.09e-2	4.74	6.12e-2
	0.5	-9.28e-1	1.71e-2	< .001	[-9.76e-1, -8.81e-1]	0.19	-6.54e-2	7.57	9.48e-2

Prop = Proportion of missingness; Std. error = Standard error; CI = Confidence interval; CV = Coverage; RB = Raw bias; PB = Percent bias; AW = Average width

Results simulations assuming MAR

Table 3 shows the REM results of the simulations averaged assuming MAR. The simulations are also conducted with different proportion of missingness: 10%, 30% and 50%. The raw bias of the estimates for all missingness proportions are close to zero, indicating that

the bias is negligible. The percent bias is below 5% for almost all statistics, which indicates that there is very little bias. Only *reciprocity* with a missingness proportion of 50% and *same location* with a missing proportion of 30% and 50% have a percent bias above 5%. The percent bias is overall lower for MAR than by assuming MCAR, except for *same location*. The coverage rate of *reciprocity* and *indegree sender* with 10% and 30% missingness and *outdegree receiver* with 10% missingness are slightly below 90% which is suboptimal. It is not extremely low, but there is a chance of creating a type I error. A chance of creating a type I error is high with *same location*, as the coverage rate of *same location* is extremely low for all missingness proportions. The coverage rate of *indegree sender* with 30% and 50% missingness, *outdegree receiver* with 50% missingness are relatively high. The confidence interval is small for all the statistics. The average width is close to zero for all effects. The estimates of the imputed model are all significant, which differs from the true model in which only the estimate of *indegree sender* and *same location* is significant.

Table 3

Averaged REM results on the simulations under the MAR assumption for each missingness proportion

Statistic	Prop	Estimate	Std. error	P-value	CI	CV	RB	PB	AW
Reciprocity	0.1	2.35e-2	3.85e-4	< .001	[2.24e-2, 2.46e-2]	0.83	1.95e-4	2.00	2.14e-3
	0.3	2.42e-2	7.57e-4	< .001	[2.22e-2, 2.64e-2]	0.85	9.35e-4	4.43	4.20e-3
	0.5	2.50e-2	1.20e-3	< .001	[2.17e-2, 2.83e-2]	0.91	1.70e-3	7.35	6.67e-3
Indegree sender	0.1	4.33e-4	4.34e-6	< .001	[4.21e-4, 4.45e-4]	0.92	1.47e-6	0.87	2.41e-5
	0.3	4.34e-4	8.46e-6	< .001	[4.11e-4, 4.58e-4]	0.97	2.73e-5	1.54	4.70e-5
	0.5	4.32e-4	1.13e-5	< .001	[4.01e-4, 4.64e-4]	0.99	8.18e-7	1.90	6.25e-5
Outdegree receiver	0.1	-8.93e-5	1.82e-6	< .001	[-9.43e-5, -8.42e-5]	0.87	9.68e-7	2.41	1.01e-5
	0.3	-8.89e-5	3.63e-6	< .001	[-9.90e-5, -7.88e-5]	0.93	1.34e-6	3.79	2.02e-5
	0.5	-8.90e-5	4.73e-6	< .001	[-1.02e-4, -7.59e-5]	0.96	1.24e-6	4.74	2.63e-5
Same location	0.1	-8.82e-1	5.90e-3	< .001	[-8.98e-1, -8.65e-1]	0.37	-1.88e-2	2.18	3.27e-2
	0.3	-9.13e-1	1.29e-2	< .001	[-9.49e-1, -8.77e-1]	0.16	-4.99e-2	5.79	7.14e-2
	0.5	-9.38e-1	1.80e-2	< .001	[-9.88e-1, -8.88e-1]	0.13	-7.52e-2	8.71	9.97e-2

Prop = Proportion of missingness; Std. error = Standard error; CI = Confidence interval; CV = Coverage; RB = Raw bias; PB = Percent bias; AW = Average width

Comparison results

The simulated results assuming both MCAR and MAR show a significant effect of *reciprocity* and *outdegree receiver* with much smaller standard errors compared to the true model. A higher precision is achieved in these models with smaller standard errors compared to the true analysis, because multiple imputation uses the underlying data distribution of the observed data for the imputation. This can reduce the variability which was present before amputation and thus result in a higher precision. The estimates of both simulations are close to the true model's estimate, but the significance differs, making the simulations show significant results while the true model does not. This can be due to using a fixed set to keep all actors in the data and thus missingness is only generated in a part of the data, resulting in a small

between variance, smaller confidence intervals and smaller p-values (Van Buuren, 2018). The sampling variance has not been considered in the simulation evaluations of the imputation performance due to the use of finite populations (Vink & Van Buuren, 2014). A constant was introduced by this variance, which represents the conserved part of the dataset. This led to an underestimation of the variance and can thus be the reason for the significant effects observed in the analysis of the simulations.

For *indegree sender* and *same location*, the REM results under both mechanisms show significant effects with estimates very close to the true model's estimate with also much smaller standard errors than for the true model. A higher proportion of missingness resulted in an increased bias for both MCAR and MAR, likely because there is less data available for multiple imputation. This can also be the reason for increase in the standard error and confidence interval when the missingness proportion increases with MAR, indicating less accurate estimates. For MCAR this is not found for all statistics, which can be explained that the imputation of missing values is not dependent on the observed values. However, the bias in MAR simulations was mostly below 5%, indicating better results than assuming MCAR. Multiple imputation thus produces unbiased inferences under the MAR assumptions, however unexpected significant results are found.

5. Discussion and conclusion

This study provides insights into the use of multiple imputation for producing valid inferences when applied on REH data. The estimates of a REM of the complete data are compared to the REM analysis results of the simulations where the data is MCAR or MAR. The statistics considered were *reciprocity*, *indegree sender*, *outdegree receiver* and *same location*. The true model showed non-significant effects for *reciprocity* and *outdegree receiver*, the analysis on the simulated models showed for both MAR and MCAR significant effects for all four predictors. This can be due to the missingness being only generated in a part of the data, resulting in a reduced variability, smaller confidence intervals and smaller p-values in combination with not accounting for sampling variance due to the finite population, which leads to an underestimation of the variance (Van Buuren 2018; Vink & Van Buuren, 2014). However, despite the inconsistency in the significance of the results, the use of MI under both MCAR and MAR shows potential benefits, such as improving the accuracy and reliability of analyses, with handling missing data in REH data.

Additionally, the proportion of missing data was set at 10%, 30% and 50% for both MCAR and MAR. Previous research already suggests that the percentage of missingness should not be

too high when using imputation methods (Li et al., 2015; Sterne et al., 2009). This study confirms these findings: the bias increased significantly with a higher proportion of missingness, making the use of multiple imputation less effective when the missing data percentage is high. This is not a characteristic specifically for multiple imputation but holds for all imputation methods.

This study has several limitations that are important to consider. The choice of the statistics included in the REM can have an effect on the outcomes. Three out of four statistics used in this study are endogenous statistics, which used internal information until a given time point. Using more exogenous statistics can give more insight into the effect of external information on REH data. More attribute variables can also be added to the data and used as exogenous predictors for imputation. This can potentially improve the imputation method. Another limitation is the preserving of a part of the dataset, which can have a negative impact on the variance estimates. It can potentially have introduced or preserved some selective biases or inaccuracies in the results. Further studies need to examine the impact of this to ensure more accurate and reliable results. Additionally, the number of imputations was set to five due to the computational efficiency, increasing the number of imputations can possibly improve the outcomes. Besides, in this study a right-tailed MAR mechanism was used, as this is known to mimic a severe missingness scenario realistic in real-world datasets. However, to get more insight into the effectiveness of MI under different conditions, it is valuable to consider other directions, such as left- and mid-tailed MAR mechanisms. Finally, it is important to consider the unexpected significant results and improved precision that are found in the analysis of the simulation.

Further research is needed to investigate the effectiveness of multiple imputation to produce valid inferences assuming all possible missingness patterns and mechanisms. Additionally, increasing the number of imputations and iterations, considering more statistics and using the entire dataset instead of just a part of the dataset could provide more robust findings. The conservation of only part of the dataset negatively impact the variance of the estimates, which could be addressed in future studies. Future research should also explore the effects of varying the proportion of missing data and different weights to determine the frequency of each pattern.

In conclusion, this study provides valuable insight into the use of MI on REH data. It shows that under the assumption of MCAR and MAR, MI can be a valid method for missing data in REH data when the percentage of missingness is not too high. The limitations highlight the need for further research to confirm and expand upon the results obtained in this study.

References

- Apollo 13 Real-time, (n.d.). <http://apollo13realtime.org/>.
- Arena, G., Lakdawala, R., Meijerink-Bosman, M., Karimova, D., Generoso Vieira, F., Shafiee Kamalabad, M., Leenders, R., Mulder, J. (2023). *_remify: Processing and transforming REH to formats suitable for the remverse packages and more_*. R package version 3.0.0, <<https://github.com/TilburgNetworkGroup/remify>>.
- Back, M. D. (2021). Social interaction processes and personality. In Elsevier eBooks (pp. 183–226). <https://doi.org/10.1016/b978-0-12-813995-0.00008-x>
- Borgatti, S. P., Carley, K. M., & Krackhardt, D. (2006). On the robustness of centrality measures under conditions of imperfect data. *Social Networks*, 28(2), 124–136. <https://doi.org/10.1016/j.socnet.2005.05.001>
- Butts, C. T. (2008). 4. A Relational Event Framework for Social Action. *Sociological Methodology*, 38(1), 155–200. <https://doi.org/10.1111/j.1467-9531.2008.00203.x>
- Butts, C. T. (2009). Revisiting the Foundations of Network Analysis. *Science*, 325(5939), 414–416. <https://doi.org/10.1126/science.1171022>
- Butts, C. T., & Marcum, C. S. (2017). A Relational Event Approach to Modeling Behavioral Dynamics. In *Computational social sciences* (pp. 51–92). https://doi.org/10.1007/978-3-319-48941-4_4
- Carrington, P.J., Scott, J., Wasserman, S. (2005). Models and Methods in Social Network Analysis. In *Cambridge University Press eBooks*. <https://doi.org/10.1017/cbo9780511811395>
- Collins, L. M., Schafer, J. L., & Kam, C. M. (2001). A comparison of inclusive and restrictive strategies in modern missing data procedures. *Psychological Methods*, 6(4), 330–351. <https://doi.org/10.1037/1082-989x.6.4.330>
- Hand, D. J. (2020). Dark data: Why what you don't know matters. <http://scholarvox.library.inseec-u.com/catalog/book/88876396>
- Hammersley, J. M., & Handscomb, D. C. (1964). Monte Carlo methods. In Springer eBooks. <https://doi.org/10.1007/978-94-009-5819-7>
- Huisman, M. (2009). Imputation of missing network data: Some simple procedures. DOAJ (DOAJ: Directory Of Open Access Journals). <https://doi.org/10.21307/joss-2019-050>
- Huisman, M., & Krause, R. W. (2017). Imputation of Missing Network Data. In *Springer eBooks* (pp. 1–10). https://doi.org/10.1007/978-1-4614-7163-9_394-1

- Kolmogorov, A. N., & Smirnov, N. V. (1933). On the empirical determination of a distribution function. *Mathematische Annalen*, 109(1), 461-472.
<https://doi.org/10.1007/BF01449206>
- Krivitsky, P. N., & Handcock, M. S. (2014). A Separable Model for Dynamic Networks. *Journal Of The Royal Statistical Society. Series B, Statistical Methodology*, 76(1), 29–46.
<https://doi.org/10.1111/rssb.12014>
- Li, P., Stuart, E. A., & Allison, D.B. (2015). Multiple Imputation: A Flexible Tool for Handling Missing Data. *JAMA*. 314(18), doi:10.1001/jama.2015.15281.
- Little, T. D., Jorgensen, T. D., Lang, K. M., & Moore, E. W. G. (2013). On the Joys of Missing Data. *Journal Of Pediatric Psychology*, 39(2), 151–162. <https://doi.org/10.1093/jpepsy/jst048>
- Mack, C., Su, Z., & Westreich, D. (2018). Analytic Implications and Management Strategies For Missing Data. *Managing Missing Data in Patient Registries - NCBI Bookshelf*.
<https://www.ncbi.nlm.nih.gov/books/NBK493610/>
- Marcum, C. S., & Butts, C. T. (2015). Constructing and Modifying Sequence Statistics for relevant Using informr. *Journal Of Statistical Software*, 64(5).
<https://doi.org/10.18637/jss.v064.i05>
- Meijerink-Bosman, M., Leenders, R., & Mulder, J. (2022). Dynamic relational event modeling: Testing, exploring, and applying. *PloS One*, 17(8), e0272309.
<https://doi.org/10.1371/journal.pone.0272309>
- Meijerink-Bosman M, Arena G, Karimova D, Lakdawala R, Shafiee Kamalabad M, Generoso Vieira F. (2023). TilburgNetworkGroup/remstats: Computes Statistics For Relational Event History Data. *rdrr.io*. <https://rdrr.io/github/TilburgNetworkGroup/remstats/>
- Newman, M. F. (2018). *Networks*. In Oxford University Press eBooks.
<https://doi.org/10.1093/oso/9780198805090.001.0001>
- Nguyen, C., Carlin, J. B., & Lee, K. J. (2017). Model checking in multiple imputation: an overview and case study. *Emerging Themes in Epidemiology*, 14(1).
<https://doi.org/10.1186/s12982-017-0062-6>
- Oberman, H. I., Van Buuren, S., & Vink, G. (2021). Missing the Point: Non-Convergence in Iterative Imputation Algorithms. *arXiv (Cornell University)*.
<https://doi.org/10.48550/arxiv.2110.11951>
- R Core Team (2023). *_R: A Language and Environment for Statistical Computing_*. R Foundation for Statistical Computing, Vienna, Austria. <https://www.R-project.org/>.
- Rubin, D. B. (1976). Inference and missing data. *Biometrika*, 63(3), 581–592.
<https://doi.org/10.1093/biomet/63.3.581>

- Rubin, D. B. (1987). Multiple Imputation for Nonresponse in Surveys. In Wiley series in probability and statistics. <https://doi.org/10.1002/9780470316696>
- Schafer, J. L., & Graham, J. W. (2002). Missing data: Our view of the state of the art. *Psychological Methods*, 7(2), 147–177. <https://doi.org/10.1037/1082-989x.7.2.147>
- Schouten, R. M., Lugtig, P., & Vink, G. (2018). Generating missing values for simulation purposes: a multivariate amputation procedure. *Statistical Computation And Simulation/Journal Of Statistical Computation And Simulation*, 88(15), 2909–2930. <https://doi.org/10.1080/00949655.2018.1491577>
- Schouten, R. M., & Vink, G. (2018). The Dance of the Mechanisms: How Observed Information Influences the Validity of Missingness Assumptions. *Sociological Methods & Research*, 50(3), 1243–1258. <https://doi.org/10.1177/0049124118799376>
- Shafiee Kamalabad, M. S., Leenders, R., & Mulder, J. (2023). What is the Point of Change? Change Point Detection in Relational Event Models. *Social Networks*, 74, 166–181. <https://doi.org/10.1016/j.socnet.2023.03.004>
- Serrat, O. (2017). Social network analysis. *Knowledge solutions: Tools, methods, and approaches to drive organizational performance*, 39–43. https://doi.org/10.1007/978-981-10-0983-9_9
- Sterne, J. A., White, I. R., Carlin, J. B., Spratt, M., Royston, P., Kenward, M. G., Wood, A. M., & Carpenter, J. R. (2009). Multiple imputation for missing data in epidemiological and 21 clinical research: potential and pitfalls. *BMJ (Clinical research ed.)*, 338, b2393. <https://doi.org/10.1136/bmj.b2393>
- Smith, J. A., Moody, J., & Morgan, J. H. (2017). Network sampling coverage II: The effect of non-random missing data on network measurement.
- Snijders, T. A. B. (2001). The Statistical Evaluation of Social Network Dynamics. *Sociological Methodology*, 31(1), 361–395. <https://doi.org/10.1111/0081-1750.00099> Social Networks, 48, 78–99. <https://doi.org/10.1016/j.socnet.2016.04.005>
- Therneau T (2023). *_A Package for Survival Analysis in R_*. R package version 3.5-5., <<https://CRAN.R-project.org/package=survival>>.
- Tseng, I., 2017. Apollo 13 Real-Time. GitHub, <https://github.com/issa-tseng/apollo13rt>
- Van Buuren, S. (2018). Flexible imputation of missing data. CRC press.
- Van Buuren, S. & Groothuis-Oudshoorn, K. (2011). “mice: Multivariate Imputation by Chained Equations in R.” *Journal of Statistical Software*, 45(3), 1–67. doi:10.18637/jss.v045.i03
- Vink, G., Van Buuren, S. (2014). Pooling multiple imputations when the sample happens to be the population. *arXiv (Cornell University)*. <https://arxiv.org/pdf/1409.8542v1>

Vink G., (2022). Strategies for simulating missingness (v1.0). Zenodo.

<https://doi.org/10.5281/zenodo.7467995>

Vink G., (2023). mice::mice.impute.pmm.conditional()

<https://www.gerkovink.com/miceVignettes/>

Wasserman, S. and Faust, K. (1994). Social Network Analysis: Methods and Applications.

Cambridge University Press, Cambridge. <https://doi.org/10.1017/CBO9780511815478>

Appendix A

Figure A1

Density plots for simulations for each statistic under MAR and MCAR assumptions with 10% missingness (blue dotted line indicates the population value of true analysis)

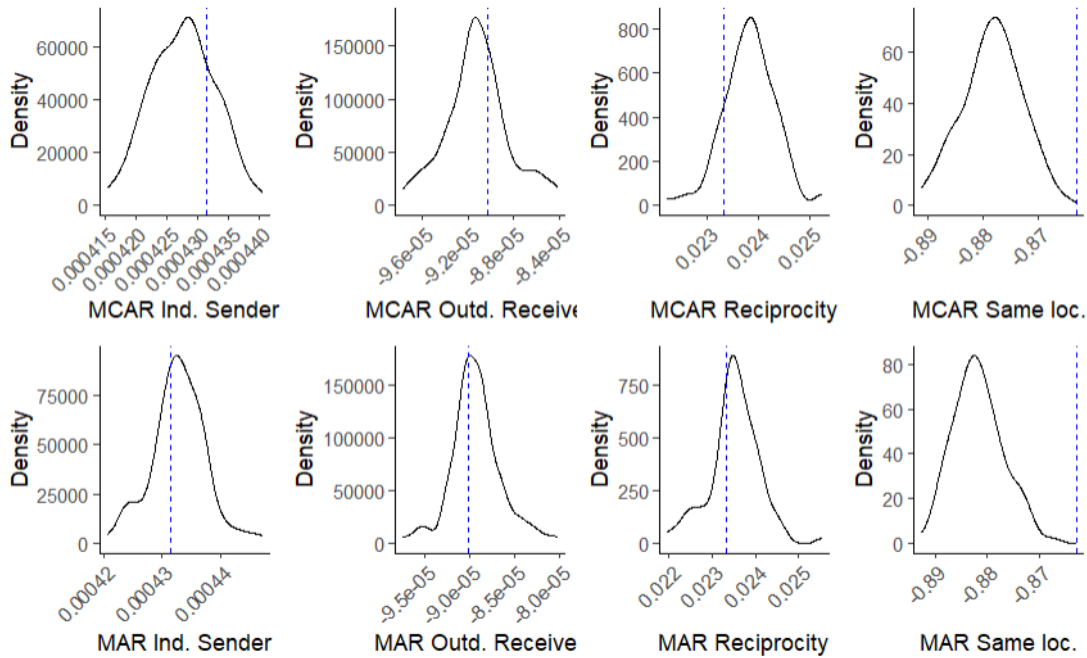


Figure A2

Density plots for simulations for each statistic under MAR and MCAR assumptions with 30% missingness (blue dotted line indicates the population value of true analysis)

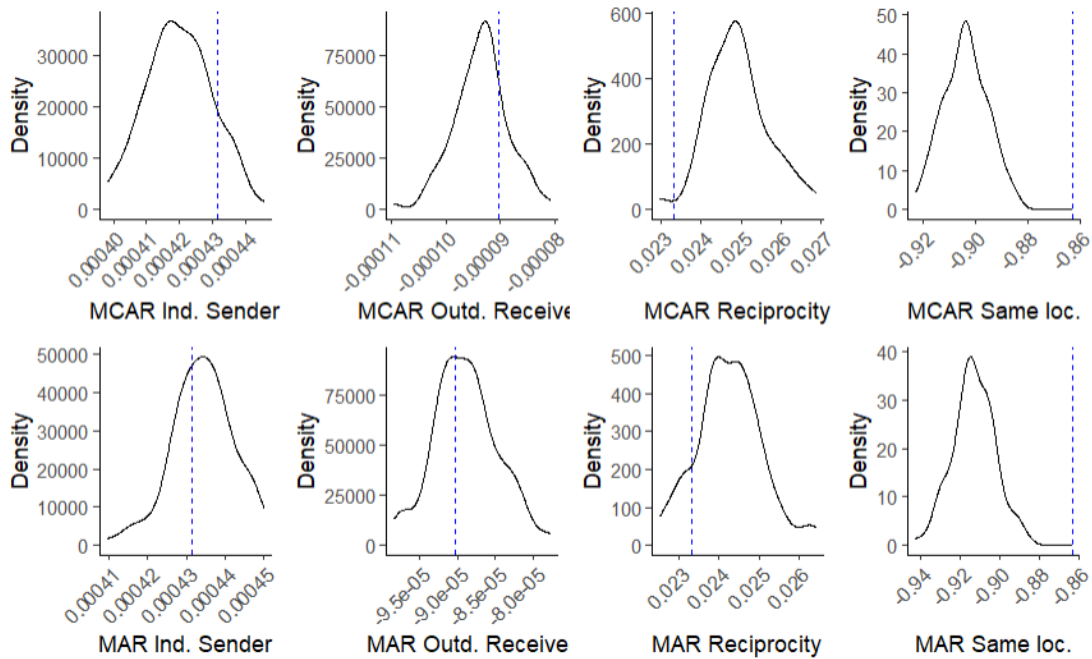
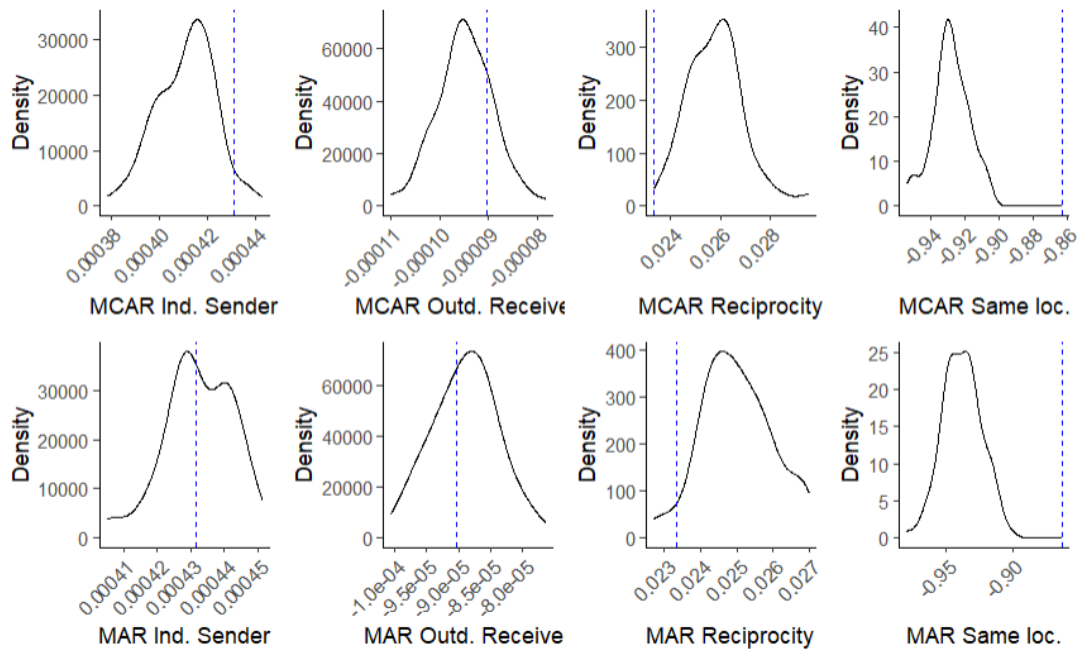


Figure A3

Density plots for simulations for each statistic under MAR and MCAR assumptions with 50% missingness (blue dotted line indicates the population value of true analysis)

**Figure A4**

Quantile-quantile plots for distribution of estimates for each statistic under MCAR and MAR assumptions with 10% missingness

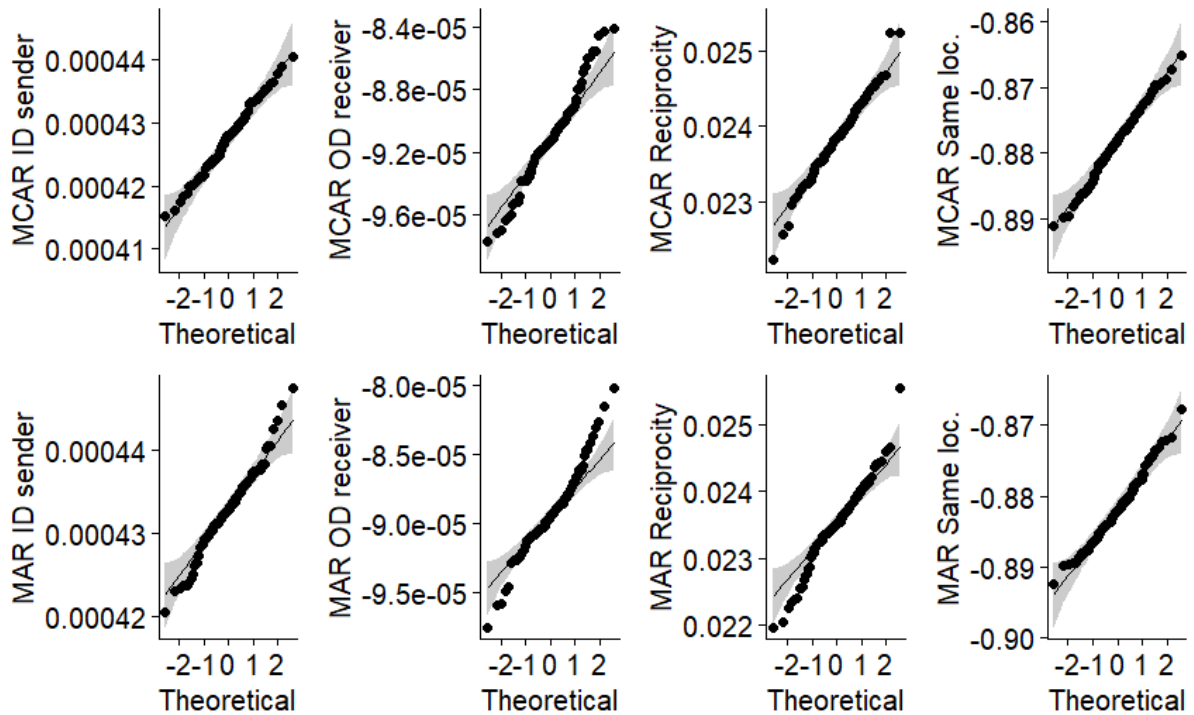
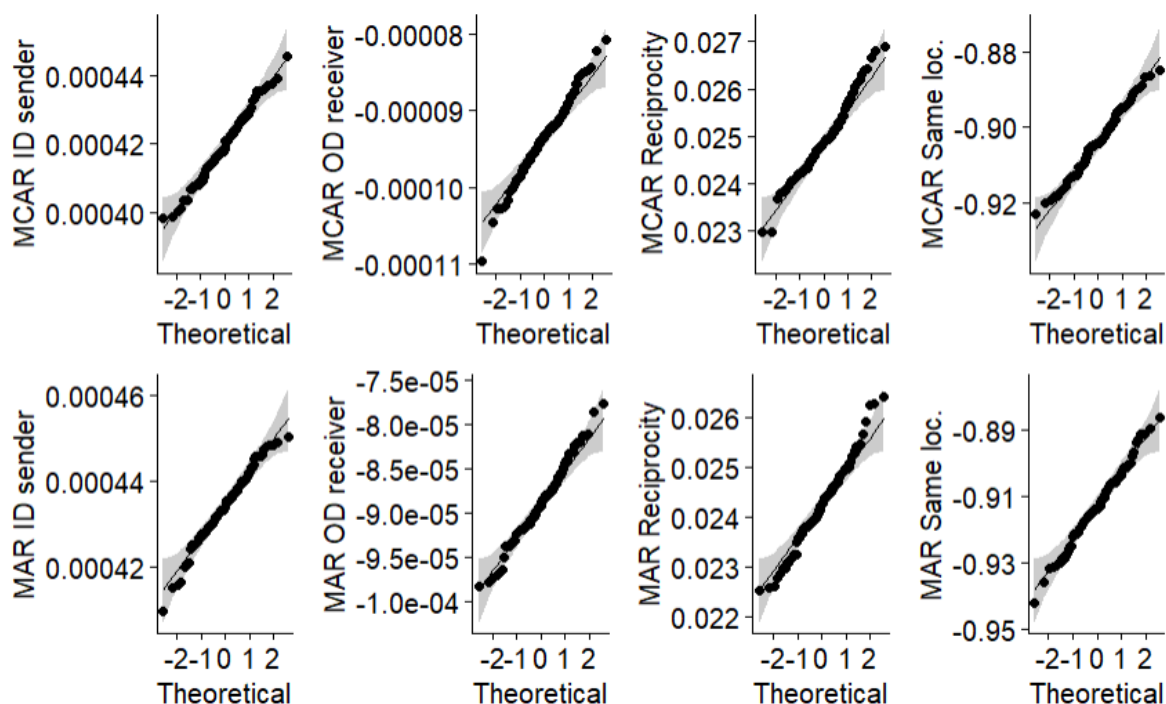


Figure A5

Quantile-quantile plots for distribution of estimates for each statistic under MCAR and MAR assumptions with 30% missingness

**Figure A6**

Quantile-quantile plots for distribution of estimates for each statistic under MCAR and MAR assumptions with 50% missingness

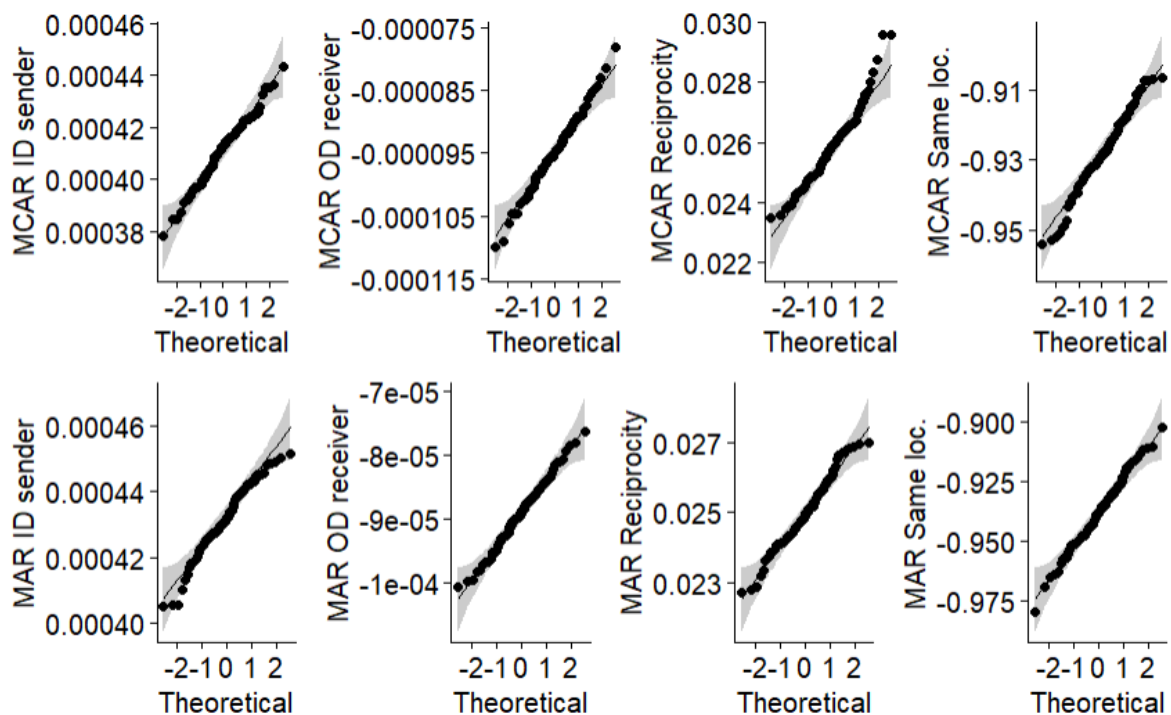
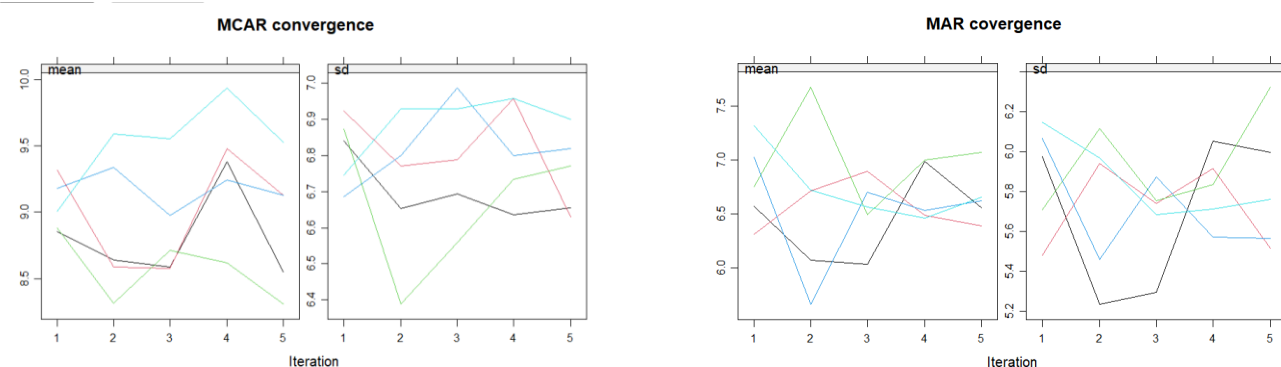
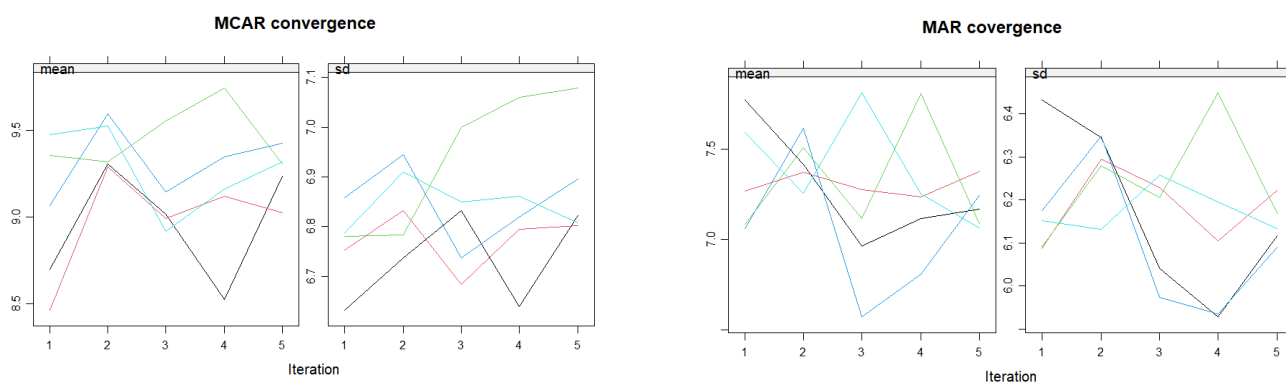


Figure A7

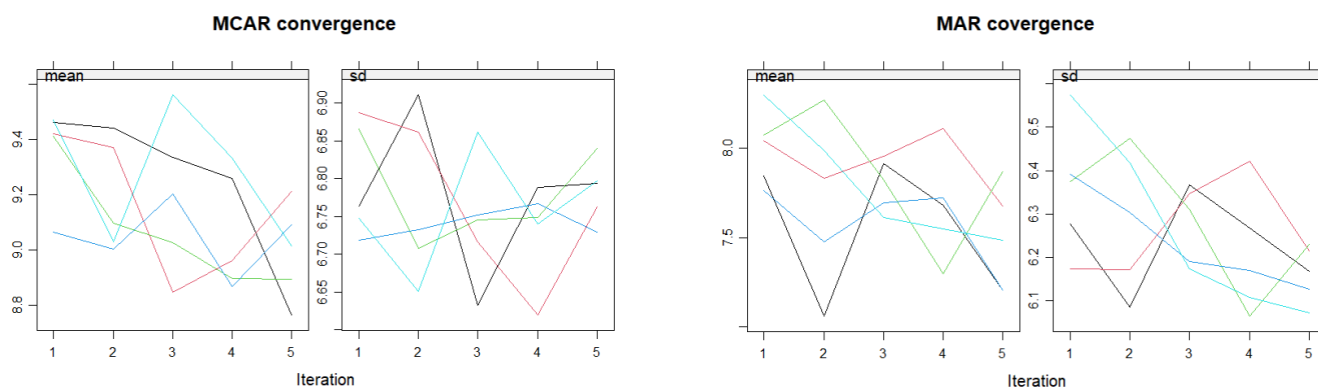
Convergence plots under MCAR and MAR assumptions with 10% missingness, derived from a single imputation

**Figure A8**

Convergence plots under MCAR and MAR assumptions with 30% missingness, derived from a single imputation

**Figure A9**

Convergence plots under MCAR and MAR assumptions with 50% missingness



Appendix B

https://github.com/mshafieek/ADS-Missing-data-social-network/tree/main/Myrthe_2024