

# Diagnosing missing always at random in multivariate data\*

Iavor I. Bojinov, Natesh S. Pillai, and Donald B. Rubin  
*Harvard University*

April 3, 2018

## Abstract

Models for analyzing multivariate data sets with missing values require strong, often unassessable, assumptions. The most common of these is that the mechanism that created the missing data is ignorable - a twofold assumption dependent on the mode of inference. The first part, which is the focus here, under the Bayesian and direct-likelihood paradigms, requires that the missing data are missing at random; in contrast, the frequentist-likelihood paradigm demands that the missing data mechanism always produces missing at random data, a condition known as missing always at random. Under certain regularity conditions, assuming missing always at random leads to an assumption that can be tested using the observed data alone namely, the missing data indicators only depend on fully observed variables. Here, we propose three different diagnostic tests that not only indicate when this assumption is incorrect but also suggest which variables are the most likely culprits. Although missing always at random is not a necessary condition to ensure validity under the Bayesian and direct-likelihood paradigms, it is sufficient, and evidence for its violation should encourage the careful statistician to conduct targeted sensitivity analyses.

**Keywords:** Missing Data, Diagnostic Tools, Sensitivity Analysis, Hypothesis Testing, Missing at Random, Row Exchangeability

## 1 Introduction

When conducting statistical analyses of data sets with missing values, researchers have to make assumptions that can not be assessed using the observed data alone, referred to as unassessable assumptions by [Liublinska and Rubin \(2014\)](#). [Rubin's \(1976\)](#) seminal paper was the first to formalize these assumptions by explicitly considering the missing data indicators as random variables as well as providing the weakest sufficient conditions that lead to correct inference about the parameter that governs the distribution of the data when ignoring the mechanism that created the

---

\*We thank Fabrizia Meall for her useful insights.

missing data. [Rubin \(1976\)](#) showed that it is appropriate to ignore the missingness mechanism (i.e., the conditional distribution of the missing data indicators given the missing and the observed data) when using direct-likelihood ([Fisher, 1956](#); [Edwards, 1984](#)) or Bayesian inference for the parameter governing the distribution of the data,  $\theta$ , when the missing data are missing at random (MAR) and parameter of the missingness mechanism,  $\phi$ , and  $\theta$  are distinct. The missing data are MAR if the missingness mechanism evaluated at the observed missing data pattern and the observed data, considered as a function of the missing data and  $\phi$ , takes the same value for all possible values of the missing data and the parameter  $\phi$ . The definition of MAR depends on the missing data, and it is, therefore, generally unassessable as its validity cannot be falsified using the observed data alone. The much stronger assumption of missing completely at random (MCAR), which is sufficient, although not necessary, to ignore the missing data mechanism for valid frequentist inference, however, can be assessed ([Little, 1988](#); [Park and Davis, 1993](#)). This is true because MCAR further imposes that, for all values of  $\phi$ , the missingness mechanism, evaluated at the observed missing data pattern, takes the same value for all values of the observed data as well as the missing data ([Marini et al., 1980](#)).

Subsequently, there have been many misinterpretations of the conditions provided in [Rubin \(1976\)](#) (e.g., [Lu and Copas \(2004\)](#) used MAR to mean that the missingness mechanism always produces MAR data sets, and [Fitzmaurice et al. \(2012\)](#) interpreted MAR as a conditional independence statement). [Seaman et al. \(2013\)](#) and [Mealli and Rubin \(2015\)](#) are two recent papers that clarify the situation by explaining the difference between the data being missing at random and the missingness mechanism always producing missing at random data. This distinction is critical for understanding when ignoring the missingness mechanism yields valid statements from different inferential perspectives, and what assumptions can be assessed using the observed data alone and under what conditions.

Although the MCAR assumption is sufficient for valid frequentist inference, certain aspects of likelihood-based frequentist inference (that is, using maximum likelihood estimates and the observed information matrix to measure their precision) are asymptotically valid when the missingness mechanism always produces MAR data sets ([Molenberghs and Kenward, 1998](#); [Little and Rubin, 2002](#)). [Seaman et al. \(2013\)](#) referred to this type of missingness mechanism as “everywhere missing at random,” whereas [Mealli and Rubin \(2015\)](#) referred to it as “missing always at random”; we choose to follow the latter suggestion as the word “everywhere” in probability and statistics has a different mathematical meaning, which is not reflected in its use in this context. A missing always at random (MAAR) missingness mechanism always produces MAR data, but not all MAR data sets are generated from a MAAR mechanism. Therefore, assuming MAAR and parameter distinctness is a sufficient condition for ignoring the missing data mechanism when conducting direct-likelihood, frequentist likelihood, or Bayesian inference. With Bayesian inference, the distinct parameters assumption requires independent prior distributions on  $\phi$  and  $\theta$ .

[Mealli and Rubin \(2015\)](#) also formalized the general intuition that having more fully observed covariates makes the MAAR assumption more plausible; we restate the relevant theorem without proof in Section 2. Utilizing this result, we derive two corollaries that we use to construct three different diagnostic tests for diagnosing the plausibility of MAAR. Violating MAAR does not imply that Bayesian or direct-likelihood inference conducted by ignoring the missingness mechanism are

invalid; however, data sets that satisfy MAR but not MAAR are ones where the resulting inference are more sensitive to model misspecification, in the sense first explored in [Rosenbaum and Rubin \(1984\)](#). Each of our diagnostic tests identifies the variables that are more likely to break the MAAR assumption by looking for conditional dependencies between the missing data indicators and the variables with missing values, given the fully observed variables. Focusing on these variables can lead to better, more targeted, sensitivity analyses because these variables are likely creators of the violation of MAR. The diagnostic tests we propose should not be treated as formal hypothesis tests, but rather instead treated like diagnostic tools such as the ones presented in [Potthoff et al. \(2006\)](#); [Abayomi et al. \(2008\)](#); and [Bondarenko and Raghunathan \(2016\)](#).

In Section 2 we provide the notation, definitions and the main results that will be used throughout the paper. In Section 3 we introduce our three diagnostic tests and explain how to apply each of them to a data set. In Section 4 we conduct a simulation study to analyze the frequentist operating characteristics of our proposed tests. In Section 5 we provide concluding comments.

## 2 Notation and definitions

Let  $Y$  be the complete data matrix, arising from distribution  $p(Y \mid \theta)$ , with entries  $Y_{i,j}$  corresponding to the (possibly missing) response of unit  $i = 1, \dots, N$  to variable  $j = 1, \dots, J$ . Define  $Y_{\cdot j} = (Y_{1,j}, \dots, Y_{N,j})^T$  to be the random column vector of the unit responses for variable  $j$ , and let  $Y_i = (Y_{i,1}, \dots, Y_{i,J})$  be the random row vector of unit  $i$ 's responses. Define  $R$  to be the response indicator matrix such that  $R_{i,j} = 1$  if  $Y_{i,j}$  is observed and 0 if  $Y_{i,j}$  is missing. The missingness mechanism is then the conditional distribution of  $R$  given  $Y$  indexed by a parameter  $\phi$ , denoted by  $p(R \mid Y, \phi)$ . The response indicators  $R_{\cdot j} = (R_{1,j}, \dots, R_{N,j})$  partition  $Y_{\cdot j}$  into two sets  $\mathcal{S}(R_{\cdot j}, Y_{\cdot j}) = \{Y_{i,j} \text{ such that } R_{i,j} = 1\}$  the observed values of  $Y_{\cdot j}$  and  $\mathcal{S}(1 - R_{\cdot j}, Y_{\cdot j}) = \{Y_{i,j} \text{ such that } R_{i,j} = 0\}$  the missing values of  $Y_{\cdot j}$ . In general  $\mathcal{S}(R_{\cdot k}, Y_{\cdot j}) = \{Y_{i,j} \text{ such that } R_{i,k} = 1\}$  will denote the set of  $Y_{\cdot j}$  for which  $R_{i,k} = 1$ . The concatenation of  $\mathcal{S}(R_{\cdot j}, Y_{\cdot j})$  will be written as  $\mathcal{S}(R, Y)$  rather than the usual  $Y_{\text{obs}}$  or  $Y_{(1)}$ . We have purposefully kept our notation very general so that we can readily consider the partition that a missingness indicator induces on all of the columns on  $Y$ . Throughout most of the paper, we will assume that only  $J^* < J$  of the columns of  $Y$  can have missing values, then the matrix  $Y_{\cdot J^*+1:J} = (Y_{\cdot J^*+1}, \dots, Y_{\cdot J})$  will always be fully observed, and  $R$  will have  $J^*$  columns. A generic value of  $R$  and  $Y$  will be denoted by  $r$  and  $y$  respectively, whereas the realized values will be indicated by  $\tilde{r}$  and  $\tilde{y}$ , respectively.

We now provide the formal definitions of MAR and MAAR, and then state the key theorem from [Mealli and Rubin \(2015\)](#).

**Definition 1.** *The missing data are missing at random (MAR) if*

$$p(R = \tilde{r} \mid Y = \tilde{y}, \phi) = p(R = \tilde{r} \mid Y = \tilde{y}', \phi), \quad (1)$$

for all  $\phi$ , and all  $\tilde{y}$  and  $\tilde{y}'$  such that  $\mathcal{S}(\tilde{r}, \tilde{y}) = \mathcal{S}(\tilde{r}, \tilde{y}')$ .

*The missing data are missing always at random (MAAR) if*

$$p(R = r \mid Y = y, \phi) = p(R = r \mid Y = y', \phi), \quad (2)$$

for all  $\phi$ , and all  $r, y, y'$  such that  $\mathcal{S}(r, y) = \mathcal{S}(r, y')$ .

**Theorem 1** (Mealli and Rubin (2015)). *Suppose that the missing data are MAAR, and:*

1. *the rows of  $(Y, R)$  are exchangeable,*
2. *the columns of  $R$  are mutually conditionally independent given  $Y$ ,*
3. *the probability that the  $j^{\text{th}}$  variable is missing is positive for some  $\phi$ ,  $p(R_{i,j} = 0 \mid Y_{i\cdot} = y_{i\cdot}, \phi) > 0$ , and moreover this probability depends on the  $k^{\text{th}}$  component of  $Y_{i\cdot}$ ,  $Y_{i,k}$ ,*

$$p(R_{i,j} = 0 \mid Y_{i\cdot} = y_{i\cdot}, \phi) \neq p(R_{i,j} = 0 \mid Y_{i,-k} = y_{i,-k}, \phi), \quad (3)$$

where the subscript “ $-k$ ” indicates the removal of  $k^{\text{th}}$  element of the vector  $Y_{i\cdot}$ . Then  $Y_{i,k}$  must be always observed for all  $i$

$$p(R_{i,k} = 0 \mid Y_{i\cdot} = y_{i\cdot}, \phi) = 0 \quad \text{for all } i \text{ and all } \phi.$$

Theorem 1 implies that, if the vectors  $Y_{\cdot j}$  and  $Y_{\cdot k}$  have missing values, then both  $R_{\cdot j}$  and  $R_{\cdot k}$  can not depend on  $Y_{\cdot k}$  or  $Y_{\cdot j}$  given  $\phi$  and the fully observed outcomes, whenever  $J^* < J$ . Below we formally state this corollary, the proof is given in Appendix 1.

**Corollary 1.** *Assume that only the first  $J^* < J$  columns of  $Y$  have a positive probability of having missing values, and the other  $J - J^*$  columns of  $Y$  are always fully observed. Under the conditions of Theorem 1, the probability that the  $j^{\text{th}}$  variable is missing must only depend on the  $J - J^*$  fully observed variables  $Y_{\cdot J^*+1:J}$ ,*

$$p(R_{i,j} = 0 \mid Y_{i\cdot} = y_{i\cdot}, \phi) = p(R_{i,j} = 0 \mid Y_{i,J^*+1:J} = y_{i,J^*+1:J}, \phi) \quad (4)$$

for all  $\phi$ ,  $i = 1, \dots, N$  and,  $j = 1, \dots, J^*$ .

By applying Bayes theorem, we can further show that conditional distribution of  $Y_{\cdot j}$  is independent of  $R_{\cdot j'}$  for  $j, j' \in \{1, \dots, J^*\}$ , given the observed variables  $Y_{\cdot J^*+1:J}$ . Below we formally state this corollary, the proof is given in Appendix 1.

**Corollary 2.** *Under the conditions of Corollary 1 with parameters  $\theta$  and  $\phi$  distinct,*

$$p(Y_{\cdot j} = y_{\cdot j} \mid Y_{\cdot J^*+1:J} = y_{\cdot J^*+1:J}, R_{\cdot j'} = r_{\cdot j'}, \theta) = p(Y_{\cdot j} = y_{\cdot j} \mid Y_{\cdot J^*+1:J} = y_{\cdot J^*+1:J}, \theta). \quad (5)$$

### 3 Diagnostic tests for MAAR

In this section, we describe three diagnostic tests for detecting violations of Corollary 1 and 2 by introducing each in a simple case and then explaining how it extends to more general circumstances. Henceforth, we assume that the rows of  $Y$  are independent and identically distributed, which simplifies our exposition but is not always necessary.

The diagnostic tests are formulated as hypothesis tests to allow for straightforward examination of their properties through a simulation study. Practitioners should use these diagnostic tests

to identify problematic variables and conduct focused sensitivity analysis, rather than as formal accept-reject hypothesis tests. When one of our diagnostic tests “rejects,” we can conclude that some of the assumptions are likely to be violated; of course, we can not distinguish between a rejection due to a violation of the MAAR assumption as opposed to any of the other conditions, except by using domain knowledge.

### 3.1 A comparison of conditional means approach

Consider three variables;  $Y_1$  and  $Y_2$  can have missing values, and  $Y_3$  is always fully observed. Under the conditions of Corollary 2, the conditional expectation of  $Y_1$  given  $Y_3$  is the same for the two partitions induced by  $R_2$ ,

$$E[Y_{i,1} | Y_{i,3} = y_{i,3}, R_{i,2} = 0] = E[Y_{i,1} | Y_{i,3} = y_{i,3}, R_{i,2} = 1] = E[Y_{i,1} | Y_{i,3} = y_{i,3}]. \quad (6)$$

Therefore,  $\mathcal{S}(R_2, Y_1)$  and  $\mathcal{S}(1 - R_2, Y_1)$  should have the same conditional distributions, given  $Y_3$ . We can not directly test for a difference in  $\mathcal{S}(R_2, Y_1)$  and  $\mathcal{S}(1 - R_2, Y_1)$  because we have not accounted for possible differences due to  $Y_3$ . But, we can regress  $Y_1$  on  $Y_3$  and compare it to the regression of  $Y_1$  on  $Y_3 \times R_2$ , for example using a likelihood ratio test. This is relatively straightforward; for example, if  $Y_1$  is Gaussian, then the maximum likelihood estimates are obtained by applying least squares regression on the  $R_1$  fully observed variables. The rejection of the smaller model in favor of the larger model, which includes the interaction of  $R_2$ , implies that the assumptions of Corollary 2 do not hold.

Similarly,  $\mathcal{S}(R_1, Y_2)$  and  $\mathcal{S}(1 - R_1, Y_2)$  conditional on  $Y_3$ , should also have the same distributions, and we can perform the analogous test. In this example, we have two hypotheses to test; to obtain valid  $p$ -values, we need to account for the multiple testing.

This diagnostic test does not work whenever two variables have the same missingness pattern, but if that is the case, it seems unlikely that the assumption that the columns of  $R$  are mutually conditionally independent given  $Y$  holds.

The above diagnostic test extends to general situations where we have  $J$  variables,  $J^*$  of which have missing values. As  $J^*$  grows, the number of hypothesis tests to be conducted grows exponentially.

### 3.2 Directly testing a postulated missingness mechanism approach

With three variables,  $Y_1$  and  $Y_2$  that can potentially have missing values, and  $Y_3$  that is always fully observed, by Corollary 1, a MAAR missingness mechanism can only be a function of the fully observed variable,  $Y_3$ :

$$p(R_1 = r_1 | Y = y, \phi) = f_1^\phi(y_3) \quad \text{and} \quad p(R_2 = r_2 | Y = y, \phi) = f_2^\phi(y_3), \quad (7)$$

for some functions  $f_1^\phi$  and  $f_2^\phi$ . With specific forms for  $f_1^\phi$  and  $f_2^\phi$ , we can directly test (7). For example, focusing on  $R_1$ , assume that the missingness mechanism follows a logistic regression

linear in  $Y_{.3}$ , and we test,

$$\begin{aligned} H_0 : p(R_{.1} = r_{.1} \mid Y = \tilde{y}, \phi) &= \text{logit}^{-1}(\alpha + \beta_3 \tilde{y}_{.3}) \\ \text{versus } H_A : p(R_{.1} = r_{.1} \mid Y = \tilde{y}, \phi) &= \text{logit}^{-1}MAAR(\alpha + \beta_1 \tilde{y}_{.1} + \beta_2 \tilde{y}_{.2} + \beta_3 \tilde{y}_{.3}). \end{aligned} \quad (8)$$

“Rejection” of  $H_0$  implies that the assumptions of Corollary 1 are violated.

Performing such a hypothesis test, with a specified missingness mechanism, can be challenging. One possible direction is to use multiple imputation (Rubin, 1987) to generate imputations under the null hypothesis. Using the completed data sets, we can conduct the hypothesis test, specified in (8), and obtain a  $p$ -value, using the appropriate adjustment as described in Meng and Rubin (1992). Once we have tested  $R_{.1}$ , we could also test a, possibly different, postulated missingness mechanism for  $R_{.2}$ . This diagnostic tests generalizes easily to any value of  $J$  and  $J^*$  as well as to any postulated missingness mechanism.

### 3.3 Gaussian copula approach

Suppose we model the joint distribution of  $(Y, R)$  using a Gaussian copula, which implies a simple procedure that can be used both for diagnostic purposes and for generating multiple imputations (Hoff, 2007; Hollenbach et al., 2017). Copulas factor the joint distribution of a multivariate variable into the univariate marginal distributions times a factor called the copula. This approach was initially developed in Sklar (1959); see Nelson (2006) for a more recent treatment.

The main benefit of the copula approach is the decoupling of the correlation from each of the marginal distribution, which can allow for simple modeling of continuous and categorical variables. The appeal of a semiparametric methods for copula models is that they treat the marginal distributions as nuisance parameters, and thereby reduce the amount of information that must be specified *a priori*. In particular, we propose using the approach of Hoff (2007), which employs the rank likelihood for semiparametric copula estimation.

For unit  $i$ , define  $W_i = (Y_i, R_i) = (Y_{i,1}, \dots, Y_{i,J}, R_{i,1}, \dots, R_{i,J^*})$ , and assume

$$W_{i,j} = F_j(\Phi(Z_{i,j})) \quad Z_{i,j} \sim N(0, C), \quad j = 1, \dots, J, J+1, \dots, J+J^*,$$

where  $F_j$  is the univariate CDF of variable  $j$ ,  $\Phi(x)$  is the CDF of the standard Gaussian evaluated at  $x$ , and  $C$  is the correlation matrix, which we assume has an inverse-Wishart prior distribution. Details of the algorithm for estimation are in Hoff (2007).

Once we obtain a posterior distribution of  $C$ , we can obtain the posterior distribution of  $\text{cor}(W_i, W_{i'} \mid W_{\mathcal{K}})$  for any  $i, i'$  and  $\mathcal{K} \subset \{1, \dots, J+J^*\}$ , which we use to test the hypothesis that  $Y_{.j}$  is correlated with  $R_{.j'}$  conditional on  $Y_{.J^*+1:.J}$ . If we reject this hypothesis, meaning that the variables are dependent in the latent scale, then we can conclude that the assumptions for Corollary 1 are likely violated.



## 4 Simulation study

### 4.1 Factors and their levels

We evaluate each of the three proposed diagnostic tests under an array of different scenarios by tracking the number of correct decisions made, where a decision is labeled correct when a procedure rejects an incorrect null hypothesis or fails to reject a correct null hypothesis. Our outcome of interest is the proportion of correct decisions made, as opposed to the type I error and power, because statisticians typically do not know the true missingness mechanism and are rarely interested in studying its properties; their primary aim is to conduct an appropriate analysis under explicitly stated plausible assumptions.

Suppose  $Y$  follows a multivariate normal distribution with mean 0 and covariance matrix  $\Sigma = \rho \mathbb{1}_J \mathbb{1}_J^T + (1 - \rho) I_J$ , where  $\mathbb{1}_J$  is a vector of length  $J$  with entries all equal to 1, and  $I_J$  is the  $J \times J$  identity matrix, where  $J = 5$ ;  $Y_{.1}$ ,  $Y_{.2}$  and  $Y_{.3}$  can potentially have missing values, whereas  $Y_{.4}$  and  $Y_{.5}$  are fully observed.

We partition the five factors of our simulation study into three categories: nature's factors (uncontrolled by the investigator); nature's estimable factors (unknown to the investigator but estimable from the data) and; factors over which the investigator has direct control.

Nature's unknown factors:

Factor 1: the missingness mechanism with levels: a MAAR mechanism that satisfies the conditions of Theorem 1 (denoted by MAAR), a MAAR mechanism where the columns of  $R$  are not mutually conditionally independent given  $Y$  (denoted by MAAR2), and a missing always not at random mechanism that satisfies the other conditions of Theorem 1 (denoted by MNAAR). Table 1, provides the specifications of each of the missingness mechanisms.

Nature's known factors:

Factor 2: sample size,  $N = 100, 500, 1000$ ;

Factor 3: average proportion of missing values per variable,  $m = 0.2, 0.4, 0.6$ .

Nature's estimable factors:

Factor 4: correlation between the columns of  $Y$ ,  $\rho = 0.2, 0.4, 0.6, 0.8$ .

Factors under control of the investigator:

Factor 5: diagnostic test: comparison of conditional means approach (CCM), directly testing a postulated missingness mechanism (DTPMM), Gaussian copula approach (GC).

Table 1: The dependence of the missingness mechanism for the simulation study. For example, the first entry in the  $R_{i,1}$  column means that  $p(R_{i,1} = 1 | Y = y, R_{i,-1} = r_{i,-1}, \phi) = \text{logit}^{-1}(\alpha + y_{i,4} - y_{i,5})$ . The constant  $\alpha$  is used to determine the average proportion of missing values per variable.

Mechanism	$R_{i,1}$	$R_{i,2}$	$R_{i,3}$
MAAR	$y_{i,4} - y_{i,5}$	$y_{i,4} - y_{i,5}$	$y_{i,4} - y_{i,5}$
MAAR2	$y_{i,4} - y_{i,5}$	$y_{i,1}r_{i,1} + y_{i,4} - y_{i,5}$	$y_{i,1}r_{i,1} + y_{i,2}r_{i,2} + y_{i,4} - y_{i,5}$
MNAAR	$y_{i,4} - y_{i,5}$	$\frac{1}{2}y_{i,1} + \frac{1}{2}y_{i,3} + y_{i,4} - y_{i,5}$	$y_{i,1} + y_{i,2} + y_{i,4} - y_{i,5}$

Table 2: Correct decision rates average over the correlation and missingness coefficient factors.

Diagnostic Test	Sample Size	$m =$	$\rho = 0.2$			$\rho = 0.4$			$\rho = 0.6$			$\rho = 0.8$			Average
			0.2	0.4	0.6	0.2	0.4	0.6	0.2	0.4	0.6	0.2	0.4	0.6	
CCM	100		.70	.74	.61	.69	.71	.60	.60	.62	.52	.47	.49	.44	.60
	500		.98	.99	.99	.98	.98	.98	.98	.98	.98	.97	.98	.94	.98
	1000		.98	.99	.98	.98	.98	.98	.98	.98	.98	.99	.98	.99	.98
DTPMM	100		.38	.33	.33	.36	.33	.33	.35	.33	.33	.34	.33	.33	.34
	500		.99	.99	.80	.99	.98	.74	.98	.93	.56	.82	.65	.37	.82
	1000		1	1	1	1	1	.99	1	1	.98	.99	.97	.69	.97
GC	100		.78	.83	.68	.79	.82	.68	.73	.76	.63	.58	.60	.52	.70
	500		.98	.98	.98	.98	.98	.98	.98	.98	.99	.98	.99	.94	.98
	1000		.98	.99	.98	.98	.98	.98	.98	.98	.98	.98	.98	.99	.98

## 4.2 Results

For each possible combination of the factors, we generated 1,000 independent replications. Table 4 in Appendix 2 presents an analysis of variance of the 5-factor study to suggest which tables we should examine. The major sources of variation are the sample size  $N$ , Factor 2, followed by the missingness mechanism, Factor 1, and the diagnostic test used, Factor 5. Because the missingness mechanism is unknown, we focus the discussion of the results on the behavior of the diagnostic tests under nature's known and estimable factors. An analysis of the results across different levels of the missingness mechanism factor is included in Appendix 2.

Table 2 shows the proportion of correct decisions made by the three diagnostic tests for different sample sizes, correlations  $\rho$ , and proportion of missing values per variable ( $m$ ). For large sample sizes, all three diagnostic tests reach the correct decision over 97% of the time, for all levels of the other factors. For small samples sizes, directly testing a postulated missingness mechanism performs poorly because the likelihood ratio test used is based on the asymptotic distribution of the test statistic and has low power in small samples (Liu and Enders, 2017). In small samples, the Gaussian copula slightly outperforms the comparison of conditional means approach, but the difference becomes negligible as the sample size increases.

The main advantage of running diagnostics individually for each of the missingness indicators is to identify variables that are likely to violate the assumptions of Theorem 1. Table 3 presents



Table 3: Correct decision rates for each of the three missingness indicators averaged over the missingness mechanism, sample size and correlation.

Diagnostic Test	$N$	$R_{.1}$	$R_{.2}$	$R_{.3}$
CCM	100	.98	.42	.55
	500	.98	.85	.98
	1000	.98	.97	.99
DTPMM	100	1	.33	.34
	500	1	.47	.82
	1000	1	.74	.97
GC	100	.99	.47	.66
	500	.98	.91	.98
	1000	.97	.98	.99

the correct decision rates for each of the missingness indicators as well as the overall rate for the three diagnostic tests across the different sample sizes. We present the results averaged over the correlation and proportion of missing values because, as reflected in Table 4 the effect of their interaction is small. Under all three missingness mechanisms,  $R_{.1}$  always satisfies the conditions of Theorem 1. For all sample sizes, all three diagnostic tests make the correct decisions regarding  $R_{.1}$  at least 97% of the time. Under the MAAR2 and MNAAR missingness mechanisms, both  $R_{.2}$  and  $R_{.3}$  violate the conditional independence assumption ( $R_{.2}$  depends on  $R_{.1}$ , and  $R_{.3}$  depends on both  $R_{.1}$  and  $R_{.2}$ ), and the MAAR mechanism assumption ( $R_{.3}$  has a greater dependence than  $R_{.2}$  on variables with missing values), respectively. As expected, the diagnostic tests have a higher correct decision rate for  $R_{.3}$  relative to  $R_{.2}$ ; however, the difference reduces as the sample size increases. Again, the comparison of conditional means and Gaussian copula diagnostic tests have relatively similar performance, and directly testing a postulated missingness mechanism does not perform well when the sample size are small.

### 4.3 Discussion

Our simulations suggest that statisticians should use either the Gaussian copula or the comparison of conditional means approach. The choice between the two diagnostic tests depends on the statistician’s beliefs about the correlation structures in  $Y$ . Directly testing a postulated missingness mechanism should be avoided for small samples and only used if the statistician can provide a plausible model for the missingness mechanism.

In this simulation study, we did not consider scenarios for which the global modeling assumptions are violated. We expect to see a reduction in performance due to model misspecification; however, we still believe that these tests can provide the statistician with useful information regarding which variables are likely to violate the MAAR missingness assumptions.

## 5 Conclusion

Using Corollary 1 and 2, we proposed three diagnostic tests for detecting the validity of the assumptions in Theorem 1, using the observed data alone. We showed, through simulation, that all three diagnostic tests had high discriminatory power and can identify the variables that violate the assumptions. As expected, we cannot distinguish between violations of the MAAR assumption and the assumption that the columns of  $R$  are mutually conditionally independent given  $Y$ .

We encourage practitioners to use these diagnostic tests in the diagnosis stage to identify problematic variables and then conduct a targeted sensitivity analysis to assess the impact of violations of suspect assumptions on the scientific conclusions. The identification of problematic variables is especially important when the number of variables is large because most sensitivity analysis methods are computationally intensive.

## References

- Abayomi, K., Gelman, A., and Levy, M. (2008). Diagnostics for multivariate imputations. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 57(3):273–291.
- Bondarenko, I. and Raghunathan, T. (2016). Graphical and numerical diagnostic tools to assess suitability of multiple imputations and imputation models. *Statistics in medicine*, 35(17):3007–3020.
- Edwards, A. W. F. (1984). *Likelihood*. CUP Archive.
- Fisher, R. A. (1956). *Statistical methods and scientific inference*. Oxford, England: Hafner Publishing Co.
- Fitzmaurice, G. M., Laird, N. M., and Ware, J. H. (2012). *Applied longitudinal analysis*, volume 998. John Wiley & Sons.
- Hoff, P. D. (2007). Extending the rank likelihood for semiparametric copula estimation. *The Annals of Applied Statistics*, pages 265–283.
- Hollenbach, M. F., Bojinov, I., Minhas, S., Metternich, N. W., Ward, M. D., and Volfovsky, A. (2017). Principled imputation made simple: Multiple imputation using gaussian copulas. Duke University, Working Paper.
- Little, R. J. (1988). A test of missing completely at random for multivariate data with missing values. *Journal of the American Statistical Association*, 83(404):1198–1202.
- Little, R. J. and Rubin, D. B. (2002). *Statistical analysis with missing data*. New York: Wiley, 2nd edition.
- Liu, Y. and Enders, C. K. (2017). Evaluation of multi-parameter test statistics for multiple imputation. *Multivariate Behavioral Research*, pages 1–20.

- Liublinska, V. and Rubin, D. B. (2014). Sensitivity analysis for a partially missing binary outcome in a two-arm randomized clinical trial. *Statistics in medicine*, 33(24):4170–4185.
- Lu, G. and Copas, J. B. (2004). Missing at random, likelihood ignorability and model completeness. *The Annals of Statistics*, 32(2):754–765.
- Marini, M. M., Olsen, A. R., and Rubin, D. B. (1980). Maximum-likelihood estimation in panel studies with missing data. *Sociological methodology*, 11:314–357.
- Mealli, F. and Rubin, D. B. (2015). Clarifying missing at random and related definitions, and implications when coupled with exchangeability. *Biometrika*, 102(4):995–1000.
- Meng, X.-L. and Rubin, D. B. (1992). Performing likelihood ratio tests with multiply-imputed data sets. *Biometrika*, 79(1):103–111.
- Molenberghs, G. and Kenward, M. G. (1998). Likelihood-based frequentist inference. *Missing Data in Clinical Studies*, pages 145–162.
- Nelson, R. B. (2006). *An Introduction to Copulas*, volume 2. Springer, New York, N.Y.
- Park, T. and Davis, C. S. (1993). A test of the missing data mechanism for repeated categorical data. *Biometrics*, pages 631–638.
- Potthoff, R. F., Tudor, G. E., Pieper, K. S., and Hasselblad, V. (2006). Can one assess whether missing data are missing at random in medical studies? *Statistical methods in medical research*, 15(3):213–234.
- Rosenbaum, P. R. and Rubin, D. B. (1984). Sensitivity of Bayes inference with data-dependent stopping rules. *The American Statistician*, 38(2):106–109.
- Rubin, D. B. (1976). Inference and missing data. *Biometrika*, 63(3):581–592.
- Rubin, D. B. (1987). *Multiple imputation for nonresponse in surveys*, volume 81. John Wiley & Sons.
- Seaman, S., Galati, J., Jackson, D., Carlin, J., et al. (2013). What is meant by “missing at random”? *Statistical Science*, 28(2):257–268.
- Sklar, M. (1959). *Fonctions de répartition à  $n$  dimensions et leurs marges*. Université Paris 8.

# Appendix 1

## Proofs of Corollary 1 and 2

of Corollary 1. By contradiction. Suppose that,

$$p(R_{i,j} = 0 \mid Y_i = y_i, \phi) \neq p(R_{i,j} = 0 \mid Y_{i,J^*+1:J} = y_{i,J^*+1:J}, \phi).$$

Then, there must exist a set  $S_j \subset \{1, \dots, J^*\}$  such that

$$p(R_{i,j} = 0 \mid Y_i = y_i, \phi) = p(R_{i,j} = 0 \mid Y_{i,S_j \cup J^*+1:J} = y_{i,S_j \cup J^*+1:J}, \phi),$$

and removing any element from  $S_j$  breaks the equality. Theorem 1 implies that for all  $k \in S_j$ ,

$$p(R_{i,k} = 0 \mid Y_i = y_i, \phi) = 0 \quad \text{for all } i \text{ and all } \phi,$$

which contradicts the assumption that the  $k^{\text{th}}$  column of  $Y$  has positive probability of having missing values.  $\square$

of Corollary 2. By Bayes theorem,

$$\begin{aligned} & p(Y_{\cdot j} = y_{\cdot j} \mid Y_{\cdot J^*+1:J} = y_{\cdot J^*+1:J}, R_{\cdot j'} = r_{\cdot j'}, \theta) \\ &= \frac{p(Y_{\cdot j} = y_{\cdot j} \mid Y_{\cdot J^*+1:J} = y_{\cdot J^*+1:J}, \theta) p(R_{\cdot j'} = r_{\cdot j'} \mid Y_{\cdot j} = y_{\cdot j}, Y_{\cdot J^*+1:J} = y_{\cdot J^*+1:J})}{p(R_{\cdot j'} = r_{\cdot j'} \mid Y_{\cdot J^*+1:J} = y_{\cdot J^*+1:J})} \\ &= \frac{p(Y_{\cdot j} = y_{\cdot j} \mid Y_{\cdot J^*+1:J} = y_{\cdot J^*+1:J}, \theta) \int p(R_{\cdot j'} = r_{\cdot j'} \mid Y_{\cdot j} = y_{\cdot j}, Y_{\cdot J^*+1:J} = y_{\cdot J^*+1:J}, \phi) p(\phi) d\phi}{p(R_{\cdot j'} = r_{\cdot j'} \mid Y_{\cdot J^*+1:J} = y_{\cdot J^*+1:J})} \\ &= \frac{p(Y_{\cdot j} = y_{\cdot j} \mid Y_{\cdot J^*+1:J} = y_{\cdot J^*+1:J}, \theta) \int p(R_{\cdot j'} = r_{\cdot j'} \mid Y_{\cdot J^*+1:J} = y_{\cdot J^*+1:J}, \phi) p(\phi) d\phi}{p(R_{\cdot j'} = r_{\cdot j'} \mid Y_{\cdot J^*+1:J} = y_{\cdot J^*+1:J})} \\ &= p(Y_{\cdot j} = y_{\cdot j} \mid Y_{\cdot J^*+1:J} = y_{\cdot J^*+1:J}, \theta). \end{aligned}$$

The third equality holds by Corollary 1.  $\square$

## Appendix 2

### Further discussion of the results from the simulation study

Table 5 shows the correct decision rates for  $R_{\cdot 1}$ ,  $R_{\cdot 2}$  and  $R_{\cdot 3}$  as well as the overall rate, for different sample sizes and the three different missingness mechanisms. Under the MAAR mechanism, directly testing a postulated missingness mechanism approach never rejected the null hypothesis; Liu and Enders (2017) similarly showed that the likelihood ratio test has low empirical Type I error for small samples with a significant proportion of missing data. The Gaussian copula and comparison of conditional means approaches had rejection rates close to the specified nominal level. For low values of  $N$ , the Gaussian copula and the comparison of conditional means approaches had

Table 4: Analysis of Variance of the 5-Factor Simulation Study

Source	Degrees of Freedom	Mean Square $\times 10^4$
Diagnostic Test	2	9465
$N$	2	60044
$\rho$	3	1657
$m$	2	1556
Mechanism	2	9187
Diagnostic Test $\times N$	4	2906
Diagnostic Test $\times \rho$	6	108
$N \times \rho$	6	183
Diagnostic Test $\times m$	4	428
$N \times m$	4	220
$\rho \times m$	6	42
Diagnostic Test $\times$ Mechanism	4	3766
$N \times$ Mechanism	4	5559
$\rho \times$ Mechanism	6	490
$m \times$ Mechanism	4	432
Diagnostic Test $\times N \times \rho$	12	394
Diagnostic Test $\times N \times m$	8	481
Diagnostic Test $\times \rho \times m$	12	54
$N \times \rho \times m$	12	48
Diagnostic Test $\times N \times$ Mechanism	8	644
Diagnostic Test $\times \rho \times$ Mechanism	12	33
$N \times \rho \times$ Mechanism	12	47
Diagnostic Test $\times m \times$ Mechanism	8	103
$N \times m \times$ Mechanism	8	67
$\rho \times m \times$ Mechanism	12	13
Diagnostic Test $\times N \times \rho \times m$	24	29
Diagnostic Test $\times N \times \rho \times$ Mechanism	24	104
Diagnostic Test $\times N \times m \times$ Mechanism	16	120
Diagnostic Test $\times \rho \times m \times$ Mechanism	24	14
$N \times \rho \times m \times$ Mechanism	24	17
Diagnostic Test $\times N \times \rho \times m \times$ Mechanism	48	11

Table 5: Correct choice rate averaged over the correlation and missingness factors.

Missingness Mechanism	Diagnostic Test	$N =$	$R_1$			$R_2$			$R_3$			Overall		
			100	500	1000	100	500	1000	100	500	1000	100	500	1000
MAAR	CCM		.98	.99	.98	.99	.99	.99	.98	.98	.98	.95	.96	.95
	DTPMM		1	1	1	1	1	1	1	1	1	1	1	1
	GC		.99	.99	.98	.99	.99	.98	.99	.99	.98	.98	.96	.95
MAAR2	CCM		.98	.98	.98	.17	.88	.99	.32	.98	1	.42	.99	1
	DTPMM		1	1	1	0	.24	.66	.01	.68	.94	.01	.69	.94
	GC		.98	.97	.95	.28	.94	1	.48	.98	1	.57	.99	1
MNAAR	CCM		.98	.98	.98	.1	.68	.94	.36	.99	1	.42	.99	1
	DTPMM		1	1	1	0	.16	.55	.02	.77	.97	.02	.77	.97
	GC		.99	.98	.98	.13	.78	.96	.49	.99	1	.55	.99	1

higher rejection rates for the MNAAR and MAAR2 mechanisms than directly testing a postulated missingness mechanism; however, as the sample size increased, the difference decreased.

All three diagnostic tests were able to identify the variables that were violating the MAAR/conditional independence assumption, although they needed a larger sample size to detect violations in  $R_2$  relative to  $R_3$ . The correct decision rates under the MAAR2 and MNAAR missingness mechanisms are similar because there is no information in the data that distinguishes between them.