# Estimation and imputation in Probabilistic Principal Component Analysis with Missing Not At Random data

**Aude Sportisse**
Sorbonne University
Paris, France
`aude.sportisse@sorbonne-universite.fr`

**Claire Boyer**
Sorbonne University
Paris, France
`claire.boyer@sorbonne-universite.fr`

**Julie Josse**
INRIA
Montpellier, France
`julie.josse@inria.fr`

## Abstract

Missing Not At Random (MNAR) values where the probability of having missing data may depend on the missing value itself, are notoriously difficult to account for in analyses, although very frequent in the data. One solution to handle MNAR data is to specify a model for the missing data mechanism, which makes inference or imputation tasks more complex. Furthermore, this implies a strong *a priori* on the parametric form of the distribution. However, some works have obtained guarantees on the estimation of parameters in the presence of MNAR data, without specifying the distribution of missing data [18, 25]. This is very useful in practice, but is limited to simple cases such as few self-masked MNAR variables in data generated according to linear regression models. We continue this line of research, but extend it to a more general MNAR mechanism, in a more general model of the probabilistic principal component analysis (PPCA), *i.e.*, a low-rank model with random effects. We prove identifiability of the PPCA parameters. We then propose an estimation of the loading coefficients, and a data imputation method. Both are based on estimators of means, variances and covariances of missing variables, for which consistency is discussed. These estimators have the great advantage of being calculated using only the observed information, leveraging the underlying low-rank structure of the data. We illustrate the relevance of the method with numerical experiments on synthetic data and also on two datasets, one collected from a medical register and the other one from a recommendation system.

## 1 Introduction

The problem of missing data is ubiquitous in the practice of data analysis. Theoretical guarantees of estimation strategies or imputation methods rely on assumptions regarding the missing-data mechanism, *i.e.* the cause of the lack of data. Rubin [22] introduced three missing-data mechanisms. The data are said (i) Missing Completely At Random (MCAR) if the probability of being missing does not depend on any values observed or missing, (ii) Missing At Random (MAR) if the probability of being missing only depends on observed values, (iii) Missing Not At Random (MNAR) if the unavailability of the data may depend on both observed and unobserved data such as its value itself. We focus on this later case, which is frequent in practice, and theoretically challenging. A classic example of MNAR data is surveys about salary for which rich people would be less willing to disclose their income.

When the data is MCAR or MAR, statistical inference is carried out by ignoring the missing data mechanism [10]. In the MNAR case, the observed data are no longer representative of the population, which leads to selection bias in the sample, and therefore to bias in the parameters estimation when using for instance complete case analysis. One solution to handle MNAR data, known as *selection model* [10], is to model missing data distribution; most of the time, by logistic regression models

[4, 19, 24]. This comes at the price of an important computational burden to perform inference and is often restricted to a limited number of MNAR variables. In the recommender system community, some authors [12, 3, 11, 28] suggest that not MCAR values can be handled using a joint modelling of the data and mechanism distributions by matrix factorization; then they debias existing methods for MCAR data, for instance with inverse probability weighting approaches.

In addition, a key issue of MNAR data is to establish identifiability, which is not always guaranteed [15]. The literature on this topic is abundant, both in the non-parametric [17, 16, 6, 23, 20], and semi-parametric settings [27, 14]. For parametric models, in the case of multivariate regression, Tang et al. [25] and Miao et al. [15] guarantee the identifiability of the coefficients of the conditional distribution of $Y|X$, when $Y$ is missing. Tang et al. [25] estimate them by calculating the coefficient of the distributions of $X$ and $X|Y$ using only observations with no missing values. Besides, in a linear model with self-masked missing mechanism, *i.e.*, the lack depends only on the missing variable itself, Mohan et al. [18] consider a related approach based on graphical models, adopting a causal point of view. Despite the great advantage of not modeling the distribution of missing values, the assumption of a self-masked MNAR mechanism and the restriction to a linear model are yet strong.

**Contributions.** We consider a framework where the data are generated according to a probabilistic principal components analysis (PPCA) [26] model. Contrary to available works that handle only MAR data in PPCA [5], we consider that the missing values mechanism can be MNAR (on several variables) and we also consider the possibility of having different mechanisms in the same data (MNAR and M(C)AR).

- We prove the identifiability of the PPCA model parameters in a self-masked MNAR values setting encompassing a large set of self-masked mechanism distributions.

- For more general MNAR mechanism, we give a strategy to estimate the PPCA loading parameters without any modeling of the missing-data mechanism and use it to impute missing values.

- The proposed method is based on estimators for the mean, the variance and the covariance of the variables with MNAR values. We show that they can be consistently estimated. Two strategies lead to the proposed estimators: (i) the first one uses algebraic arguments based on partial linear models derived from the PPCA model; (ii) the second one is inspired by [18] and uses graphical models and in particular the so-called missingness graph.

- We derive an algorithm implementing our proposal. We show that it outperforms the state-of-the-art methods on synthetic data and on two real datasets, collected from a medical registry (Traumabase®) and from a joke recommender system (the Jester Online Joke Recommender System [2]). The code to reproduce all the simulations and the numerical experiments is available at `https://github.com/AudeSportisse/PPCA_MNAR`.

## 2 PPCA model with informative missing values: identifiability issues

**Setting.** The data matrix $Y \in \mathbb{R}^{n \times p}$ is assumed to be generated under a fully-connected PPCA model [26] (a.k.a. a low-rank model with random effects), *i.e.* by the factorization of the loading matrix $B \in \mathbb{R}^{r \times p}$ and $r$ latent variables grouped in the matrix $W \in \mathbb{R}^{n \times r}$,

$$Y = \mathbf{1}\alpha + WB + \epsilon, \text{ with } \begin{cases} W = (W_{1.}|\dots|W_{n.})^T, \text{ with } W_{i.} \sim \mathcal{N}(0_r, \mathrm{Id}_{r \times r}) \in \mathbb{R}^r, \\ B \text{ of rank } r < \min\{n, p\}, \\ \alpha \in \mathbb{R}^p \text{ and } \mathbf{1} = (1 \dots 1)^T \in \mathbb{R}^n, \\ \epsilon = (\epsilon_{1.}|\dots|\epsilon_{n.})^T, \text{ with } \epsilon_{i.} \sim \mathcal{N}(0_p, \sigma^2 \mathrm{Id}_{p \times p}) \in \mathbb{R}^p, \end{cases} \tag{1}$$

for $\sigma^2$ and $r$ known. In the sequel, $Y_{.j}$ and $Y_{i.}$ respectively denote the column $j$ and the row $i$ of $Y$. The rows of $Y$ are identically distributed, $\forall i \in \{1, \dots, n\}$, $Y_{i.} \sim \mathcal{N}(\alpha, \Sigma)$, with $\Sigma = B^T B + \sigma^2 \mathrm{Id}_{p \times p}$. We denote $\Omega \in \{0, 1\}^{n \times p}$ the missing-data pattern (or mask) defined as follows:

$$\forall i \in \{1, \dots, n\}, \ \forall j \in \{1, \dots, p\}, \quad \Omega_{ij} = \begin{cases} 0 & \text{if } Y_{ij} \text{ is missing}, \\ 1 & \text{otherwise}. \end{cases} \tag{2}$$

Some variables $Y_{.m_1}, \dots, Y_{.m_d}$, indexed by $\mathcal{M} := \{m_1, \dots, m_d\} \subset \{1, \dots, p\}$ (with $d < p$), contain MNAR values. The other variables are considered to be observed (or M(C)AR see Appendix B.5). We define a general MNAR mechanism where the probability to have missing values may

depend on the $d$ MNAR variables but also on $p - d - r$ other variables that can be observed or M(C)AR[1]. The remaining $r$ variables are called <u>pivot variables</u> and can be observed or MCAR. More precisely, we denote the complementary of a set $\mathcal{A}$ as $\overline{\mathcal{A}} := \{1, \dots, p\} \backslash \mathcal{A}$. The general MNAR mechanism is defined as follows, with $\mathcal{J} \subset \overline{\mathcal{M}}$ the set of indices of the $r$ pivot variables ($|\mathcal{J}| = r$),

$$\forall m \in \mathcal{M}, \forall i \in \{1, \dots, n\}, \quad \mathbb{P}(\Omega_{im} = 1 | Y_{i.}) = \mathbb{P}(\Omega_{im} = 1 | (Y_{ik})_{k \in \bar{\mathcal{J}}}). \tag{3}$$

We also define a specific MNAR mechanism, called the self-masked MNAR mechanism as follows. We assume that $d$ variables are self-masked MNAR indexed by $\mathcal{M}$ and the $p - d$ other variables are MCAR (or observed), indexed by $\overline{\mathcal{M}}$, i.e, $\forall i \in \{1, \dots, n\}$,

$$\forall m \in \mathcal{M}, \quad \mathbb{P}(\Omega_{im} = 1 | Y_{i.}) = \mathbb{P}(\Omega_{im} = 1 | Y_{im}). \tag{4}$$

**Model identifiability.** We prove the identifiability of the PPCA model (see Appendix A for the complete proof), *i.e.* the joint distribution of $Y$ can be uniquely determined from the available information, in the self-masked missing values case. More particularly, assume the following

**A01.** $d$ variables are self-masked MNAR as in (4) and the $p - d$ other variables are MCAR (or observed). The missing-data distributions $(F_m)_{m \in \mathcal{M}}$ and $(F_j)_{j \in \overline{\mathcal{M}}}$ are known strictly monotone functions with a finite support, defined as follows, $\forall i \in \{1, \dots, n\}$,

$$\forall m \in \mathcal{M}, \quad \mathbb{P}(\Omega_{im} = 1 | Y_{i.}) = F_m(\phi_m^0 + \phi_m^1 Y_{im}),$$
$$\forall j \in \overline{\mathcal{M}}, \quad \mathbb{P}(\Omega_{ij} = 1 | Y_{i.}) = \mathbb{P}(\Omega_{ij} = 1) = F_j(\phi_j),$$

with $\phi_j \in \mathbb{R}$ and $\phi_m = (\phi_m^0, \phi_m^1) \in \mathbb{R}^2$ the mechanism parameters.

**A02.** $\forall (k, \ell) \in \{1, \dots, p\}^2, \quad k \neq \ell, \qquad \Omega_{.k} \perp\!\!\!\perp \Omega_{.\ell} | Y$

Note that under Assumption **A01.**, any function $F_m, m \in \mathcal{M}$ can be considered, as a logistic function while [15] presented many counterexamples when identification fails considering the logistic distribution. **A02.** requires that the missing-data patterns are independent conditionally to the data.

**Proposition 1.** *Under Assumptions **A01.** and **A02.**, the parameters $(\alpha, \Sigma)$ of the PPCA model (1) and the mechanism parameters $\phi = (\phi_\ell)_{\ell \in \{1, \dots p\}}$ are identifiable. Assuming that the noise level $\sigma^2$ is known, the parameter $B$ is identifiable up to a row permutation.*

## 3 Estimators with theoretical guarantees

In this section, we provide estimators of the means, variances and covariances for the MNAR variables, when data are generated under the PPCA model described in (1). These estimators are used to derive an estimator of the loading matrix $B$ in (1). This makes it possible to derive a new imputation method with MNAR data as detailed in Algorithm 1.

We denote $\mathcal{J}_{-j} := \mathcal{J} \backslash \{j\}$ and assume

**A1.** $\forall m \in \mathcal{M}, \forall j \in \mathcal{J}, \; \left( B_{.m} \quad (B_{.j'})_{j' \in \mathcal{J}_{-j}} \right)$ is invertible,

**A2.** $\forall m \in \mathcal{M}, \forall j \in \mathcal{J}, \; Y_{.j} \perp\!\!\!\perp \Omega_{.m} | (Y_{.k})_{k \in \overline{\{j\}}}$.

Note that Assumption **A1.** implies that $B$ has a full rank $r$ and that any variable in $Y$ is generated by all the latent variables[2] (named a "fully-connected" PPCA). Assumption **A2.** is implied by the general MNAR mechanism in (3).

We start by illustrating the methodology and the assumptions using an example in small dimension, before turning to the general case.

### 3.1 Estimation of the mean of a MNAR variable

Consider a toy dataset where $p = 3, r = 2$, in which only one variable is missing, $\mathcal{M} = \{1\}$ and there are two pivots variables $\mathcal{J} = \{2, 3\}$. Note that the MNAR mechanism is self-masked in such a context, because Equation (3) leads to $\mathbb{P}(\Omega_{.1} = 1 | Y_{.1}, Y_{.2}, Y_{.3}) = \mathbb{P}(\Omega_{.1} = 1 | Y_{.1})$, but the method can be extended to more general cases. Our aim is to estimate the mean of $Y_{.1}$, without specifying the distribution of the missing-data mechanism.

---

[1] Note that it implies that $d < p - r$.

[2] It does not require that the linear combination coefficients are non-zero.

**Using algebraic arguments.** We proceed in three steps: (i) **A1.** allows to obtain linear link between the pivot variables $(Y_{.2}, Y_{.3})$ and the MNAR variable $Y_{.1}$. For instance,

$$Y_{.2} = \mathcal{B}_{2\to1,3[0]} + \mathcal{B}_{2\to1,3[1]}Y_{.1} + \mathcal{B}_{2\to1,3[3]}Y_{.3} + \zeta, \qquad (5)$$

with $\zeta$ a noise term, $\mathcal{B}_{2\to1,3[0]}$, $\mathcal{B}_{2\to1,3[1]}$ and $\mathcal{B}_{2\to1,3[3]}$ the intercept and the coefficients in the model (the arrow $2 \to 1, 3$ indicates the regression model of $Y_{.2}$ on $Y_{.1}$ and $Y_{.3}$, while the squared bracket represents the coefficient, for instance 3 for the coefficient of $Y_{.3}$) ; (ii) Assumption **A2.**, *i.e.* $Y_{.2} \perp\!\!\!\perp \Omega_{.1}|Y_{.1}, Y_{.3}$, is required to obtain identifiable and consistent parameters of the distribution of $Y_{.2}$ given $Y_{.1}, Y_{.3}$ in the complete-case when $\Omega_{.1} = 1$, denoted as $\mathcal{B}^c_{2\to1,3[0]}$, $\mathcal{B}^c_{2\to1,3[1]}$ and $\mathcal{B}^c_{2\to1,3[3]}$,

$$(Y_{.2})_{|\Omega_{.1}=1} = \mathcal{B}^c_{2\to1,3[0]} + \mathcal{B}^c_{2\to1,3[1]}Y_{.1} + \mathcal{B}^c_{2\to1,3[3]}Y_{.3} + \zeta^c, \qquad (6)$$

(note that the regression of $Y_{.1}$ on $(Y_{.2}, Y_{.3})$ is prohibited, as **A2.** does not hold); (iii) using again **A2.**,

$$\mathbb{E}\left[Y_{.2}|Y_{.1}, Y_{.3}, \Omega_{.1} = 1\right] = \mathbb{E}\left[\mathcal{B}^c_{2\to1,3[0]} + \mathcal{B}^c_{2\to1,3[1]}Y_{.1} + \mathcal{B}^c_{2\to1,3[3]}Y_{.3}|Y_{.1}, Y_{.3}\right],$$

and taking the expectation leads to

$$\mathbb{E}\left[Y_{.2}\right] = \mathcal{B}^c_{2\to1,3[0]} + \mathcal{B}^c_{2\to1,3[1]}\mathbb{E}\left[Y_{.1}\right] + \mathcal{B}^c_{2\to1,3[3]}\mathbb{E}\left[Y_{.3}\right].$$

The latter expression can be reshuffled so that the expectation of $Y_{.1}$ can be estimated: the means of $Y_{.2}$ and $Y_{.3}$ are estimated by standard empirical estimators (it will be Assumption **A4.** in the sequel).

**Using graphical arguments.** The PPCA model can be represented with structural causal graphs [21], as illustrated in Figure 1. The top left graph in which each variable is generated by a combination of all latent variables, see Assumption **A1.**, can be represented as the top right one, as $Y_{.1} \leftarrow W_{.1} \to Y_{.2}$ is equivalent to $Y_{.1} \leftrightarrow Y_{.2}$ (see [21, page 52]). Then, six reduced graphical models can be derived from the top right graph (two instances are represented in the bottom). Indeed, a bidirected edge $Y_{.1} \leftrightarrow Y_{.2}$ can be interchanged (see [21, rule 1, page 147]) with an oriented edge $Y_{.1} \to Y_{.2}$, if each neighbor of $Y_{.2}$ (*i.e.* $Y_{.1}$ or $Y_{.3}$) is inseparable of $Y_{.1}$ (see [21, page 17]). The bottom left graph can also be represented by Equation (6), which gives a connection between the algebraic and graphical approaches.
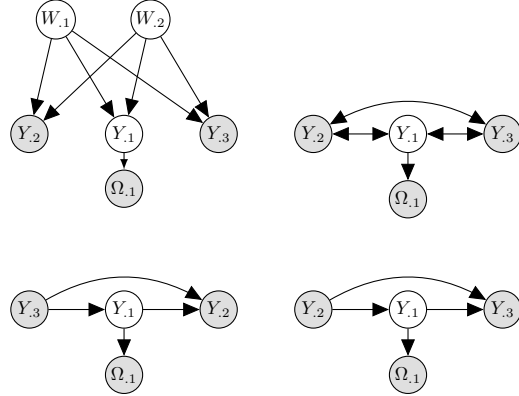


Figure 1: Graphical models for the toy example with one missing variable $Y_{.1}$, $p = 3$ and $r = 2$.

## 3.2 Estimation of the mean, variance and covariances of the MNAR variables

In a general case, estimators of the mean, variance and covariances of the variables with MNAR values can be computed one by one. We detail the results only for one variable $Y_{.m}, m \in \mathcal{M}$, but the results hold for several variables with MNAR values. In addition, the other variables are considered to be observed for simplicity but they could contain MCAR and MAR values as well, as explained in Appendix B.5. We adopt the algebraic strategy here to derive estimators (see Appendix B for proofs) but graphical arguments can also be used to obtain similar results (see Appendix F). The starting point is to exploit the linear links between variables, as described in the next lemma.

**Lemma 2.** *Under the PPCA model* (1) *and Assumption* **A1.***, choose* $j \in \mathcal{J}$. *One has*

$$Y_{.j} = \mathcal{B}_{j\to m,\mathcal{J}_{-j}[0]} + \sum_{j'\in\mathcal{J}_{-j}} \mathcal{B}_{j\to m,\mathcal{J}_{-j}[j']}Y_{.j'} + \mathcal{B}_{j\to m,\mathcal{J}_{-j}[m]}Y_{.m} + \zeta, \qquad (7)$$

*where* $\zeta = -\sum_{j'\in\mathcal{J}_{-j}} \mathcal{B}_{j\to m,\mathcal{J}_{-j}[j']}\epsilon_{.j'} - \mathcal{B}_{j\to m,\mathcal{J}_{-j}[m]}\epsilon_{.m} + \epsilon_{.j}$. *is a noise term.*

$\mathcal{B}_{j\to m,\mathcal{J}_{-j}[0]}$, $\mathcal{B}_{j\to m,\mathcal{J}_{-j}[j']}$ *and* $\mathcal{B}_{j\to m,\mathcal{J}_{-j}[m]}$ *are given in Appendix B.1 and depend on the coefficients of* $B$ *given in* (1).

4

Then we define the regression coefficients of $Y_{.j}$ on $Y_{.m}$ and $Y_{.k}$, for $k \in \mathcal{J}_{-j}$ in the complete case, that will be used to express the mean of a variable with MNAR values.

**Definition 3** (Coefficients in the complete case). *For $j \in \mathcal{J}$ and $k \in \mathcal{J}_{-j}$, let $\mathcal{B}^c_{j \to m, \mathcal{J}_{-j}[0]}$, $\mathcal{B}^c_{j \to m, \mathcal{J}_{-j}[m]}$ and $\mathcal{B}^c_{j \to m, \mathcal{J}_{-j}[j']}$ be respectively the intercept and the coefficients standing for the effects of $Y_{.j}$ on $(Y_{.m}, (Y_{.j'})_{j' \in \mathcal{J}_{-j}})$ in the complete case, i.e. when $\Omega_{.m} = 1$:*

$$(Y_{.j})_{|\Omega_{.m}=1} := \mathcal{B}^c_{j \to m, \mathcal{J}_{-j}[0]} + \sum_{j' \in \mathcal{J}_{-j}} \mathcal{B}^c_{j \to m, \mathcal{J}_{-j}[j']} Y_{.j'} + \mathcal{B}^c_{j \to m, \mathcal{J}_{-j}[m]} Y_{.m} + \zeta^c, \qquad (8)$$

*with $\zeta^c = -\sum_{j' \in \mathcal{J}_{-j}} \mathcal{B}^c_{j \to m, \mathcal{J}_{-j}[j']} \epsilon_{.j'} - \mathcal{B}^c_{j \to m, \mathcal{J}_{-j}[m]} \epsilon_{.m} + \epsilon_{.j}$.*

Then, we make the two following assumptions:

**A3.** For all $j \in \mathcal{J}$, for all $m \in \mathcal{M}$, the complete-case coefficients $\mathcal{B}^c_{j \to m, \mathcal{J}_{-j}[0]}$, $\mathcal{B}^c_{j \to m, \mathcal{J}_{-j}[m]}$ and $\mathcal{B}^c_{j \to m, \mathcal{J}_{-j}[k]}$, $k \in \mathcal{J}_{-j}$ can be consistently estimated.

**A4.** The means $(\alpha_j)_{j \in \mathcal{J}}$, variances $(\mathrm{Var}(Y_{.j}))_{j \in \mathcal{J}}$ and covariances $(\mathrm{Cov}(Y_{.j}, Y_{.j'}))_{j \in \mathcal{J}, j' \in \mathcal{J}_{-j}}$ of the $r$ pivot variables can be consistently estimated.

Note that Assumption **A4.** is met whether the $r$ pivot variables are fully observed.

**Proposition 4** (Mean estimator). *Consider the PPCA model* (1). *Under Assumptions **A1.** and **A2.**, an estimator of the mean of a MNAR variable $Y_{.m}$, for $m \in \mathcal{M}$, can be constructed as follows: choose $j \in \mathcal{J}$, and compute*

$$\hat{\alpha}_m := \frac{\hat{\alpha}_j - \hat{\mathcal{B}}^c_{j \to m, \mathcal{J}_{-j}[0]} - \sum_{j' \in \mathcal{J}_{-j}} \hat{\mathcal{B}}^c_{j \to m, \mathcal{J}_{-j}[j']} \hat{\alpha}_{j'}}{\hat{\mathcal{B}}^c_{j \to m, \mathcal{J}_{-j}[m]}}, \qquad (9)$$

*with $(\hat{\mathcal{B}}^c_{j \to m, \mathcal{J}_{-j}[k]})_{k \in \{0,m\} \cup \mathcal{J}_{-j}}$ estimators of the coefficients obtained from Definition 3.*

*Under the additional Assumptions **A3.** and **A4.**, this estimator is consistent.*

The proof is given in Appendix B.2. Proposition 4 provides an estimator easily computable from all observed cells. Furthermore, different choices of $Y_{.j}$, $j \in \mathcal{J}$ can be done in Equation (9) and all the resulting estimators may be aggregated to stabilize the estimation of $\alpha_m$.

**Proposition 5** (Variance and covariances estimators). *Consider the PPCA model* (1). *Under Assumptions **A1.** and **A2.**, an estimator of the variance of a MNAR variable $Y_{.m}$, for $m \in \mathcal{M}$, and its covariances with the pivot variables, can be constructed as follows: choose a pivot variable $Y_{.j}$ for $j \in \mathcal{J}$ and compute*

$$\left( \widehat{\mathrm{Var}}(Y_{.m}) \quad \widehat{\mathrm{Cov}}(Y_{.m}, (Y_{.j'})_{j' \in \mathcal{J}}) \right)^T := (\widehat{M_j})^{-1} \hat{P}_j, \qquad (10)$$

*assuming that $\sigma^2$ tends to zero, with $\widehat{M_j}^{-1} \in \mathbb{R}^{(r+1) \times (r+1)}$, $\hat{P}_j \in \mathbb{R}^{r+1}$ detailed in Appendix B.3. These quantities depend on $(\hat{\alpha}_{j'})_{j' \in \mathcal{J}}$, $\hat{\alpha}_m$ given in Proposition 4, on $(\widehat{\mathrm{Var}}(Y_{.j'}))_{j' \in \mathcal{J}}$ and on complete-case coefficients such as $(\hat{\mathcal{B}}^c_{j' \to m, \mathcal{J}_{-j'}[k]})_{k \in \{m\} \cup \mathcal{J}_{-j'}}$ for $j' \in \mathcal{J}$.*

*Under the additional Assumptions **A3.** and **A4.**, the estimators of the variance of $Y_{.m}$ and its covariances with the pivot variables given in* (10) *are consistent.*

The proof is given in Appendix B.3. Note that to estimate the variance of a MNAR variable, only $r$ pivot variables are required to solve (10) and $r$ tasks have to be performed for estimating the coefficients of the effects of $Y_{.k}$ on $(Y_{.\ell})_{\ell \in \{m\} \cup \mathcal{J}_{-k}}$ for all $k \in \mathcal{J}$.

All the ingredients can be combined to form an estimator $\hat{\Sigma}$ for the covariance matrix $\Sigma$. Define

$$\hat{\Sigma} := \left( \widehat{\mathrm{Cov}}(Y_{.k}, Y_{.\ell}) \right)_{k, \ell \in \{1, \dots, p\}}, \qquad (11)$$

- if $Y_{.k}$ and $Y_{.\ell}$ have both consistent mean/variance estimators, then $\widehat{\mathrm{Cov}}(Y_{.k}, Y_{.\ell})$ can be trivially evaluated by standard empirical covariance estimators.

- if $Y_{.k}$ is a MNAR variable and $Y_{.\ell}$ is a pivot variable, then $\widehat{\mathrm{Cov}}(Y_{.k}, Y_{.\ell})$ is given by (10),
- if $Y_{.k}$ is a MNAR variable and $Y_{.\ell}$ is not a pivot variable, *i.e.* $\ell \in \bar{\mathcal{J}} \backslash \{k\}$, a similar strategy as the one above can be devised. Then $\widehat{\mathrm{Cov}}(Y_{.k}, Y_{.\ell})$ is given by (48) detailed in Appendix B.4 and for which some additional assumptions similar as the ones above are required. This estimator relies on the choice of $r-1$ pivot variables indexed by $j$ and $\mathcal{H} \subset \mathcal{J}$, and only necessitates to evaluate the effects of $Y_{.j}$ on $(Y_{.j'})_{j' \in \{k, \ell\} \cup \mathcal{H}}$ in the complete case.

### 3.3 Performing PPCA with MNAR variables

With the estimator $\hat{\Sigma}$ in (11) at hand, one can perform the estimation of the loading matrix $B$ in (1).

**Definition 6** (Estimation of the loading matrix). *Given the estimator $\hat{\Sigma}$ of the covariance matrix in (11), let the orthogonal matrix $\hat{U} = (\hat{u}_1 | \ldots | \hat{u}_p) \in \mathbb{R}^{p \times p}$ and the diagonal matrix $\hat{D} = \mathrm{diag}(\hat{d}_1, \hat{d}_2, \ldots, \hat{d}_p) \in \mathbb{R}^{p \times p}$ with $\hat{d}_1 \geqslant \hat{d}_2 \geqslant \ldots \geqslant \hat{d}_p \geqslant 0$ form the singular value decomposition of the following matrix $\hat{\Sigma} - \sigma^2 \mathrm{Id}_{p \times p} =: \hat{U} \hat{D} \hat{U}^T$. An estimator $\hat{B}$ of $B$ can be defined using the $r$ first singular values and vectors, as follows*

$$\hat{B} = \hat{D}_{|r}^{1/2} \hat{U}_{|r}^T = \mathrm{diag}(\hat{d}_1, \ldots, \hat{d}_r)^{1/2} (\hat{u}_1^T | \ldots | \hat{u}_r^T)^T \tag{12}$$

The estimation of the loading matrix is used to impute the variables with missing values. More precisely, a classical strategy to impute missing values is to estimate their conditional expectation given the observed values. One can note that with $\Sigma = B^T B + \sigma^2 \mathrm{Id}_{p \times p}$, the conditional expectation of $Y_{.m}$ for $m \in \mathcal{M}$ given $(Y_{.k})_{k \in \overline{\mathcal{M}}}$ reads as follows

$$\mathbb{E}[Y_{.m} | (Y_{.k})_{k \in \overline{\mathcal{M}}}] = \alpha_m + \Sigma_{m, \overline{\mathcal{M}}} \Sigma_{\overline{\mathcal{M}}, \overline{\mathcal{M}}}^{-1} \left( Y_{.\overline{\mathcal{M}}}^T - \alpha_{\overline{\mathcal{M}}} \right),$$

with $\Sigma_{m, \overline{\mathcal{M}}} := (\Sigma_{m,k})_{k \in \overline{\mathcal{M}}}^T$, $\Sigma_{\overline{\mathcal{M}}, \overline{\mathcal{M}}} := (\Sigma_{k,k'})_{k,k' \in \overline{\mathcal{M}}}$, $Y_{.\overline{\mathcal{M}}} := (Y_{.k})_{k \in \overline{\mathcal{M}}}$, and $\alpha_{\overline{\mathcal{M}}} := (\alpha_k)_{k \in \overline{\mathcal{M}}}$.

**Definition 7** (Imputation of a MNAR variable). *Set $\hat{\Gamma} := \hat{B}^T \hat{B} + \sigma^2 \mathrm{Id}_{p \times p}$ for $\hat{B}$ given in Definition 6. The MNAR variable $Y_{.m}$ with $m \in \mathcal{M}$ can be imputed as follows: for $i$ such that $\Omega_{i,m} = 0$,*

$$\hat{Y}_{im} = \hat{\alpha}_m + \hat{\Gamma}_{m, \overline{\mathcal{M}}} \hat{\Gamma}_{\overline{\mathcal{M}}, \overline{\mathcal{M}}}^{-1} \left( Y_{i, \overline{\mathcal{M}}}^T - \hat{\alpha}_{\overline{\mathcal{M}}} \right) \tag{13}$$

*with $\hat{\Gamma}_{m, \overline{\mathcal{M}}} := (\hat{\Gamma}_{m,k})_{k \in \overline{\mathcal{M}}}^T$, $\hat{\Gamma}_{\overline{\mathcal{M}}, \overline{\mathcal{M}}} := (\hat{\Gamma}_{k,k'})_{k,k' \in \overline{\mathcal{M}}}$, $Y_{.\overline{\mathcal{M}}} := (Y_{.k})_{k \in \overline{\mathcal{M}}}$ and $\hat{\alpha}_{\overline{\mathcal{M}}} := (\hat{\alpha}_k)_{k \in \overline{\mathcal{M}}}$.*

### 3.4 Algorithm

The proposed imputation method described in Algorithm 1 can handle the different MNAR mechanisms, the self-masked MNAR case and the general MNAR cases where the probability to have missing values on variables depends on both the underlying values and values of other variables (observed or missing).

---
**Algorithm 1** PPCA with MNAR variables.

---
**Require:** $r$ (number of latent variables), $\sigma^2$ (noise level), $\mathcal{J}$ (pivot variables indices), $\Omega$ (mask).

1: **for** each MNAR variable $(Y_{.m})_{m \in \mathcal{M}}$ **do**
2:      Evaluate $\hat{\alpha}_m$ the estimator of its mean given in (9) using the $r$ pivot variables indexed by $\mathcal{J}$.
3:      Evaluate $\widehat{\mathrm{Var}}(Y_{.m})$, and $\widehat{\mathrm{Cov}}(Y_{.m}, Y_{.\ell})$ for $\ell \in \mathcal{J}$, using (10).
4:      Evaluate $\widehat{\mathrm{Cov}}(Y_{.m}, Y_{.\ell})$ for $\ell \in \bar{\mathcal{J}} \backslash \{m\}$ using Proposition 8.
5: **end for**
6: Form $\hat{\Sigma}$, covariance matrix estimator in (11).
7: Compute the loading matrix estimator $\hat{B}$ given in (12).
8: Compute $\hat{\Gamma} = \hat{B}^T \hat{B} + \sigma^2 \mathrm{Id}_{p \times p}$.
9: **for** each missing variable $(Y_{.j})$ **do**
10:      **for** $i$ such that $\Omega_{ij} = 0$ **do**
11:          $\hat{Y}_{ij} \leftarrow$ Impute $Y_{ij}$ as in (13).
12:      **end for**
13: **end for**

---

Algorithm 1 requires the set $\mathcal{J}$, *i.e.* the selection of $r$ pivot variables on which the regressions in Propositions 4, 5 and 8 will be performed. If there are more than $r$ variables that can be pivot, we suggest selecting a bigger set ($> r$) and computing the final estimator with the median of the estimators over all possible combinations. The efficiency of this strategy is illustrated in Appendix C.

The estimators associated to any missing variable in the steps 1 to 5 are computed in the complete case, i.e. with the rows for which the missing variable is observed. When the pivot variables are also missing, the complete case corresponds to discarding all rows where the pivot variables or the MNAR one are missing and not all rows containing missing values. This could be problematic in the high-dimensional setting, but here the low-rank assumption ($r < \min\{n, p\}$) ensures that the number of pivot variables is small enough, so that the complete case analysis will not result in discarding many rows of the dataset.

In order to estimate the coefficients in Definition 3, we use ordinary least squares despite that the exogeneity assumption, *i.e.* the noise term is independent of the covariates, does not hold. It still leads to accurate estimation in numerical experiments as shown in Section 4. Actually, the consistency required by Assumption **A3.** holds as the variance of the noise tends to 0.

# 4 Numerical experiments

## 4.1 Synthetic data

We empirically compare Algorithm 1 (**MNAR**) to the state-of-the-art methods, including

  (i) **MAR**: our method which has been adapted to handle MAR data (inspired by [18, Theorems 1, 2, 3] in linear models), see Appendix G for details;

 (ii) **EMMAR**: which consists in an EM algorithm to perform PPCA with MAR values [5];

(iii) **SoftMAR**: a matrix completion method using an iterative soft-thresholding singular value decomposition algorithm [13] relevant only for M(C)AR values;

(iv) **MNARparam**: a matrix completion technique modeling the MNAR mechanism with a parametric logistic model [24].

Note that Method (ii) is specially designed to estimate the PPCA loading matrix and not to perform imputation, but this is possible combining Method (ii) with steps 8 and 9 in Algorithm 1. This is the other way around for completion Methods (iii) and (iv), but the loading matrix can be computed as in (12). Note also that Methods (iii) and (iv) are developed in a context of low-rank models with fixed effects. They require tuning a regularization parameter $\lambda$: we consider an oracle value minimizing the true imputation error. We also use oracle values for the noise level and the rank in Algorithm 1. These methods are compared with the imputation by the mean (**Mean**), which serves as a benchmark, and the naive listwise deletion method (**Del**) which consists in estimating the parameters empirically with the fully-observed data only. A comparison of the methods in terms of computational times is given in Appendix D.

**Measuring the performance.** For the loading matrix, the RV coefficient [8], which is a measure of relationship between two random vectors, is computed between the estimate $\hat{B}$ and the true $B$. An RV coefficient close to one means high correlation between the image spaces of $\hat{B}$ and $B$. Denoting the Frobenius norm as $\|.\|_F$, the quality of imputation is measured with the normalized imputation error given by $\|(\hat{Y} - Y) \odot (1 - \Omega)\|_F^2 / \|Y \odot (1 - \Omega)\|_F^2$.

**Setting.** We generate a data matrix of size $n = 1000$ and $p = 10$ from a PPCA model (1) with two latent variables ($r = 2$) and with a noise level $\sigma = 0.1$. Missing values are introduced on seven variables $(Y_{.k})_{k \in [1:7]}$ according to a logistic self-masked MNAR mechanism, leading to 35% of missing values in total. Results are presented[3] for one missing variable $Y_{.1}$ (same results hold for other missing variables). All the observed variables $(Y_{.k})_{k \in [8:10]}$ are considered to be pivot.
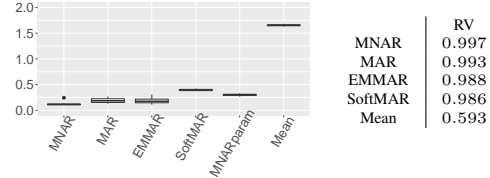


| | RV |
|---|---|
| MNAR | 0.997 |
| MAR | 0.993 |
| EMMAR | 0.988 |
| SoftMAR | 0.986 |
| Mean | 0.593 |

Figure 2: Imputation error (left) and median of the RV coefficients for the loading matrix (right).

---

[3]For a given set of PPCA parameters, the stochasticity comes from the process of drawing 20 times the latent variables, the additive noise and the missing-data pattern.

Figure 3 shows that Algorithms 1 is the only one which always gives unbiased estimators of the mean, variance and associated covariances of $Y_{.1}$. As expected, the listwise deletion method provides biased estimates inasmuch as the observed sample is not representative of the population with MNAR data. Method (ii), specifically designed for PPCA models but assuming MAR missing values, provides biased estimators. Method (iv) improves on the benchmark mean imputation and on Method (iii) as well as it explicitly takes into account the MNAR mechanism, but it still leads to biased estimates probably because of the fixed effects model assumption. Figure 2 shows that Algorithm 1 gives the best estimate of the loading matrix and the smallest imputation error. Method (i), based on the same arguments as Algorithm 1 but considering MAR data, may be considered as a second choice for this low-dimensional example as the biais is quite small (yet not in higher dimension, see Appendix C).



Figure 3: Mean (top left) and variance (top right) estimations of the missing variable and covariances (bottom) estimations of $\text{Cov}(Y_{.1}, Y_{.2})$ (*i.e.* covariance between two missing variables) and of $\text{Cov}(Y_{.1}, Y_{.8})$ (*i.e.* between one missing variable and one pivot variable). True values are indicated by red lines.

**Misspecification to the PPCA model.** The data matrix $Y \in \mathbb{R}^{n \times p}$ of size $n = 200$ and $p = 10$ is now generated under the fixed effects model such that $Y = \Theta + \epsilon$, with $\Theta \in \mathbb{R}^{n \times p}$ a low-rank matrix with $r = 2$ and $\epsilon \in \mathbb{R}^{n \times p}$ a Gaussian noise matrix with $\sigma = 0.1$. Figure 4 shows that mean and variance estimators given by Algorithm 1 have a larger variance than those given by Method (iv) precisely dedicated to this specific setting. But surprisingly, Algorithm 1 provides less biased estimates than Method (iv).
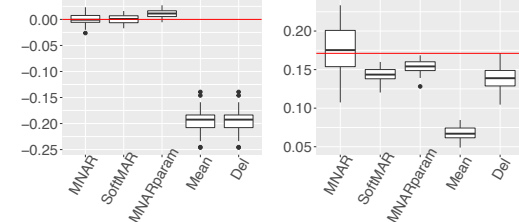


Figure 4: Mean (left) and variance (right) estimations of $Y_{.1}$ when data are generated under the fixed effects model.

In Appendix C, we report further simulation results, where we vary the features dimension ($p = 50$), the rank ($r = 5$), the missing values mechanism using probit self-masking and also multivariate MNAR (when the probability to be missing for a variable depends on its underlying values and on values of other variables that can be missing) and the percentage of missing values (10%, 50%). We obtain similar results as before, and as expected, all the methods deteriorate with an increasing percentage of missing values but our method remains stable.

In addition to the model misspecification experiment (assuming a fixed effect model), we assess the robustness of the methods in terms of noise level and we evaluate the impact of under- or overestimating the number $r$ of latent variables. When the level of noise increases, our method is very robust in terms of mean and variance estimations, and despite a bias for some covariances estimations for large noise it outperforms competitors regarding the imputation error. It also turns out that the procedure remains stable at a wrong specification of the number $r$ of latent variables.

## 4.2 Application to recommendation system data

To show the extent and feasibility of our methodology on real data, we detail the methodology on the Jester dataset [2] of 5000 users who rated 100 jokes, with 27% of missing values.

**Discussion on the assumptions.** First, considering MNAR and self-masking values is plausible because users only rate jokes they like or dislike strongly or might be ashamed to assume their taste for sexual jokes for instance. Then, Assumption **A1.**, which can be viewed as a low-rank assumption for the loading matrix, makes sense in the rating context: any variable (i.e. user preferences) can be

expressed as a linear combination of $r$ latent variables. In particular, the first latent variable opposes individuals who like jokes about physics but dislike jokes about sexuality, and conversely. Finally, Assumption **A2.** means that a user's non-response for a sexual joke given all jokes may depend on the scores of the sexual and physical jokes but not on the musical and computer jokes.

**Selecting the number $r$ of latent variables and estimating the noise variance.** In practice, to select $r$, one could use complete observations only but this is not possible when the number of features is large. As an alternative, we use a cross-validation strategy assuming M(C)AR mechanism as detailed in [7]. Algorithm 1 is robust to a misspecification of the rank (see Appendix C) and thus a reasonable heuristic may already be enough. With $r$ at hand, the noise variance is obtained directly using weighted residual sum of squares as in [9]. Without further information on the missing mechanisms, we select the $r$ pivot variables with the lowest missing rate.

**Imputation performances.** To assess the quality of our method, we introduce additional MNAR values using a logistic self-masked mechanism in a chosen variable with an initial rate of 33% and a final one of 65%. The other variables are considered M(C)AR. The process is repeated 10 times. We compare our method to the EMMAR, SoftMAR and add an imputation method based on deep generative models `Deep` [1][4]. The parametric method `MNARparam` is not performed as it does not scale on such large data. Figure 5 shows that Algorithm 1 outperforms the competitors (mean imputation corresponds to an error of 1).
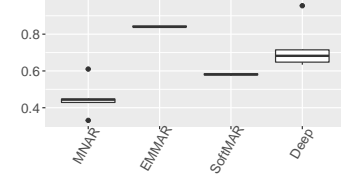
Figure 5: Imputation error for the Jester dataset.

### 4.3 Application to clinical data

We illustrate our method on the TraumaBase® dataset containing the clinical measurements of 3159 patients with brain trauma injury (see Appendix E for more information). Nine quantitative variables, selected by doctors, contain missing values (11% in the whole dataset). After discussion with doctors, some variables can be considered to have MNAR values, such as the variable *HR.ph*, which denotes the heart rate. Indeed, when the patient's condition is too critical and therefore his heart rate is either high or low, the heart rate may not be measured, as doctors prefer to provide emergency care.

As for the Jester dataset, we introduce additional MNAR values in the variable *HR.ph* (which has an initial missing rate of 1%) using a logistic self-masked mechanism leading to 50% missing values. Both the rank and the noise level are estimated using the complete-case analysis (1862 observations). The selection of the pivot variables was discussed with experts (doctors) who identified M(C)AR variables. In Figure 6, Algorithm 1 gives significantly smaller imputation error than other methods. In addition, a supervised learning task is also performed in Appendix E for which Algorithm 1 also gives the smallest prediction error.
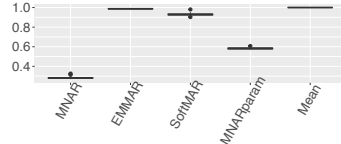
Figure 6: Imputation error for the TraumaBase dataset.

## Conclusion

In this work, we propose a new estimation and imputation method to perform PPCA with MNAR data (possibly coupled with M(C)AR data), without any need of modeling the missing mechanism. This comes with strong theoretical guarantees as identifiability and consistency, but also with an efficient algorithm. Estimating the rank in the PPCA setting with MNAR data remains non trivial. Once the number of latent variables is estimated, the noise variance can be estimated. A cross-validation strategy by additionally adding some MNAR values is a first solution, but this definitely requires further research. Another ambitious prospect would be to extend work to the exponential family to process count data, for example, which is prevalent in many application fields such as genomics.

---

[4]Note that this method requires to be trained on a complete dataset.

## Broader impact

Our goal is to provide a rigorous and consistent method for processing MNAR missing values, in data with an underlying low-rank structure. The low-rank assumption has become widespread in applications in recent years and it plays a key modeling role in many scientific and engineering tasks, such as collaborative filtering, genome-wide studies, or even functional magnetic resonance imaging.

The problem of missing data is particularly evident for large data, possibly aggregated from multiple sources, that is why we illustrate this work on a real dataset such as the medical register TraumaBase, coming from different hospitals.

Managing informative missing data is a double challenge: on the one hand, because most of the available data contains missing values, preventing the use of standard machine learning techniques; and on the other hand, because the MNAR data can introduce large bias in the statistical analysis of databases.

Because of the PPCA hypothesis and the processing of informative missing data, this work has a wide range of applications.

## References

[1] L. Gondara and K. Wang. Mida: Multiple imputation using denoising autoencoders. In PAKDD, 2018.

[2] M. Hahsler. recommenderlab: A framework for developing and testing recommendation algorithms. Technical report, 2015.

[3] José Miguel Hernández-Lobato, Neil Houlsby, and Zoubin Ghahramani. Probabilistic matrix factorization with non-random missing data. In International Conference on Machine Learning, pages 1512–1520, 2014.

[4] Joseph G Ibrahim, Stuart R Lipsitz, and M-H Chen. Missing covariates in generalized linear models when the missing data mechanism is non-ignorable. Journal of the Royal Statistical Society: Series B (Statistical Methodology), 61(1):173–190, 1999.

[5] Alexander Ilin and Tapani Raiko. Practical approaches to principal component analysis in the presence of missing values. Journal of Machine Learning Research, 11(Jul):1957–2000, 2010.

[6] Shpitser Ilya, Mohan Karthika, and Pearl Judea. Missing data as a causal and probabilistic problem in proceedings of the thirty-first conference on uncertainty in artificial intelligence, 2015.

[7] Julie Josse and François Husson. Selecting the number of components in principal component analysis using cross-validation approximations. Computational Statistics & Data Analysis, 56 (6):1869–1879, 2012.

[8] Julie Josse, Jérome Pagès, and François Husson. Testing the significance of the rv coefficient. Computational Statistics & Data Analysis, 53(1):82–91, 2008.

[9] Julie Josse, Sylvain Sardy, and Stefan Wager. denoiser: A package for low rank matrix estimation. Journal of Statistical Software, 2016.

[10] Roderick JA Little and Donald B Rubin. Statistical analysis with missing data, volume 333. John Wiley & Sons, 2014.

[11] Wei Ma and George H Chen. Missing not at random in matrix completion: The effectiveness of estimating missingness probabilities under a low nuclear norm assumption. In Advances in Neural Information Processing Systems, pages 14871–14880, 2019.

[12] Benjamin M Marlin and Richard S Zemel. Collaborative prediction and ranking with non-random missing data. In Proceedings of the third ACM conference on Recommender systems, pages 5–12, 2009.

[13] Rahul Mazumder, Trevor Hastie, and Robert Tibshirani. Spectral regularization algorithms for learning large incomplete matrices. Journal of machine learning research, 11(Aug):2287–2322, 2010.

[14] Wang Miao and Eric Tchetgen Tchetgen. Identification and inference with nonignorable missing covariate data. Statistica Sinica, 28(4):2049–2067, 2018.

[15] Wang Miao, Peng Ding, and Zhi Geng. Identifiability of normal and normal mixture models with nonignorable missing data. Journal of the American Statistical Association, 111(516): 1673–1683, 2016.

[16] Karthika Mohan and Judea Pearl. Graphical models for recovering probabilistic and causal queries from missing data. In Advances in Neural Information Processing Systems, pages 1520–1528, 2014.

[17] Karthika Mohan, Judea Pearl, and Jin Tian. Graphical models for inference with missing data. In Advances in neural information processing systems, pages 1277–1285, 2013.

[18] Karthika Mohan, Felix Thoemmes, and Judea Pearl. Estimation with incomplete data: The linear case. In IJCAI, pages 5082–5088, 2018.

[19] Kosuke Morikawa, Jae Kwang Kim, and Yutaka Kano. Semiparametric maximum likelihood estimation with data missing not at random. Canadian Journal of Statistics, 45(4):393–409, 2017.

[20] Razieh Nabi, Rohit Bhattacharya, and Ilya Shpitser. Full law identification in graphical models of missing data: Completeness results. arXiv preprint arXiv:2004.04872, 2020.

[21] Judea Pearl. Causality: models, reasoning, and inference. Econometric Theory, 19(675-685): 46, 2003.

[22] Donald B Rubin. Inference and missing data. Biometrika, 63(3):581–592, 1976.

[23] Ilya Shpitser. Consistent estimation of functions of data missing non-monotonically and not at random. In Advances in Neural Information Processing Systems, pages 3144–3152, 2016.

[24] Aude Sportisse, Claire Boyer, and Julie Josse. Imputation and low-rank estimation with missing non at random data. arXiv preprint arXiv:1812.11409, 2018.

[25] Gong Tang, Roderick JA Little, and Trivellore E Raghunathan. Analysis of multivariate missing data with nonignorable nonresponse. Biometrika, 90(4):747–764, 2003.

[26] Michael E Tipping and Christopher M Bishop. Probabilistic principal component analysis. Journal of the Royal Statistical Society: Series B (Statistical Methodology), 61(3):611–622, 1999.

[27] Sheng Wang, Jun Shao, and Jae Kwang Kim. An instrumental variable approach for identification and estimation with nonignorable nonresponse. Statistica Sinica, pages 1097–1116, 2014.

[28] Xiaojie Wang, Rui Zhang, Yu Sun, and Jianzhong Qi. Doubly robust joint learning for recommendation on data missing not at random. In International Conference on Machine Learning, pages 6638–6647, 2019.

# A  Proof of Proposition 1

**Proposition 1.** *Under Assumptions A01. and A02., the parameters $(\alpha, \Sigma)$ of the PPCA model* (1) *and the mechanism parameters $\phi = (\phi_\ell)_{\ell \in \{1,...p\}}$ are identifiable. Assuming that the noise level $\sigma^2$ is known, the parameter $B$ is identifiable up to a row permutation.*

For the sake of readability, we first present the proof of Proposition 1 in the case of the toy example presented in Section 3.1 with $p = 3$ and $r = 2$. The proof in the general setting follows.

## A.1 Proof of Proposition 1 in the case of the toy example presented in Section 3.1

Consider the setting of the toy example presented in Section 3.1 with $p = 3$ and $r = 2$. The PPCA model in (1) reads

$$\begin{cases} Y & = (Y_1 \quad Y_2 \quad Y_3) = (\alpha_1 \quad \alpha_2 \quad \alpha_3) + (W_1 \quad W_2) B + \epsilon, \\ Y & \sim \mathcal{N}(\alpha, \Sigma), \ \Sigma = B^T B + \sigma^2 I. \end{cases}$$

$Y_2$ and $Y_3$ are assumed to be observed and $Y_1$ is self-masked MNAR, *i.e.*

$$\mathbb{P}(\Omega_1 = 1 | Y_1, Y_2, Y_3; \phi_1) = \mathbb{P}(\Omega_1 = 1 | Y_1; \phi_1) = F_1(\phi_1^0 + \phi_1^1 y_1), \tag{14}$$

where $F_1$ is strictly monotone with a positive finite support and where $\phi_1 = (\phi_1^0, \phi_1^1)$.

*Proof.* Assume that $(Y, \Omega)$ and $(Y', \Omega')$ have distributions respectively parameterized by $(\alpha, \Sigma, \phi_1)$ and $(\alpha', \Sigma', \phi_1')$. Assume that $Y$ and $Y'$ have the same observed distribution, *i.e.*

$$\mathcal{L}(Y_1, \Omega_1 = 1; \alpha_1, \Sigma_{11}, \phi_1) = \mathcal{L}(Y_1', \Omega_1' = 1; \alpha_1', \Sigma_{11}', \phi_1') \tag{15}$$

$$\mathcal{L}(Y_1, Y_j, \Omega_1 = 1; \alpha_1, \alpha_j, \Sigma_{(1j)}, \phi_1) = \mathcal{L}(Y_1', Y_j', \Omega_1' = 1; \alpha_1', \alpha_j', \Sigma_{(1j)}', \phi_1') \qquad j \in \{2, 3\}, \tag{16}$$

where $\Sigma_{(1j)}$ is the covariance matrix $\begin{pmatrix} \Sigma_{11} & \Sigma_{1j} \\ \Sigma_{1j} & \Sigma_{jj} \end{pmatrix}$. In order to show that parameters identifiability holds, we need to show that (15) and (16) imply that $\alpha = \alpha'$, $\Sigma = \Sigma'$ and $\phi_1 = \phi_1'$. Then, under a known noise level $\sigma^2$, we prove that $B$ and $B'$ are equal up to a row permutation.

As $(Y_2, Y_3)$ and $(Y_2', Y_3')$ are fully observed, the parameters of the distributions $\mathcal{L}(Y_2), \mathcal{L}(Y_2'), \mathcal{L}(Y_3)$, $\mathcal{L}(Y_3'), \mathcal{L}(Y_2, Y_3)$ and $\mathcal{L}(Y_2', Y_3')$ are identifiable. It trivially implies that $\alpha_2 = \alpha_2', \Sigma_{22} = \Sigma_{22}'$, $\alpha_3 = \alpha_3', \Sigma_{33} = \Sigma_{33}'$ and $\Sigma_{23} = \Sigma_{23}'$.

**Identifiability of the MNAR variable variance.** Equation (15) can be rewritten in terms of density function as follows

$$f_{Y_1, \Omega_1 = 1}(y_1; \alpha_1, \Sigma_{11}, \phi_1) = f_{Y_1', \Omega_1' = 1}(y_1; \alpha_1', \Sigma_{11}', \phi_1') \qquad \forall y_1 \in \mathbb{R}.$$

Given the missing mechanism in (14) and that $Y_{.1} \sim \mathcal{N}(\alpha_1, \Sigma_{11})$, [15, Theorem 1 a)] ensures that $\Sigma_{11} = \Sigma_{11}'$.

**Identifiability of the Mean and the MNAR mechanism parameter.** Using (15) and (16), the previous computations entail that

$$\mathcal{L}(Y_2 | Y_1, \Omega_1 = 1; \alpha_1, \alpha_2, \Sigma_{(12)}, \phi_1) = \mathcal{L}(Y_2' | Y_1', \Omega_1' = 1; \alpha_1', \alpha_2', \Sigma_{(12)}', \phi_1'),$$

noting that

$$f_{Y_2 | Y_1 = y_1, \Omega_1 = 1}(y_2; \alpha_1, \alpha_2, \Sigma_{(12)}, \phi_1) = \frac{f_{Y_1, Y_2, \Omega_1 = 1}(y_1, y_2; \alpha_1, \alpha_2, \Sigma_{(12)}, \phi_1)}{f_{Y_1, \Omega_1 = 1}(y_1; \alpha_1, \Sigma_{11}, \phi_1)} \qquad \forall (y_1, y_2) \in \mathbb{R}^2$$

One obtains

$$\frac{\mathbb{P}(\Omega_1 = 1 | Y_1 = y_1, Y_2 = y_2; \phi_1) f_{Y_2 | Y_1 = y_1}(y_2; \alpha_1, \alpha_2, \Sigma_{(12)})}{\mathbb{P}(\Omega_1 = 1 | Y_1 = y_1; \phi_1)}$$

$$= \frac{\mathbb{P}(\Omega_1' = 1 | Y_1' = y_1, Y_2' = y_2; \phi_1') f_{Y_2' | Y_1' = y_1}(y_2; \alpha_1', \alpha_2', \Sigma_{(12)}')}{\mathbb{P}(\Omega_1' = 1 | Y_1 = y_1; \phi_1')} \qquad \forall (y_1, y_2) \in \mathbb{R}^2$$

Yet,

$$\begin{aligned} \mathbb{P}(\Omega_1 = 1 | Y_1 = y_1, Y_2 = y_2; \phi_1) &= \mathbb{E}[\mathbb{E}[1_{\Omega_1 = 1} | Y_1 = y_1, Y_2 = y_2, Y_3 = y_3; \phi_1] | Y_1 = y_1, Y_2 = y_2] \\ &= \mathbb{E}[\mathbb{P}(\Omega_1 = 1 | Y = y; \phi_1) | Y_1 = y_1, Y_2 = y_2] \\ &= \mathbb{E}[\mathbb{P}(\Omega_1 = 1 | Y_1 = y_1; \phi_1) | Y_1 = y_1, Y_2 = y_2] \\ &= \mathbb{P}(\Omega_1 = 1 | Y = y_1; \phi_1) \end{aligned} \tag{17}$$

by measurability. It implies for all $y_1 \in \mathbb{R}$ and $y_2 \in \mathbb{R}$

$$f_{Y_2|Y_1=y_1}(y_2; \alpha_1, \alpha_2, \Sigma_{(12)}) = f_{Y_2'|Y_1'=y_1}(y_2; \alpha_1', \alpha_2', \Sigma_{(12)}')$$

which leads to the equality of the conditional expectations and variances associated to the above densities:

$$\alpha_2 + \Sigma_{12}\Sigma_{11}^{-1}(\alpha_1 - y_1) = \alpha_2 + \Sigma_{12}'\Sigma_{11}^{-1}(\alpha_1' - y_1) \qquad \forall y_1 \in \mathbb{R}$$
$$\Sigma_{22} - \Sigma_{12}^2\Sigma_{11}^{-1} = \Sigma_{22} - (\Sigma_{12}')^2\Sigma_{11}^{-1}.$$

It implies that

$$\Sigma_{12}^2 = (\Sigma_{12}')^2 \implies |\Sigma_{12}| = |\Sigma_{12}'| \tag{18}$$
$$\frac{\Sigma_{21}}{\Sigma_{21}'} = \frac{(\alpha_1' - y_1)}{(\alpha_1 - y_1)} \implies |\alpha_1 - y_1| = |\alpha_1' - y_1| \qquad \forall y_1 \in \mathbb{R} \tag{19}$$

Equation (19) implies that $\alpha_1 = \alpha_1'$, since for $y_1 = \alpha_1'$, one has $\alpha_1 - \alpha_1' = 0$.

Using (16), one has

$$\mathbb{P}(\Omega_1 = 1 | Y_1 = y_1, Y_2 = y_2; \phi_1)f_{(Y_1,Y_2)}(y_1, y_2; \alpha_1, \alpha_2, \Sigma_{(12)})$$
$$= \mathbb{P}(\Omega_1' = 1 | Y_1' = y_1, Y_2' = y_2; \phi_1')f_{(Y_1',Y_2')}(y_1, y_2; \alpha_1', \alpha_2', \Sigma_{(12)}') \qquad \forall (y_1, y_2) \in \mathbb{R}^2 \tag{20}$$

Using (17),

$$\frac{\exp\left(-\frac{1}{2}\begin{pmatrix} y_1 - \alpha_1 & y_2 - \alpha_2 \end{pmatrix}\Sigma_{(12)}^{-1}\begin{pmatrix} y_1 - \alpha_1 \\ y_2 - \alpha_2 \end{pmatrix}\right)}{\exp\left(-\frac{1}{2}\begin{pmatrix} y_1 - \alpha_1 & y_2 - \alpha_2 \end{pmatrix}(\Sigma_{(12)}')^{-1}\begin{pmatrix} y_1 - \alpha_1 \\ y_2 - \alpha_2 \end{pmatrix}\right)}\frac{\mathbb{P}(\Omega_1 = 1 | Y_1 = y_1; \phi_1)}{\mathbb{P}(\Omega_1' = 1 | Y_1' = y_1; \phi_1')} = \frac{\sqrt{\det(\Sigma_{(12)})}}{\sqrt{\det(\Sigma_{(12)}')}},$$

where $\det(\Sigma_{(12)})$ denotes the determinant of the matrix $\Sigma_{(12)}$.

With (18), one has $\Sigma_{11}\Sigma_{22} - \Sigma_{12}^2 = \Sigma_{11}\Sigma_{22} - (\Sigma_{12}')^2$ and $\frac{\sqrt{\det(\Sigma_{(12)})}}{\sqrt{\det(\Sigma_{(12)}')}} = 1$.

It leads to $\forall (y_1, y_2) \in \mathbb{R}^2$,

$$K \cdot \frac{\mathbb{P}(\Omega_1 = 1 | Y_1 = y_1; \phi_1)}{\mathbb{P}(\Omega_1' = 1 | Y_1' = y_1; \phi_1')} = 1,$$

with

$$K := \frac{\exp\left(-\frac{1}{2\det(\Sigma_{(12)})}\left((y_1 - \alpha_1)^2\Sigma_{11} + (y_2 - \alpha_2)^2\Sigma_{22} - 2(y_1 - \alpha_1)(y_2 - \alpha_2)\Sigma_{12}\right)\right)}{\exp\left(-\frac{1}{2\det(\Sigma_{(12)})}\left((y_1 - \alpha_1)^2\Sigma_{11} + (y_2 - \alpha_2)^2\Sigma_{22} - 2(y_1 - \alpha_1')(y_2 - \alpha_2)\Sigma_{12}'\right)\right)}.$$

The quantity $K$ is equal to one, because

$$(y_2 - \alpha_2)\left((y_1 - \alpha_1)\Sigma_{12} - (y_1 - \alpha_1')\Sigma_{12}'\right) = 0$$

using (19). Thus,

$$\frac{\mathbb{P}(\Omega_1 = 1 | Y_1 = y_1; \phi_1)}{\mathbb{P}(\Omega_1' = 1 | Y_1' = y_1; \phi_1')} = 1 \iff F_1(\phi_1^0 + \phi_1^1 y_1) = F_1((\phi')_1^0 + (\phi')_1^1 y_1) \qquad \forall y_1 \in \mathbb{R}$$

As $F_1$ is strictly monotone, it is an injective function. Thus,

$$\phi_1^0 + \phi_1^1 y_1 = (\phi')_1^0 + (\phi')_1^1 y_1 \qquad \forall y_1 \in \mathbb{R} \iff (\phi_1^0 - (\phi')_1^0) + ((\phi')_1^1 - \phi_1^1)y_1 = 0 \qquad \forall y_1 \in \mathbb{R}$$

It implies $\phi_1 = \phi_1'$.

**Identifiability of the Covariances of the MNAR variable.** Equation (20) thus leads to

$$f_{(Y_1,Y_2)}(y_1,y_2;\alpha_1,\alpha_2,\Sigma_{(12)}) = f_{(Y_1',Y_2')}(y_1,y_2;\alpha_1',\alpha_2',\Sigma_{(12)}') \qquad \forall(y_1,y_2) \in \mathbb{R}^2$$

One can conclude that $\Sigma_{12} = \Sigma_{12}'$. The same reasoning may be done for the covariance between $Y_1$ and $Y_3$.

**Identifiability of the loading matrix.** One wants to prove $B = B'$ up to row permutation. One has

$$\Sigma = \Sigma' \Leftrightarrow \Sigma - \sigma^2 I_{p\times p} = \Sigma' - \sigma^2 I_{p\times p}$$
$$\Leftrightarrow B^T B = (B')^T B' \qquad (21)$$

As $B^T B$ is a positive symetric matrix of rank 2, one has the following singular value decomposition,

$$B^T B = (B')^T B' = U D U^T,$$

where $U = (u_1|u_2|u_3) \in \mathbb{R}^{3\times 3}$ the orthogonal matrix of singular vector and

$$D = \begin{pmatrix} \sqrt{d_1} & 0 & 0 \\ 0 & \sqrt{d_2} & 0 \\ 0 & 0 & 0 \end{pmatrix} \in \mathbb{R}^{3\times 3}$$

with $d_1 \geqslant d_2 \geqslant 0$. One can choose

$$B = \begin{pmatrix} \sqrt{d_1} u_1^T \\ \sqrt{d_2} u_2^T \end{pmatrix}$$

noting that a row permutation of B would not change the product $B^T B$. Therefore, $B = B'$ up to a row permutation.

$\square$

## A.2 Proof of Proposition 1 in the general case

We present the proof of Proposition 1 in the general case where $d$ variables are self-masked MNAR and $p - d$ variables are MCAR.

*Proof.* Assume that $(Y, \Omega)$ and $(Y', \Omega')$ have distributions respectively parameterized by $(\alpha, \Sigma, \phi)$ and $(\alpha', \Sigma', \phi')$. Assume that $Y$ and $Y'$ have the same following observed distributions

$$\mathcal{L}(Y_j, \Omega_j = 1; \alpha_j, \Sigma_{jj}, \phi_j) = \mathcal{L}(Y_j', \Omega_j' = 1; \alpha_j', \Sigma_{jj}', \phi_j') \qquad \forall j \in \{1, \ldots, p\}, \qquad (22)$$

$$\mathcal{L}(Y_j, Y_k, \Omega_j = 1, \Omega_k = 1; \alpha_j, \alpha_k, \Sigma_{(jk)}, \phi_j, \phi_k)$$
$$= \mathcal{L}(Y_j', Y_k', \Omega_j' = 1, \Omega_k' = 1; \alpha_j', \alpha_k', \Sigma_{(jk)}', \phi_j', \phi_k') \qquad \forall j \neq k \in \{1, \ldots, p\}, \quad (23)$$

where $\Sigma_{(jk)}$ denotes the covariance matrix $\begin{pmatrix} \Sigma_{jj} & \Sigma_{jk} \\ \Sigma_{jk} & \Sigma_{kk} \end{pmatrix}$.

In order to show that parameters identifiability holds, we need to show that (22) and (23) implies that $\alpha = \alpha'$, $\Sigma = \Sigma'$ and $\phi = \phi'$. Then, under a known noise level $\sigma^2$, we will prove that $B$ and $B'$ are equal up to row permutations.

In what follows, $f_{Y_{.j}}$ or $f_{(Y_{.j}, Y_{.k})}$ respectively denote the density function of $Y_{.j}$, and of $(Y_{.j}, Y_{.k})$.

In the following, we will use the following tip, for any $l \in \{1, \ldots, p\}$ and $\mathcal{K} \subset \{1, \ldots, p\}\backslash\{l\}$ such that $0 \leqslant |\mathcal{K}| \leqslant p-1$,

$$\mathbb{P}(\Omega_l = 1|Y_l = y_l, Y_{\mathcal{K}} = y_{\mathcal{K}}; \phi_l) = \mathbb{E}[\mathbb{E}[\mathbf{1}_{\Omega_l=1}|Y; \phi_l]|Y_l = y_l, Y_{\mathcal{K}} = y_{\mathcal{K}}]$$
$$= \mathbb{E}[\mathbb{P}(\Omega_l = 1|Y = y; \phi_l)|Y_l = y_l, Y_{\mathcal{K}} = y_{\mathcal{K}}]$$

14

Thus, using the mechanisms in **A01.**,

$$\mathbb{P}(\Omega_l = 1 | Y_l = y_l, Y_{\mathcal{K}} = y_{\mathcal{K}}; \phi_l)$$
$$= \begin{cases} \mathbb{E}[\mathbb{P}(\Omega_l = 1 | Y_l = y_l; \phi_l) | Y_l = y_l, Y_{\mathcal{K}} = y_{\mathcal{K}}] & \text{if } Y_l \text{ is self-masked MNAR} \\ \mathbb{E}[\mathbb{P}(\Omega_l = 1; \phi_l) | Y_l = y_l, Y_{\mathcal{K}} = y_{\mathcal{K}}] & \text{if } Y_l \text{ is MCAR} \end{cases}$$

Thus,

$$\mathbb{P}(\Omega_l = 1 | Y_l = y_l, Y_{\mathcal{K}} = y_{\mathcal{K}}; \phi_l) = \begin{cases} \mathbb{P}(\Omega_l = 1 | Y_l = y_l; \phi_l) & \text{if } Y_l \text{ is self-masked MNAR} \quad (24) \\ \mathbb{P}(\Omega_l = 1; \phi_l) & \text{if } Y_l \text{ is MCAR} \quad (25) \end{cases}$$

by measurability if $Y_l$ is self-masked MNAR and by independence if $Y_l$ is MCAR.

**Identifiability of the parameters for the not-MNAR variables $(Y_j)_{j \in \overline{\mathcal{M}}}$.**

**Mechanism parameter, Mean and Variance of $Y_j$, $j \in \overline{\mathcal{M}}$.** Equation (22) leads to

$$\mathbb{P}(\Omega_j = 1 | Y_j = y_j; \phi_j) f_{Y_j}(y_j; \alpha_j, \Sigma_{jj}) = \mathbb{P}(\Omega_j' = 1 | Y_j' = y_j; \phi_j') f_{Y_j'}(y_j; \alpha_j', \Sigma_{jj}') \qquad \forall y_j \in \mathbb{R}.$$

Using (25), $P(\Omega_j = 1) = \mathbb{P}(\Omega_j = 1 | Y_j = y_j; \phi_j) = F_j(\phi_j)$. This distribution is identifiable since it pertains to a conditional distribution of the observed data. As $F_j$ is strictly monotone, it implies that

$$F_j(\phi_j) = F_j(\phi_j') \iff \phi_j = \phi_j'.$$

As $\phi_j = \phi_j'$, one obtains

$$f_{Y_j}(y_j; \alpha_j, \Sigma_{jj}) = f_{Y_j'}(y_j; \alpha_j', \Sigma_{jj}') \qquad \forall y_j \in \mathbb{R}$$

which directly implies that $\alpha_j = \alpha_j'$ and $\Sigma_{jj} = \Sigma_{jj}'$, since $Y_j$ and $Y_j'$ are Gaussian variables.

**Covariance between two not MNAR variables $Y_j$ and $Y_k$, $j \neq k \in \overline{\mathcal{M}}$.** Equation (23) gives that for all $(y_j, y_k) \in \mathbb{R}^2$

$$\mathbb{P}(\Omega_j = 1, \Omega_k = 1 | Y_j = y_j, Y_k = y_k; \phi_j, \phi_k) f_{(Y_j, Y_k)}(y_j, y_k; \alpha_j, \alpha_k, \Sigma_{(j,k)})$$
$$= \mathbb{P}(\Omega_j' = 1, \Omega_k' = 1 | Y_j' = y_j, Y_k' = y_k; \phi_j', \phi_k') f_{(Y_j', Y_k')}(y_j, y_k; \alpha_j', \alpha_k', \Sigma_{(j,k)}'), \quad (26)$$

and one has as well that

$$\mathbb{P}(\Omega_j = 1, \Omega_k = 1 | Y_j = y_j, Y_k = y_k; \phi_j, \phi_k) = \mathbb{P}(\Omega_j = 1 | Y_j = y_j; \phi_j) \mathbb{P}(\Omega_k = 1 | Y_k = y_k; \phi_k),$$

using **A02.**. Likewise,

$$\mathbb{P}(\Omega_j' = 1, \Omega_k' = 1 | Y_j' = y_j, Y_k' = y_k; \phi_j', \phi_k') = \mathbb{P}(\Omega_j' = 1 | Y_j' = y_j; \phi_j') \mathbb{P}(\Omega_k' = 1 | Y_k' = y_k; \phi_k').$$

Given that $\phi_j = \phi_j'$ and $\phi_k = \phi_k'$, one obtains

$$\mathbb{P}(\Omega_j = 1, \Omega_k = 1 | Y_j = y_j, Y_k = y_k; \phi_j, \phi_k) = \mathbb{P}(\Omega_j' = 1, \Omega_k' = 1 | Y_j' = y_j, Y_k' = y_k; \phi_j, \phi_k).$$

Thus, Equation (26) leads to, for all $(y_j, y_k) \in \mathbb{R}^2$,

$$f_{(Y_j, Y_k)}(y_j, y_k; \alpha_j, \alpha_k, \Sigma_{(j,k)}) = f_{(Y_j', Y_k')}(y_j, y_k; \alpha_j', \alpha_k', \Sigma_{(j,k)}'),$$

and $\Sigma_{jk} = \Sigma_{jk}'$.

**Identifiability of the parameters for the MNAR variables.**

**Variance of $Y_m$, $m \in \mathcal{M}$.** Equation (22) gives that

$$f_{(Y_m, \Omega_m = 1)}(y_m; \alpha_m, \Sigma_{mm}, \phi_m) = f_{(Y_m', \Omega_m' = 1)}(y_m; \alpha_m', \Sigma_{mm}', \phi_m') \qquad \forall y_m \in \mathbb{R}.$$

Given the self-masked missing mechanism in **A01.** and that $Y_{.m} \sim \mathcal{N}(\alpha_m, \Sigma_{mm})$, [15, Theorem 1 a)] ensures that $\Sigma_{mm} = \Sigma_{mm}'$.

15

**Mean and mechanism parameter of $Y_m, m \in \mathcal{M}$.** Let $j \in \overline{\mathcal{M}}$ (a not MNAR variable). One has

$$\mathcal{L}(Y_j, \Omega_j = 1 | Y_m, \Omega_m = 1; \alpha_j, \alpha_m, \Sigma_{(jm)}, \phi_j, \phi_m)$$
$$= \mathcal{L}(Y_j', \Omega_j' = 1 | Y_m', \Omega_m' = 1; \alpha_j', \alpha_m', \Sigma_{(jm)}', \phi_j', \phi_m') \quad (27)$$

using (22) and (23) and noting that

$$f_{(Y_j, \Omega_j=1)|Y_m=y_m, \Omega_m=1}(y_j; \alpha_j, \alpha_m, \Sigma_{(jm)}, \phi_j, \phi_m)$$
$$= \frac{f_{(Y_j, \Omega_j=1, Y_m, \Omega_m=1)}(y_j, y_m; \alpha_j, \alpha_m, \Sigma_{(jm)}, \phi_j, \phi_m)}{f_{(Y_m, \Omega_m=1)}(y_m; \alpha_m, \Sigma_{mm}, \phi_m)} \qquad \forall (y_j, y_m) \in \mathbb{R}^2.$$

Equation (27) implies that $\forall (y_j, y_m) \in \mathbb{R}^2$,

$$\mathbb{P}(\Omega_j = 1 | Y_j = y_j, Y_m = y_m, \Omega_m = 1; \phi_j) \frac{\mathbb{P}(\Omega_m = 1 | Y_j = y_j, Y_m = y_m; \phi_m) f_{Y_j | Y_m = y_m}(y_j; \alpha_j, \alpha_m, \Sigma_{(jm)})}{\mathbb{P}(\Omega_m = 1 | Y_m = y_m; \phi_m)}$$
$$= \mathbb{P}(\Omega_j' = 1 | Y_j' = y_j, Y_m' = y_m, \Omega_m' = 1; \phi_j') \frac{\mathbb{P}(\Omega_m' = 1 | Y_j' = y_j, Y_m' = y_m; \phi_m') f_{Y_j' | Y_m' = y_m}(y_j; \alpha_j', \alpha_m', \Sigma_{(jm)}')}{\mathbb{P}(\Omega_m' = 1 | Y_m' = y_m; \phi_m')}$$
$$(28)$$

One can note that

$$\mathbb{P}(\Omega_j = 1 | Y_j = y_j, Y_m = y_m, \Omega_m = 1; \phi_j) = \mathbb{P}(\Omega_j = 1 | Y_j = y_j; \phi_j).$$

Indeed,

$$\mathbb{P}(\Omega_j = 1 | Y_j = y_j, Y_m = y_m, \Omega_m = 1; \phi_j) = \frac{\mathbb{P}(\Omega_j = 1 \cap \Omega_m = 1 | Y_j = y_j, Y_m = y_m; \phi_j, \phi_m)}{\mathbb{P}(\Omega_m = 1 | Y_j = y_j, Y_m = y_m; \phi_m)}$$
$$= \frac{\mathbb{P}(\Omega_j = 1 | Y_j = y_j; \phi_j) \mathbb{P}(\Omega_m = 1 | Y_m = y_m; \phi_m)}{\mathbb{P}(\Omega_m = 1 | Y_j = y_j, Y_m = y_m; \phi_m)}$$
$$= \mathbb{P}(\Omega_j = 1 | Y_j = y_j; \phi_j),$$

using **A02.** in the second step. Likewise,

$$\mathbb{P}(\Omega_j' = 1 | Y_j' = y_j, Y_m' = y_m, \Omega_m' = 1; \phi_j') = \mathbb{P}(\Omega_j' = 1 | Y_j' = y_j; \phi_j').$$

Given that $\phi_j = \phi_j'$,

$$\mathbb{P}(\Omega_j = 1 | Y_j = y_j, Y_m = y_m, \Omega_m = 1; \phi_j) = \mathbb{P}(\Omega_j' = 1 | Y_j' = y_j, Y_m' = y_m, \Omega_m' = 1; \phi_j')$$

Thus, Equation (28) leads to

$$\frac{\mathbb{P}(\Omega_m = 1 | Y_j = y_j, Y_m = y_m; \phi_m) f_{Y_j | Y_m = y_m}(y_j; \alpha_j, \alpha_m, \Sigma_{(jm)})}{\mathbb{P}(\Omega_m = 1 | Y_m = y_m; \phi_m)}$$
$$= \frac{\mathbb{P}(\Omega_m' = 1 | Y_j' = y_j, Y_m' = y_m; \phi_m') f_{Y_j' | Y_m' = y_m}(y_j; \alpha_j', \alpha_m', \Sigma_{(jm)}')}{\mathbb{P}(\Omega_m' = 1 | Y_m' = y_m; \phi_m')} \qquad \forall (y_j, y_m) \in \mathbb{R}^2.$$

As $\mathbb{P}(\Omega_m = 1 | Y_j = y_j, Y_m = y_m; \phi_m) = \mathbb{P}(\Omega_m = 1 | Y_m = y_m; \phi_m)$ by using (A.2), one obtains

$$f_{Y_j | Y_m = y_m}(y_j; \alpha_j, \alpha_m, \Sigma_{(jm)}) = f_{Y_j' | Y_m' = y_m}(y_j; \alpha_j', \alpha_m', \Sigma_{(jm)}') \qquad \forall (y_j, y_m) \in \mathbb{R}^2,$$

which leads to the equality of the conditional expectation and variance, as follows:

$$\alpha_j + \Sigma_{mj} \Sigma_{mm}^{-1}(\alpha_m - y_m) = \alpha_j' + \Sigma_{mj}'(\Sigma_{mm}')^{-1}(\alpha_m' - y_m) \qquad \forall (y_j, y_m) \in \mathbb{R}^2$$
$$\Sigma_{jj} - \Sigma_{mj}^2 \Sigma_{mm}^{-1} = \Sigma_{jj}' - (\Sigma_{mj}')^2 (\Sigma_{mm}')^{-1}$$

As $\alpha_j = \alpha_j'$ and $\Sigma_{mm} = \Sigma_{mm}'$,

$$\Sigma_{mj}^2 = (\Sigma_{mj}')^2 \implies |\Sigma_{mj}| = |\Sigma_{mj}'| \tag{29}$$
$$\frac{\Sigma_{mj}}{\Sigma_{mj}'} = \frac{(\alpha_m' - y_m)}{(\alpha_m - y_m)} \implies |\alpha_m - y_m| = |\alpha_m' - y_m| \qquad \forall y_m \in \mathbb{R} \tag{30}$$

Equation (30) implies that $\alpha_m = \alpha'_m$, since for $y_m = \alpha'_m$, one has $\alpha_m - \alpha'_m = 0$.

In addition, using (23), one has for all $(y_j, y_m) \in \mathbb{R}^2$,

$$\mathbb{P}(\Omega_j = 1, \Omega_m = 1 | Y_j = y_j, Y_m = y_m; \phi_j, \phi_m) f_{(Y_j, Y_m)}(y_j, y_m; \alpha_j, \alpha_m, \Sigma_{(jm)})$$
$$= \mathbb{P}(\Omega'_j = 1, \Omega'_m = 1 | Y'_j = y_j, Y'_m = y_m; \phi'_j, \phi'_m) f_{(Y'_j, Y'_m)}(y_j, y_m; \alpha'_j, \alpha'_m, \Sigma'_{(jm)}) \quad (31)$$

One can note that

$$\mathbb{P}(\Omega_j = 1, \Omega_m = 1 | Y_j = y_j, Y_m = y_m; \phi_j, \phi_m)$$
$$= \mathbb{P}(\Omega_j = 1; \phi_j) \mathbb{P}(\Omega_m = 1 | Y_m = y_m; \phi_m),$$

using **A02.** and the tips given in and (25). The same equation holds for $(Y'_j, Y'_m, \Omega'_j, \Omega'_m)$ with the parameters $(\phi'_j, \phi'_m)$. Using $\phi_j = \phi'_j$, Equation (31) leads to

$$\mathbb{P}(\Omega_m = 1 | Y_m = y_m; \phi_m) f_{(Y_j, Y_m)}(y_j, y_m; \alpha_j, \alpha_m, \Sigma_{(jm)}) =$$
$$\mathbb{P}(\Omega'_m = 1 | Y'_m = y_m; \phi'_m) f_{(Y'_j, Y'_m)}(y_j, y_m; \alpha'_j, \alpha'_m, \Sigma'_{(jm)}) \qquad \forall (y_j, y_m) \in \mathbb{R}^2. \quad (32)$$

It implies that, $\forall (y_j, y_m) \in \mathbb{R}^2$,

$$\frac{\exp\left(-\frac{1}{2}\begin{pmatrix} y_j - \alpha_j & y_m - \alpha_m \end{pmatrix} \Sigma_{(jm)}^{-1} \begin{pmatrix} y_j - \alpha_j \\ y_m - \alpha_m \end{pmatrix}\right)}{\exp\left(-\frac{1}{2}\begin{pmatrix} y_j - \alpha'_j & y_m - \alpha'_m \end{pmatrix} (\Sigma'_{(jm)})^{-1} \begin{pmatrix} y_j - \alpha'_j \\ y_m - \alpha'_m \end{pmatrix}\right)} \frac{\mathbb{P}(\Omega_m = 1 | Y_m = y_m; \phi_m)}{\mathbb{P}(\Omega'_m = 1 | Y'_m = y_m; \phi'_m)} = \frac{\sqrt{\det(\Sigma_{(jm)})}}{\sqrt{\det(\Sigma'_{(jm)})}},$$

where $\det(\Sigma_{(jm)})$ denotes the determinant of the covariance matrix $\Sigma_{(jm)}$.

With $\Sigma_{jj} = \Sigma'_{jj}$, $\Sigma_{mm} = \Sigma'_{mm}$ and Equation (29), one has

$$\Sigma_{jj}\Sigma_{mm} - \Sigma_{mj}^2 = \Sigma_{jj}\Sigma_{mm} - (\Sigma'_{mj})^2 \qquad \implies \qquad \frac{\sqrt{\det(\Sigma_{(jm)})}}{\sqrt{\det(\Sigma'_{(jm)})}} = 1.$$

Besides, using $\alpha_j = \alpha'_j$, $\Sigma_{jj} = \Sigma'_{jj}$ and $\Sigma_{mm} = \Sigma'_{mm}$, one obtains that for all $(y_j, y_m) \in \mathbb{R}^2$,

$$K \cdot \frac{\mathbb{P}(\Omega_m = 1 | Y_m = y_m; \phi_m)}{\mathbb{P}(\Omega'_m = 1 | Y'_m = y_m; \phi'_m)} = 1,$$

with

$$K := \frac{\exp\left(-\frac{1}{2\det(\Sigma_{(jm)})}\left((y_j - \alpha_j)^2\Sigma_{jj} + (y_m - \alpha_m)^2\Sigma_{mm} - 2(y_j - \alpha_j)(y_m - \alpha_m)\Sigma_{mj}\right)\right)}{\exp\left(-\frac{1}{2\det(\Sigma_{(jm)})}\left((y_j - \alpha_j)^2\Sigma_{jj} + (y_m - \alpha_m)^2\Sigma_{mm} - 2(y_j - \alpha_j)(y_m - \alpha'_m)\Sigma'_{mj}\right)\right)}.$$

The quantity $K$ is equal to one, because

$$(y_j - \alpha_j)((y_m - \alpha_m)\Sigma_{mj} - (y_m - \alpha'_m)\Sigma'_{mj}) = 0$$

using (30). Thus, for all $y_m \in \mathbb{R}$,

$$\frac{\mathbb{P}(\Omega_m = 1 | Y_m = y_m; \phi_m)}{\mathbb{P}(\Omega'_m = 1 | Y'_m = y_m; \phi'_m)} = 1 \quad \iff \quad F_m(\phi_m^0 + \phi_m^1 y_m) = F_m((\phi')_m^0 + (\phi')_m^1 y_m).$$

As F is strictly monotone, it is an injective function. Thus,

$$\phi_m^0 + \phi_m^1 y_m = (\phi')_m^0 + (\phi')_m^1 y_m \Leftrightarrow ((\phi')_m^0 - \phi_m^0) + ((\phi')_m^1 - \phi_m^1)y_m = 0 \qquad \forall y_1 \in \mathbb{R}$$

It implies that $\phi_m = \phi'_m$.

**Covariance between $Y_j$ and $Y_m$ with $j \in \overline{\mathcal{M}}, m \in \mathcal{M}$.** Using (32) and $\phi_m = \phi'_m$, one has

$$f_{(Y_j, Y_m)}(y_j, y_m; \alpha_j, \alpha_m, \Sigma_{(jm)}) = f_{(Y'_j, Y'_m)}(y_j, y_m; \alpha'_j, \alpha'_m, \Sigma'_{(jm)}) \qquad \forall (y_j, y_m) \in \mathbb{R}^2$$

One can conclude that $\Sigma_{mj} = \Sigma'_{mj}$.

**Covariance between $Y_\ell$ and $Y_m$ with $\ell \neq m \in \mathcal{M}$.**   Using (23), one has for all $(y_\ell, y_m) \in \mathbb{R}^2$,

$$\mathbb{P}(\Omega_\ell = 1, \Omega_m = 1 | Y_j = y_j, Y_m = y_m; \phi_\ell, \phi_m) f_{(Y_\ell, Y_m)}(y_\ell, y_m; \alpha_\ell, \alpha_m, \Sigma_{(\ell m)})$$
$$= \mathbb{P}(\Omega'_\ell = 1, \Omega'_m = 1 | Y'_\ell = y_\ell, Y'_m = y_m; \phi'_\ell, \phi'_m) f_{(Y'_\ell, Y'_m)}(y_\ell, y_m; \alpha'_\ell, \alpha'_m, \Sigma'_{(\ell m)}) \quad (33)$$

One can note that

$$\mathbb{P}(\Omega_\ell = 1, \Omega_m = 1 | Y_\ell = y_\ell, Y_m = y_m; \phi_\ell, \phi_m)$$
$$= \mathbb{P}(\Omega_\ell = 1 | Y_\ell = y_\ell; \phi_\ell) \mathbb{P}(\Omega_m = 1 | Y_m = y_m; \phi_m),$$

using **A02.** and the tip given in (A.2). The same equation holds for $(Y'_\ell, Y'_m, \Omega'_\ell, \Omega'_m)$ with the parameters $(\phi'_\ell, \phi'_m)$. Yet $\phi_\ell = \phi'_\ell$ and $\phi_m = \phi'_m$, which gives, for all $(y_j, y_m) \in \mathbb{R}^2$,

$$\mathbb{P}(\Omega_\ell = 1, \Omega_m = 1 | Y_\ell = y_\ell, Y_m = y_m; \phi_\ell, \phi_m) = \mathbb{P}(\Omega'_\ell = 1, \Omega'_m = 1 | Y'_\ell = y_\ell, Y'_m = y_m; \phi_\ell, \phi'_m).$$

Equation (33) leads to

$$f_{(Y_\ell, Y_m)}(y_\ell, y_m; \alpha_\ell, \alpha_m, \Sigma_{(\ell m)}) = f_{(Y'_\ell, Y'_m)}(y_\ell, y_m; \alpha'_\ell, \alpha'_m, \Sigma'_{(\ell m)}) \qquad \forall (y_\ell, y_m) \in \mathbb{R}^2,$$

which implies that $\Sigma_{\ell m} = \Sigma'_{\ell m}$.

**Identifiability of the loading matrix.**   One wants to prove that $B = B'$ up to a row permutation. One has

$$\Sigma = \Sigma' \iff \Sigma - \sigma^2 I_{p \times p} = \Sigma' - \sigma^2 I_{p \times p}$$
$$\iff B^T B = (B')^T B' \quad (34)$$

As $B^T B$ is a positive symetric matrix of rank $r$, its singular value decomposition reads

$$B^T B = (B')^T B' = U D U^T,$$

where $U = (u_1 | \dots | u_p) \in \mathbb{R}^{p \times p}$ is an orthogonal matrix containing the singular vectors and

$$D = \begin{pmatrix} \sqrt{d_1} & & & & & \\ & \ddots & & & 0 & \\ & & \sqrt{d_r} & & & \\ & 0 & & 0 & & \\ & & & & \ddots & \\ & & & & & 0 \end{pmatrix} \in \mathbb{R}^{p \times p}$$

with $d_1 \geqslant \dots \geqslant d_r \geqslant 0$. One can choose

$$B = \left( \begin{array}{c} \sqrt{d_1} u_1^T \\ \hline \vdots \\ \hline \sqrt{d_r} u_r^T \end{array} \right)$$

A row permutation of B does not change the product $B^T B$. Therefore, $B = B'$ up to a row permutation.

$\square$

# B   Proof for Section 3

## B.1   Proof of Lemma 2

**Lemma 2.** *Under the PPCA model* (1) *and Assumption **A1.**, choose $j \in \mathcal{J}$. Denote $B^{-1} \in \mathbb{R}^{r \times r}$ the inverse of $\left( B_{\cdot m} \quad (B_{\cdot j'})_{j' \in \mathcal{J}_{-j}} \right)$. One has*

$$Y_{\cdot j} = \mathcal{B}_{j \to m, \mathcal{J}_{-j}}[0] + \sum_{j' \in \mathcal{J}_{-j}} \mathcal{B}_{j \to m, \mathcal{J}_{-j}}[j'] Y_{\cdot j'} + \mathcal{B}_{j \to m, \mathcal{J}_{-j}}[m] Y_{\cdot m} + \zeta$$

*with:*

$$\mathcal{B}_{j\to m,\mathcal{J}_{-j}}[j'] := \sum_{k\in\{m\}\cup\mathcal{J}_{-j}} B_{kj'}^{-1}B_{jk}, \forall j' \in \mathcal{J}_{-j}$$

$$\mathcal{B}_{j\to m,\mathcal{J}_{-j}}[m] := \sum_{k\in\{m\}\cup\mathcal{J}_{-j}} B_{km}^{-1}B_{jk},$$

$$\mathcal{B}_{j\to m,\mathcal{J}_{-j}}[0] := \mathbf{1}\alpha_j - \sum_{j'\in\mathcal{J}_{-j}} \mathcal{B}_{j\to m,\mathcal{J}_{-j}}[j']\mathbf{1}\alpha_{j'} - \mathcal{B}_{j\to m,\mathcal{J}_{-j}}[m]\mathbf{1}\alpha_m$$

$$\zeta = -\sum_{j'\in\mathcal{J}_{-j}} \mathcal{B}_{j\to m,\mathcal{J}_{-j}}[j']\epsilon_{.j'} - \mathcal{B}_{j\to m,\mathcal{J}_{-j}}[m]\epsilon_{.m} + \epsilon_{.j}$$

*Proof.* Starting from the PPCA model written in (1) and recalled here

$$Y = \mathbf{1}\alpha + WB + \epsilon$$

and the matrix $B \in \mathbb{R}^{r\times p}$ being of full rank $r$, solving this linear system is the same as solving the following reduced system

$$\begin{pmatrix} Y_{.m} & (Y_{.j'})_{j'\in\mathcal{J}_{-j}} \end{pmatrix} = \mathbf{1}\alpha_{|r} + \begin{pmatrix} W_{.1} & \dots & W_{.r} \end{pmatrix} B_{|r} + \epsilon_{|r},$$

where $B_{|r} \in \mathbb{R}^{r\times r}$ denotes the reduced matrix $\begin{pmatrix} B_{.m} & (B_{.j'})_{j'\in\mathcal{J}_{-j}} \end{pmatrix}$ of $B$. Similarly, $\alpha_{|r} \in \mathbb{R}^r$ and $\epsilon_{|r} \in \mathbb{R}^{n\times r}$ denote the reduced matrices of $\alpha$ and $\epsilon$. With a slight abuse of notation, $B^{-1}$ denotes the inverse of the reduced matrix $\begin{pmatrix} B_{.m} & (B_{.j'})_{j'\in\mathcal{J}_{-j}} \end{pmatrix}$ which exists using **A1.**.

Then, one can derive that

$$\begin{pmatrix} W_{.1} & \dots & W_{.r} \end{pmatrix} = \left( \begin{pmatrix} Y_{.m} & (Y_{.j'})_{j'\in\mathcal{J}_{-j}} \end{pmatrix} - \mathbf{1}\alpha_{|r} - \epsilon_{|r} \right) B^{-1}.$$

The expression of $Y_{.j}$ as a function of the latent variables is

$$\begin{aligned} Y_{.j} &= \mathbf{1}\alpha_j + \begin{pmatrix} W_{.1} & \dots & W_{.r} \end{pmatrix} B_{j.} + \epsilon_{.j} \\ &= \mathbf{1}\alpha_j + \left( \begin{pmatrix} Y_{.m} & (Y_{.j'})_{j'\in\mathcal{J}_{-j}} \end{pmatrix} - \mathbf{1}\alpha_{|r} - \epsilon_{|r} \right) B^{-1} B_{j.} + \epsilon_{.j}, \end{aligned}$$

so that

$$Y_{.j} = \sum_{\ell\in\{m\}\cup\mathcal{J}_{-j}} \left( \sum_{k\in\{m\}\cup\mathcal{J}_{-j}} B_{k\ell}^{-1}B_{jk} \right) Y_{.\ell}$$

$$- \sum_{\ell\in\{m\}\cup\mathcal{J}_{-j}} \left( \sum_{k\in\{m\}\cup\mathcal{J}_{-j}} B_{k\ell}^{-1}B_{jk} \right)(\mathbf{1}\alpha_\ell + \epsilon_{.\ell}) + \epsilon_{.j} + \mathbf{1}\alpha_j.$$

which leads to the desired solution.

$\square$

## B.2  Proof of Proposition 4

**Proposition 4** (Mean estimator). *Consider the PPCA model* (1). *Under Assumptions **A1.** and **A2.**, an estimator of the mean of a MNAR variable $Y_{.m}$, for $m \in \mathcal{M}$, can be constructed as follows: choose $j \in \mathcal{J}$, and compute*

$$\hat{\alpha}_m := \frac{\hat{\alpha}_j - \hat{\mathcal{B}}_{j\to m,\mathcal{J}_{-j}}^c[0] - \sum_{j'\in\mathcal{J}_{-j}} \hat{\mathcal{B}}_{j\to m,\mathcal{J}_{-j}}^c[j']\hat{\alpha}_{j'}}{\hat{\mathcal{B}}_{j\to m,\mathcal{J}_{-j}}^c[m]},$$

*with the $(\hat{\mathcal{B}}_{j\to m,\mathcal{J}_{-j}}[k])$'s estimators of the coefficients given in Definition 3 and assuming that the coefficient $\mathcal{B}_{j\to m,\mathcal{J}_{-j}}^c[m]$ estimated by $\hat{\mathcal{B}}_{j\to m,\mathcal{J}_{-j}}^c[m]$ is non zero.*

*Under the additional Assumptions **A3.** and **A4.**, this estimator is consistent.*

19

*Proof.* The main goal is to obtain a formula for $\alpha_{.m}$, *i.e.*

$$\alpha_m = \frac{\alpha_j - \mathcal{B}^c_{j\to m,\mathcal{J}_{-j}}[0] - \sum_{j'\in\mathcal{J}_{-j}}\mathcal{B}^c_{j\to m,\mathcal{J}_{-j}}[j']\alpha_{j'}}{\mathcal{B}^c_{j\to m,\mathcal{J}_{-j}}[m]}, \tag{35}$$

from which an estimator can be deduced. The idea is to express $\alpha_j$ from $\alpha_m$ and $(\alpha_{j'})_{j'\in\mathcal{J}_{-j}}$. Note that $\mathbb{E}[Y_{.j}] = \mathbb{E}[\mathbb{E}[Y_{.j}|(Y_{.k})_{k\in\overline{\{j\}}}]]$. Assumption **A2.** leads to

$$\mathbb{E}[Y_{.j}|(Y_{.k})_{k\in\overline{\{j\}}}] = \mathbb{E}[Y_{.j}|(Y_{.k})_{k\in\overline{\{j\}}}, \Omega_{.m} = 1].$$

Then, by Definition 3 which gives $(Y_{.j})_{|\Omega_{.m}=1}$,

$$\mathbb{E}[Y_{.j}|(Y_{.k})_{k\in\overline{\{j\}}}, \Omega_{.m} = 1]$$

$$= \mathbb{E}\left[\mathcal{B}^c_{j\to m,\mathcal{J}_{-j}}[0] + \sum_{k\in\{m\}\cup\mathcal{J}_{-j}}\mathcal{B}^c_{j\to m,\mathcal{J}_{-j}}[k]Y_{.k} + \zeta^c\Big|(Y_{.k})_{k\in\overline{\{j\}}}\right]$$

$$= \mathcal{B}^c_{j\to m,\mathcal{J}_{-j}}[0] + \sum_{k\in\{m\}\cup\mathcal{J}_{-j}}\mathcal{B}^c_{j\to m,\mathcal{J}_{-j}}[k]Y_{.k} + \mathbb{E}\left[\zeta^c\Big|(Y_{.k})_{k\in\overline{\{j\}}}\right]$$

Thus, by taking the mean and given that $\mathbb{E}[\epsilon_{.k}] = 0, \forall k \in \{m\} \cup \mathcal{J}_{-j}$, one has

$$\alpha_j = \mathcal{B}^c_{j\to m,\mathcal{J}_{-j}}[0] + \sum_{j'\in\mathcal{J}_{-j}}\mathcal{B}^c_{j\to m,\mathcal{J}_{-j}}[j']\alpha_{j'} + \mathcal{B}^c_{j\to m,\mathcal{J}_{-j}}[m]\alpha_m,$$

implying Equation (35), provided that $\mathcal{B}^c_{j\to m,\mathcal{J}_{-j}}[m] \neq 0$.

From this formula for the mean $\alpha_m$, one define its estimator $\hat{\alpha}_m$ as in (9). It is trivially consistent as the linear combination of consistent quantities under **A3.** and **A4.** $\qquad\square$

## B.3  Proof of Proposition 5

**Proposition 5** (Variance and covariances estimators)**.** *Consider the PPCA model* (1). *Under Assumptions **A1.** and **A2.**, an estimator of the variance of a MNAR variable $Y_{.m}$ for $m \in \mathcal{M}$ and its covariances with the pivot variables, can be constructed as follows: choose $j \in \mathcal{J}$ and compute*

$$\left(\widehat{\mathrm{Var}}(Y_{.m}) \quad \widehat{\mathrm{Cov}}(Y_{.m}, (Y_{.k})_{k\in\mathcal{J}})\right)^T := (\widehat{M}_j)^{-1}\widehat{P}_j,$$

*assuming that $\sigma^2$ tends to zero and the inverse of the matrix $M_j$ estimated by $(\widehat{M}_j)^{-1}$ exists, with*

$$\widehat{M}_j = \begin{array}{c} \in\mathbb{R}\left\{\vphantom{\begin{bmatrix}\end{bmatrix}}\right. \\ \in\mathbb{R}^r\left\{\vphantom{\begin{bmatrix}\end{bmatrix}}\right. \end{array} \begin{bmatrix} \overbrace{(\hat{\mathcal{B}}^c_{j\to m,\mathcal{J}_{-j}}[m])^2}^{\in\mathbb{R}} & \overbrace{0 \quad 2\hat{\mathcal{B}}^c_{j\to m,\mathcal{J}_{-j}}[m]\left(\hat{\mathcal{B}}^c_{j\to m,\mathcal{J}_{-j}}[\mathcal{J}_{-j}]\right)^T}^{\in\mathbb{R}^r} \\ -(\hat{\mathcal{B}}^c_{k\to m,\mathcal{J}_{-k}}[m])_{k\in\mathcal{J}} & \left(\widehat{M}^k\right)_{k\in\mathcal{J}} \end{bmatrix}$$

*Let us precise that $\widehat{M}_j \in \mathbb{R}^{(r+1)\times(r+1)}$. One has $(\hat{\mathcal{B}}^c_{k\to m,\mathcal{J}_{-k}}[m])_{k\in\mathcal{J}} = \begin{pmatrix} \hat{\mathcal{B}}^c_{j_1\to m,\mathcal{J}_{-j_1}}[m] \\ \vdots \\ \hat{\mathcal{B}}^c_{j_r\to m,\mathcal{J}_{-j_r}}[m] \end{pmatrix}$.*

*One details $\widehat{M}^k$ for $k = j_1$ and the same definition is valid for all $k \in \mathcal{J}$.*

$$\widehat{M}^{j_1} = \begin{pmatrix} 1 & -\hat{\mathcal{B}}^c_{j_1\to m,\mathcal{J}_{-j_1}}[j_2] & \cdots & -\hat{\mathcal{B}}^c_{j_1\to m,\mathcal{J}_{-j_1}}[j_r] \end{pmatrix} \in \mathbb{R}^r$$

$$\widehat{P}_j = \begin{bmatrix} \overbrace{(\widehat{\mathrm{Var}}(Y_{.j}) - Q^c - (\hat{\mathcal{B}}^c_{j\to m,\mathcal{J}_{-j}}[\mathcal{J}_{-j}])^T\widehat{\mathrm{Var}}(Y_{\mathcal{J}_{-j}})\hat{\mathcal{B}}^c_{j\to m,\mathcal{J}_{-j}}[\mathcal{J}_{-j}])}^{\in\mathbb{R}} \\ \left(((\hat{\mathcal{B}}^c_{k\to m,\mathcal{J}_{-k}})^T \begin{pmatrix} 1 & \hat{\alpha}_m & (\hat{\alpha}_\ell)_{\ell\in\mathcal{J}_{-k}} \end{pmatrix}^T - \hat{\alpha}_k)\hat{\alpha}_m\right)_{k\in\mathcal{J}} \end{bmatrix} \begin{array}{l} \} \in\mathbb{R} \\ \} \in\mathbb{R}^r \end{array}$$

$$\hat{Q}^c = \left(\widehat{\mathrm{Var}}(Y_{.j})\big|\Omega_{.m}=1\right)$$
$$- \left(\widehat{\mathrm{Cov}}((Y_{.k})_{k\in\overline{\{j\}}}, Y_{.j})\widehat{\mathrm{Var}}((Y_{.k})_{k\in\overline{\{j\}}})^{-1}\widehat{\mathrm{Cov}}((Y_{.k})_{k\in\overline{\{j\}}}, Y_{.j})^T\big|\Omega_{.m}=1\right).$$

*Under the additional Assumptions **A3.** and **A4.**, the estimators for the variance of $Y_{.m}$ and its covariances with the pivot variables given in* (11) *are consistent.*

*Proof.* As for the mean, to derive some estimator of the variance and the covariances, we want to obtain a formula as

$$M_j \left(\mathrm{Var}(Y_{.m}) \quad \mathrm{Cov}(Y_{.m},(Y_{.k})_{k\in\mathcal{J}})\right)^T = \left(P_j - \mathcal{O}(\sigma^2)\right), \tag{36}$$

with

$$M_j = \begin{array}{l} \in\mathbb{R}\,\Big\{ \\ \in\mathbb{R}^r\,\Big\{ \end{array} \left[\begin{array}{cc} \overbrace{(\mathcal{B}^c_{j\to m,\mathcal{J}_{-j}}[m])^2}^{\in\mathbb{R}} & \overbrace{0 \quad 2\mathcal{B}^c_{j\to m,\mathcal{J}_{-j}}[m]\left(\mathcal{B}^c_{j\to m,\mathcal{J}_{-j}}[\mathcal{J}_{-j}]\right)^T}^{\in\mathbb{R}^r} \\ -(\mathcal{B}^c_{k\to m,\mathcal{J}_{-k}}[m])_{k\in\mathcal{J}} & (M^k)_{k\in\mathcal{J}} \end{array}\right]$$

Let us precise that $M_j \in \mathbb{R}^{(r+1)\times(r+1)}$. One has $(\mathcal{B}^c_{k\to m,\mathcal{J}_{-k}}[m])_{k\in\mathcal{J}} = \begin{pmatrix} \mathcal{B}^c_{j_1\to m,\mathcal{J}_{-j_1}}[m] \\ \vdots \\ \mathcal{B}^c_{j_r\to m,\mathcal{J}_{-j_r}}[m] \end{pmatrix}$.

One details $M^k$ for $k=j_1$ and the same definition is valid for all $k\in\mathcal{J}$.

$$M^{j_1} = \begin{pmatrix} 1 & -\mathcal{B}^c_{j_1\to m,\mathcal{J}_{-j_1}}[j_2] & \cdots & -\mathcal{B}^c_{j_1\to m,\mathcal{J}_{-j_1}}[j_r] \end{pmatrix} \in \mathbb{R}^r$$

$$P_j = \left[\begin{array}{c} \overbrace{(\mathrm{Var}(Y_{.j}) - Q^c - (\mathcal{B}^c_{j\to m,\mathcal{J}_{-j}}[\mathcal{J}_{-j}])^T\mathrm{Var}(Y_{\mathcal{J}_{-j}})\mathcal{B}^c_{j\to m,\mathcal{J}_{-j}}[\mathcal{J}_{-j}]}^{\in\mathbb{R}} \\ \left(((\mathcal{B}^c_{k\to m,\mathcal{J}_{-k}})^T \left(1 \quad \mathbb{E}[Y_{.m}] \quad (\mathbb{E}[Y_{.\ell}])_{\ell\in\mathcal{J}_{-k}}\right)^T - \mathbb{E}[Y_{.k}])\mathbb{E}[Y_{.m}]\right)_{k\in\mathcal{J}} \end{array}\right] \begin{array}{l}\Big\}\in\mathbb{R} \\ \Big\}\in\mathbb{R}^r\end{array}$$

$$\mathcal{O}(\sigma^2) = \left[\begin{array}{c} \overbrace{o_{\mathrm{var}}(\sigma^2)}^{\in\mathbb{R}} \\ -\left(o_{\mathrm{cov},k}(\sigma^2)\right)_{k\in\mathcal{J}} \end{array}\right] \begin{array}{l}\Big\}\in\mathbb{R} \\ \Big\}\in\mathbb{R}^r\end{array},$$

with $o_{\mathrm{var}}(\sigma^2)$ and $o_{\mathrm{cov},k}(\sigma^2)$ detailed in (42) and (45) respectively.

$$Q^c = \left(\mathrm{Var}(Y_{.j})\big|\Omega_{.m}=1\right)$$
$$- \left(\mathrm{Cov}((Y_{.k})_{k\in\overline{\{j\}}}, Y_{.j})\mathrm{Var}((Y_{.k})_{k\in\overline{\{j\}}})^{-1}\mathrm{Cov}((Y_{.k})_{k\in\overline{\{j\}}}, Y_{.j})^T\big|\Omega_{.m}=1\right). \tag{37}$$

The strategy is to prove each equality of the linear system in (36).

**Deriving an equation for the variance.** The idea is first to express $\mathrm{Var}(Y_{.j})$ from $\mathrm{Var}(Y_{.m})$, $(\mathrm{Var}(Y_{.j'}))_{j'\in\mathcal{J}_{-j}}$ and $(\mathrm{Cov}(Y_{.k},Y_{.\ell}))_{k\neq\ell\in\{m\}\cup\mathcal{J}_{-j}}$. The law of total variance reads as

$$\mathrm{Var}(Y_{.j}) = \mathbb{E}[\mathrm{Var}(Y_{.j}|Z)] + \mathrm{Var}(\mathbb{E}[Y_{.j}|Z]), \tag{38}$$

with $Z = (Y_{.k})_{k\in\overline{\{j\}}}$.

For the first term in (38), using Assumption **A2.**, one has

$$Y_{.j} \perp\!\!\!\perp (\Omega_{.m} = 1)|Z$$

which leads to

$$\text{Var}(Y_{.j}|Z) = \text{Var}(Y_{.j}|Z, \Omega_{.m} = 1).$$

The conditional variance for a Gaussian vector gives

$$\text{Var}(Y_{.j}|Z) = \text{Var}(Y_{.j}) - \text{Cov}(Z, Y_{.j})\text{Var}(Z)^{-1}\text{Cov}(Z, Y_{.j})^T,$$

implying that

$$\text{Var}(Y_{.j}|Z, \Omega_{.m} = 1) = \big(\text{Var}(Y_{.j}) - \text{Cov}(Z, Y_{.j})\text{Var}(Z)^{-1}\text{Cov}(Z, Y_{.j})^T\big|\Omega_{.m} = 1\big)$$

and then, as deterministic quantity,

$$\mathbb{E}[\text{Var}(Y_{.j}|Z)] = \big(\text{Var}(Y_{.j}) - \text{Cov}(Z, Y_{.j})\text{Var}(Z)^{-1}\text{Cov}(Z, Y_{.j})^T\big|\Omega_{.m} = 1\big).$$

One has

$$\text{Cov}(Z, Y_{.j})\text{Var}(Z)^{-1}\text{Cov}(Z, Y_{.j})^T =$$
$$\text{Cov}((Y_{.k})_{k\in\overline{\{j\}}}, Y_{.j})\text{Var}((Y_{.k})_{k\in\overline{\{j\}}})^{-1}\text{Cov}((Y_{.k})_{k\in\overline{\{j\}}}, Y_{.j})^T$$

leading to

$$\mathbb{E}[\text{Var}(Y_{.j}|Z)] = Q^c, \tag{39}$$

where $Q^c$ is defined in (37).

For the second term of (38), remark that **A2.** implies that

$$\text{Var}(\mathbb{E}[Y_{.j}|Z]) = \text{Var}(\mathbb{E}[Y_{.j}|Z, \Omega_{.m} = 1]),$$

and

$$\text{Var}(\mathbb{E}[Y_{.j}|Z, \Omega_{.m} = 1]) = \text{Var}\left(\mathbb{E}\left[\mathcal{B}^c_{j\to m, \mathcal{J}_{-j}}[0] + \sum_{k\in\{m\}\cup\mathcal{J}_{-j}}\mathcal{B}^c_{j\to m, \mathcal{J}_{-j}}[k]Y_{.k} + \zeta^c\bigg|Z\right]\right),$$

*i.e.*

$$\text{Var}(\mathbb{E}[Y_{.j}|Z, \Omega_{.m} = 1])$$
$$= \text{Var}\left(\sum_{k\in\{m\}\cup\mathcal{J}_{-j}}\mathcal{B}^c_{j\to m, \mathcal{J}_{-j}}[k]Y_{.k} - \sum_{k\in\{m\}\cup\mathcal{J}_{-j}}\mathcal{B}^c_{j\to m, \mathcal{J}_{-j}}[k]\mathbb{E}[\epsilon_{.k}|Z] + \mathcal{B}^c_{j\to m, \mathcal{J}_{-j}}[0] + \mathbb{E}[\epsilon_{.j}]\right)$$

In the variance, the first term is obtained using that the variables $(Y_{.k})_{k\in\{m\}\cup\mathcal{J}_{-j}}$ are $Z-$measurable. The two last terms use that $\mathcal{B}^c_{j\to m, \mathcal{J}_{-j}}[0]$ is a constant and $\epsilon_{.j}$ is independent of $Z$. To calculate the second term, involving $\mathbb{E}[\epsilon_{.k}|Z]$, one first shows that the vector $\big((Y_{.k})_{k\in\{m\}\cup\mathcal{J}_{-j}} \quad (\epsilon_{.k})_{k\in\{m\}\cup\mathcal{J}_{-j}}\big)^T$ is Gaussian. Indeed,

- $(Y_{.k})_{k\in\{m\}\cup\mathcal{J}_{-j}}$ is a Gaussian vector, using the model (1).

- $(\epsilon_{.k})_{k\in\{m\}\cup\mathcal{J}_{-j}}$ is a Gaussian vector, because its components are independent Gaussian variables.

- for $k \neq \ell \in \{m\} \cup \mathcal{J}_{-j}$, $(WB_{k.} \quad \epsilon_{.\ell})^T$ is a Gaussian vector, because $Y_{.k} \perp\!\!\!\perp \epsilon_{.\ell}$.

- for $k \in \{m\} \cup \mathcal{J}_{-j}$, $(Y_{.k} \quad \epsilon_{.k})^T$ is a Gaussian vector, given that $Y_{.k}$ is a linear combination of $(WB_{k.} \quad \epsilon_{.k})^T$ which is Gaussian, as $WB_{k.}$ and $\epsilon_{.k}$ are independent Gaussian variables.

Thus,

$$\mathbb{E}[\epsilon_{.k}|Z] = \mathbb{E}[\epsilon_{.k}] + \text{Cov}(\epsilon_{.k}, Z)\text{Var}(Z)^{-1}(Z - \mathbb{E}[Z])$$
$$= \text{Cov}(\epsilon_{.k}, Y_{.k})(\text{Var}(Z)^{-1})_{k.}(Z - \mathbb{E}[Z]),$$

using $\text{Cov}(\epsilon_{.k}, Y_{.l}) = 0$, for $k \neq l$. $\Gamma_Z = \text{Var}(Z)^{-1}$ denotes the inverse of the covariance matrix of $Z$ and $(\Gamma_Z)_{k.}$ is its k-th row. It leads to

$$\mathbb{E}[\epsilon_{.k}|Z] = \sigma^2 (\Gamma_Z)_{k.}(Z - \mathbb{E}[Z]). \tag{40}$$

given that $\text{Cov}(\epsilon_{.k}, Y_{.k}) = \text{Cov}(\epsilon_{.k}, WB_{k.} + \epsilon_{.k}) = \text{Var}(\epsilon_{.k})$.

Therefore,

$$\text{Var}(\mathbb{E}[Y_{.j}|Z, \Omega_{.m} = 1]) = \sum_{k \in \{m\} \cup \mathcal{J}_{-j}} (\mathcal{B}^c_{j \to m, \mathcal{J}_{-j}[k]})^2 \text{Var}(Y_{.k})$$

$$+ \sum_{(k < \ell) \in \{m\} \cup \mathcal{J}_{-j}} 2\mathcal{B}^c_{j \to m, \mathcal{J}_{-j}[k]} \mathcal{B}^c_{j \to m, \mathcal{J}_{-j}[\ell]} \text{Cov}(Y_{.k}, Y_{.\ell}) + o_{\text{var}}(\sigma^2), \tag{41}$$

where

$$o_{\text{var}}(\sigma^2) = -2\sigma^2 \sum_{(k,\ell) \in \{m\} \cup \mathcal{J}_{-j}} \mathcal{B}^c_{j \to m, \mathcal{J}_{-j}[k]} \mathcal{B}^c_{j \to m, \mathcal{J}_{-j}[\ell]} \sum_{\ell' \in \{m\} \cup \mathcal{J}_{-j}} (\Gamma_Z)_{\ell\ell'} \text{Cov}(Y_{.k}, Y_{.\ell'})$$

$$+ \sigma^4 \sum_{k \in \{m\} \cup \mathcal{J}_{-j}} (\mathcal{B}^c_{j \to m, \mathcal{J}_{-j}[k]})^2 \left( \sum_{(\ell < \ell') \in \{m\} \cup \mathcal{J}_{-j}} (\Gamma_Z)^2_{k\ell} \text{Var}(Y_{.\ell}) - 2(\Gamma_Z)_{k\ell}(\Gamma_Z)_{k\ell'} \text{Cov}(Y_{.\ell}, Y_{.\ell'}) \right)$$

$$- 2\sigma^4 \sum_{(k < \ell) \in \{m\} \cup \mathcal{J}_{-j}} \mathcal{B}^c_{j \to m, \mathcal{J}_{-j}[k]} \mathcal{B}^c_{j \to m, \mathcal{J}_{-j}[\ell]} \sum_{(k', \ell') \in \{m\} \cup \mathcal{J}_{-j}} (\Gamma_Z)_{kk'}(\Gamma_Z)_{\ell\ell'} \text{Cov}(Y_{.k'}, Y_{.\ell'}) \tag{42}$$

Combining (39) with (41), one get the following expression for the first line of the linear system

$$(\mathcal{B}^c_{j \to m, \mathcal{J}_{-j}[m]})^2 \text{Var}(Y_{.m}) + \sum_{j' \in \mathcal{J}_{-j}} 2\mathcal{B}^c_{j \to m, \mathcal{J}_{-j}[j']} \mathcal{B}^c_{j \to m, \mathcal{J}_{-j}[m]} \text{Cov}(Y_{.j'}, Y_{.m})$$

$$= \text{Var}(Y_{.j}) - Q^c - (\mathcal{B}^c_{j \to m, \mathcal{J}_{-j}[\mathcal{J}_{-j}]})^T \text{Var}(Y_{\mathcal{J}_{-j}}) \mathcal{B}^c_{j \to m, \mathcal{J}_{-j}[\mathcal{J}_{-j}]} - o_{\text{var}}(\sigma^2) \tag{43}$$

**Deriving equations for the covariances.** Let $k$ be an element of $\mathcal{J}$, our objective is to express $\text{Cov}(Y_{.m}, Y_{.k})$ from $\text{Var}(Y_{.m})$, $\alpha_m$, $(\alpha_k)_{k \in \mathcal{J}}$ and $(\text{Cov}(Y_{.m}, Y_{.k}))_{k \in \{m\} \cup \mathcal{J}}$.

$$\text{Cov}(Y_{.m}, Y_{.k}) = \mathbb{E}[Y_{.m}Y_{.k}] - \mathbb{E}[Y_{.m}]\mathbb{E}[Y_{.k}]$$
$$= \mathbb{E}[\mathbb{E}[Y_{.m}Y_{.k}|Z]] - \mathbb{E}[Y_{.m}]\mathbb{E}[Y_{.k}]$$
$$= \mathbb{E}[Y_{.m}\mathbb{E}[Y_{.k}|Z]] - \mathbb{E}[Y_{.m}]\mathbb{E}[Y_{.k}], \tag{44}$$

with $Z = (Y_{.\ell})_{\ell \in \overline{\{k\}}}$.

For the first term in (44), one has

$$\mathbb{E}[Y_{.m}\mathbb{E}[Y_{.k}|Z]] \overset{(i)}{=} \mathbb{E}[Y_{.m}\mathbb{E}[Y_{.k}|Z, \Omega_{.m} = 1]]$$

$$\overset{(ii)}{=} \mathbb{E}\left[ Y_{.m}\left( \mathcal{B}^c_{k \to m, \mathcal{J}_{-k}[0]} + \sum_{\ell \in \{m\} \cup \mathcal{J}_{-k}} \mathcal{B}^c_{k \to m, \mathcal{J}_{-k}[\ell]} Y_{.\ell} + \mathbb{E}[\zeta^c_k|Z] \right) \right]$$

$$\overset{(iii)}{=} \mathcal{B}^c_{k \to m, \mathcal{J}_{-k}[0]} \mathbb{E}[Y_{.m}] + \mathcal{B}^c_{k \to m, \mathcal{J}_{-k}[m]} \mathbb{E}[Y^2_{.m}]$$

$$+ \sum_{\ell \in \mathcal{J}_{-k}} \mathcal{B}^c_{k \to m, \mathcal{J}_{-k}[\ell]} \mathbb{E}[Y_{.m}Y_{.\ell}] + o_{\text{cov},k}(\sigma^2)$$

with $\zeta^c_k = -\sum_{\ell \in \mathcal{J}_{-k}} \mathcal{B}^c_{k \to m, \mathcal{J}_{-k}[\ell]} \epsilon_{.\ell} - \mathcal{B}^c_{k \to m, \mathcal{J}_{-k}[m]} \epsilon_{.m} + \epsilon_{.k}$.

Assumption **A2.** and Definition 3 are used for (i) and (ii) respectively. For (iii), using (40), one has

$$\mathbb{E}[Y_{.m}\mathbb{E}[\zeta^c_k|Z]] = \mathbb{E}\left[ Y_{.m}\left( -\sum_{\ell \in \{m\} \cup \mathcal{J}_{-k}} \mathcal{B}^c_{k \to m, \mathcal{J}_{-k}[\ell]} \sigma^2 (\Gamma_Z)_{\ell.}(Z - \mathbb{E}[Z]) \right) \right],$$

given that $\mathbb{E}[\epsilon_{.k}|Z] = \mathbb{E}[\epsilon_{.k}] = 0$ by independence.

$$\mathbb{E}[Y_{.m}\mathbb{E}[\zeta^c_k|Z]]$$

$$= -\sigma^2 \mathbb{E}\left[ \sum_{\ell \in \{m\} \cup \mathcal{J}_{-k}} \mathcal{B}^c_{k \to m, \mathcal{J}_{-k}[\ell]} Y_{.m} \sum_{\ell' \{m\} \cup \in \mathcal{J}_{-k}} (\Gamma_Z)_{\ell\ell'} (Y_{.\ell'} - \mathbb{E}[Y_{.\ell'}]) \right].$$

In addition,

$$\mathbb{E}\left[\sum_{\ell\in\{m\}\cup\mathcal{J}_{-k}}\mathcal{B}^c_{k\to m,\mathcal{J}_{-k}[\ell]}Y_{.m}\sum_{\ell'\in\{m\}\cup\mathcal{J}_{-k}}(\Gamma_Z)_{\ell\ell'}(Y_{.\ell'}-\mathbb{E}[Y_{.\ell'}])\right]$$

$$=\sum_{\ell\in\{m\}\cup\mathcal{J}_{-k}}\sum_{\ell'\in\{m\}\cup\mathcal{J}_{-k}}(\Gamma_Z)_{\ell\ell'}\mathcal{B}^c_{k\to m,\mathcal{J}_{-k}[\ell]}\left(\mathrm{Cov}\left(Y_{.m},Y_{.\ell'}\right)+\mathbb{E}[Y_{.m}]\mathbb{E}[(Y_{.\ell'}-\mathbb{E}[Y_{.\ell'}])]\right)$$

$$=\sum_{\ell\in\{m\}\cup\mathcal{J}_{-k}}\sum_{\ell'\in\{m\}\cup\mathcal{J}_{-k}}(\Gamma_Z)_{\ell\ell'}\mathcal{B}^c_{k\to m,\mathcal{J}_{-k}[\ell]}\mathrm{Cov}\left(Y_{.m},Y_{.\ell'}\right)$$

It implies that, in (iii),

$$o_{\mathrm{cov},k}(\sigma^2)=-\sigma^2\sum_{\ell\in\{m\}\cup\mathcal{J}_{-k}}\sum_{\ell'\in\{m\}\cup\mathcal{J}_{-k}}(\Gamma_Z)_{\ell\ell'}\mathcal{B}^c_{k\to m,\mathcal{J}_{-k}[\ell]}\mathrm{Cov}\left(Y_{.m},Y_{.\ell'}\right)\qquad(45)$$

Equation (44) leads thus to

$$\mathrm{Cov}(Y_{.m},Y_{.k})=\mathcal{B}^c_{k\to m,\mathcal{J}_{-k}[0]}\mathbb{E}[Y_{.m}]+\mathcal{B}^c_{k\to m,\mathcal{J}_{-k}[m]}(\mathrm{Var}(Y_{.m})+\mathbb{E}[Y_{.m}]^2)$$

$$+\sum_{\ell\in\mathcal{J}_{-k}}\mathcal{B}^c_{k\to m,\mathcal{J}_{-k}[\ell]}(\mathrm{Cov}(Y_{.m},Y_{.\ell})+\mathbb{E}[Y_{.m}]\mathbb{E}[Y_{.\ell}])-\mathbb{E}[Y_{.m}]\mathbb{E}[Y_{.k}]+o_{\mathrm{cov},k}(\sigma^2),\quad(46)$$

which can be rewritten as

$$\mathrm{Cov}(Y_{.m},Y_{.k})-\mathcal{B}^c_{k\to m,\mathcal{J}_{-k}[m]}\mathrm{Var}(Y_{.m})-\sum_{\ell\in\mathcal{J}_{-k}}\mathcal{B}^c_{k\to m,\mathcal{J}_{-k}[\ell]}\mathrm{Cov}(Y_{.m},Y_{.\ell})$$

$$=((\mathcal{B}^c_{k\to m,\mathcal{J}_{-k}})^T\begin{pmatrix}1 & \mathbb{E}[Y_{.m}] & (\mathbb{E}[Y_{.\ell}])_{\ell\in\mathcal{J}_{-k}}\end{pmatrix}^T-\mathbb{E}[Y_{.k}])\mathbb{E}[Y_{.m}]+o_{\mathrm{cov},k}(\sigma^2),\quad(47)$$

Combining Equations (43) and (47) forms the desired matrix system (36).

From these formulae for $(\mathrm{Var}(Y_{.m})\quad \mathrm{Cov}(Y_{.m},(Y_{.k})_{k\in\mathcal{J}}))^T$, assuming that $M_j$ is invertible and that $\sigma^2$ tends to zero, one get their estimators $(\widehat{\mathrm{Var}}(Y_{.m})\quad \widehat{\mathrm{Cov}}(Y_{.m},(Y_{.k})_{k\in\mathcal{J}}))^T$ defined in (10).

As for the consistency, $\hat{\alpha}_m$ is a consistent estimator for $\alpha_m$ by using Proposition 4. The estimators in (10) are consistent, under Assumption **A3.** and **A4.**. $\qquad\square$

## B.4 Proof of Proposition 8

For deriving the covariance between a MNAR variable and a MNAR or not pivot variable, we assume the following

**A5.** $\forall m\in\mathcal{M},\forall\ell\in\bar{\mathcal{J}}_{-m}$, for all set $\mathcal{H}\subset\mathcal{J}_{-j}$ such that $|\mathcal{H}|=r-2$, $(B_{.m}\quad B_{.\ell}\quad (B_{.j'})_{j'\in\mathcal{H}})$ is invertible,

**A6.** $\forall k\in\bar{\mathcal{J}}\backslash\mathcal{M},\forall j\in\mathcal{J},\ Y_{.j}\perp\!\!\!\perp\Omega_{.k}|(Y_{.\ell})_{\ell\in\overline{\{j\}}}$.

**A7.** $\forall k,\ell\in\bar{\mathcal{J}},\quad k\neq l,\ \Omega_{.k}\perp\!\!\!\perp\Omega_{.\ell}|Y$

**A8.** $\forall j\in\mathcal{J},\forall m\in\mathcal{M},\forall\ell\in\bar{\mathcal{J}}_{-m}$, for all set $\mathcal{H}\subset\mathcal{J}_{-j}$ such that $|\mathcal{H}|=r-2$, the complete-case coefficients $\mathcal{B}^c_{j\to m,\ell,\mathcal{H}[0]}$ and $\mathcal{B}^c_{j\to m,\ell,\mathcal{H}[k]},k\neq j,k\in\{m,\ell\}\cup\mathcal{H}$ can be consistently estimated. (Here, note that the complete case is when $\Omega_{.m}=1$ and $\Omega_{.\ell}=1$.)

**A9.** For the variables neither MNAR nor pivot, their means $(\alpha_k)_{k\in\bar{\mathcal{J}}\backslash\mathcal{M}}$, variances $(\mathrm{Var}(Y_{.k}))_{k\in\bar{\mathcal{J}}\backslash\mathcal{M}}$ and covariances $(\mathrm{Cov}(Y_{.k},Y_{.k'}))_{k\neq k'\in\bar{\mathcal{J}}\backslash\mathcal{M}}$ can be consistently estimated. The covariances between these variables and the pivot variables $(\mathrm{Cov}(Y_{.j},Y_{.k}))_{j\in\mathcal{J},k\in\bar{\mathcal{J}}\backslash\mathcal{M}}$ are also consistent.

**Proposition 8** (Covariance between a MNAR variable and a MNAR or not pivot variable). *Consider the PPCA model* (1). *Under Assumptions **A2.**, **A5.**, **A6.** and **A7.**, an estimator of the covariance between a MNAR variable $Y_{.m}$, for $m\in\mathcal{M}$, and a variable $Y_{.\ell}$, for $\ell\in\bar{\mathcal{J}}\backslash\{m\}$, can be constructed*

*as follows: choose $j \in \mathcal{J}$ and $r - 2$ variable indexes in $\mathcal{J}_{-j}$ and compute:*

$$\widehat{\mathrm{Cov}}(Y_{.m}, Y_{.\ell}) = \frac{1}{\hat{K}} \widehat{\mathrm{Var}}(Y_{.j}) - \hat{q}^c - \sum_{k \in \{m,\ell\} \cup \mathcal{H}} (\hat{\mathcal{B}}^c_{j \to m,\ell,\mathcal{H}[k]})^2 \widehat{\mathrm{Var}}(Y_{.k})$$

$$- \sum_{k < k', k \in \{m,\ell\} \cup \mathcal{H}, k' \in \mathcal{H}} 2\hat{\mathcal{B}}^c_{j \to m,\ell,\mathcal{H}[k]} \hat{\mathcal{B}}^c_{j \to m,\ell,\mathcal{H}[k']} \widehat{\mathrm{Cov}}(Y_{.k}, Y_{.k'}), \quad (48)$$

*assuming that $\sigma^2$ tends to zero and with $\hat{K} = 2\hat{\mathcal{B}}^c_{j \to m,\ell,\mathcal{H}[m]} \hat{\mathcal{B}}^c_{j \to m,\ell,\mathcal{H}[\ell]}$ and*

$$\hat{q}^c = \left( \widehat{\mathrm{Var}}(Y_{.j}) \middle| \Omega_{.m} = 1, \Omega_{.\ell} = 1 \right)$$
$$- \left( \widehat{\mathrm{Cov}}((Y_{.k})_{k \in \overline{\{j\}}}, Y_{.j}) \widehat{\mathrm{Var}}((Y_{.k})_{k \in \overline{\{j\}}})^{-1} \widehat{\mathrm{Cov}}((Y_{.k})_{k \in \overline{\{j\}}}, Y_{.j})^T \middle| \Omega_{.m} = 1, \Omega_{.\ell} = 1 \right),$$

*given that $K$ estimated by $\hat{K}$ is non zero.*

*Under the additional Assumptions **A3.**, **A8.** and **A9.**. this estimator given in (48) is consistent.*

*Proof.* Let $\mathcal{H}$ be the set of the $r - 2$ variable indexes. One has $\mathcal{H} \subset \mathcal{J}_{-j}$. We use the same strategy as the proof for Proposition 5 (paragraph for deriving an equation for the variance).

To derive a formula for $\mathrm{Cov}(Y_{.m}, Y_{.\ell})$ with $m \in \mathcal{M}$ and $\ell \in \bar{\mathcal{J}}_{-m}$, the idea is to express $\mathrm{Var}(Y_{.j})$ from $(\mathrm{Var}(Y_{.k}))_{k \in \{m,\ell\} \cup \mathcal{H}}$ and $(\mathrm{Cov}(Y_{.k}, Y_{.k'}))_{k \neq k' \in \{m,\ell\} \cup \mathcal{H}}$.

The law of total variance reads as

$$\mathrm{Var}(Y_{.j}) = \mathbb{E}[\mathrm{Var}(Y_{.j}|Z)] + \mathrm{Var}(\mathbb{E}[Y_{.j}|Z]), \quad (49)$$

with $Z = (Y_{.k})_{k \in \overline{\{j\}}}$.

For the first term in (49), one uses

$$Y_{.j} \perp\!\!\!\perp \Omega_{.m}, \Omega_{.\ell} | Z.$$

If $Y_{.m}$ and $Y_{.\ell}$ are both MNAR variables, this conditional independence is obtained using Assumption **A2.** and **A7.**. Otherwise, if $Y_{.\ell}$ is not a MNAR variable, Assumption **A6.** and **A7.** lead to the desired result. It implies

$$\mathrm{Var}(Y_{.j}|Z) = \mathrm{Var}(Y_{.j}|Z, \Omega_{.m} = 1, \Omega_{.\ell} = 1).$$

The conditional variance for a Gaussian vector gives

$$\mathrm{Var}(Y_{.j}|Z) = \mathrm{Var}(Y_{.j}) - \mathrm{Cov}(Z, Y_{.j})\mathrm{Var}(Z)^{-1}\mathrm{Cov}(Z, Y_{.j})^T,$$

implying that

$$\mathrm{Var}(Y_{.j}|Z, \Omega_{.m} = 1, \Omega_{.\ell} = 1)$$
$$= \left( \mathrm{Var}(Y_{.j}) - \mathrm{Cov}(Z, Y_{.j})\mathrm{Var}(Z)^{-1}\mathrm{Cov}(Z, Y_{.j})^T \middle| \Omega_{.m} = 1, \Omega_{.\ell} = 1 \right)$$

and then, as deterministic quantity,

$$\mathbb{E}[\mathrm{Var}(Y_{.j}|Z)] = q^c \quad (50)$$

with

$$q^c = \left( \mathrm{Var}(Y_{.j}) \middle| \Omega_{.m} = 1, \Omega_{.\ell} = 1 \right)$$
$$- \left( \mathrm{Cov}((Y_{.k})_{k \in \overline{\{j\}}}, Y_{.j})\mathrm{Var}((Y_{.k})_{k \in \overline{\{j\}}})^{-1}\mathrm{Cov}((Y_{.k})_{k \in \overline{\{j\}}}, Y_{.j})^T \middle| \Omega_{.m} = 1, \Omega_{.\ell} = 1 \right).$$

For the second term of (38), remark that **A2.**, **A6.** and **A7.** implies that

$$\mathrm{Var}(\mathbb{E}[Y_{.j}|Z]) = \mathrm{Var}(\mathbb{E}[Y_{.j}|Z, \Omega_{.m} = 1, \Omega_{.\ell} = 1]),$$

and

$$\mathrm{Var}(\mathbb{E}[Y_{.j}|Z, \Omega_{.m} = 1, \Omega_{.\ell} = 1])$$
$$= \mathrm{Var}\left( \mathbb{E}\left[ \mathcal{B}^c_{j \to m,\ell,\mathcal{H}[0]} + \sum_{k \in \{m,\ell\} \cup \mathcal{H}} \mathcal{B}^c_{j \to m,\ell,\mathcal{H}[k]} Y_{.k} + \zeta^c_j \middle| Z \right] \right),$$

*i.e.*

$$\text{Var}(\mathbb{E}[Y_{.j}|Z, \Omega_{.m} = 1, \Omega_{.\ell} = 1])$$

$$= \text{Var}\left(\sum_{k \in \{m,\ell\} \cup \mathcal{H}} \mathcal{B}^c_{j \to m,\ell,\mathcal{H}[k]} Y_{.k} - \sum_{k \in \{m,\ell\} \cup \mathcal{H}} \mathcal{B}^c_{j \to m,\ell,\mathcal{H}[k]} \mathbb{E}[\epsilon_{.k}|Z] + \mathcal{B}^c_{j \to m,\ell,\mathcal{H}[0]} + \mathbb{E}[\epsilon_{.j}]\right)$$

One uses the same reasoning as in the proof of Proposition 5 (paragraph for deriving an equation for the variance) to get

$$\text{Var}(\mathbb{E}[Y_{.j}|Z, \Omega_{.m} = 1, \Omega_{.\ell} = 1]) = \sum_{k \in \{m,\ell\} \cup \mathcal{H}} (\mathcal{B}^c_{j \to m,\ell,\mathcal{H}[k]})^2 \text{Var}(Y_{.k})$$

$$+ \sum_{k < k' \in \{m,\ell\} \cup \mathcal{H}} 2\mathcal{B}^c_{j \to m,\ell,\mathcal{H}[k]} \mathcal{B}^c_{j \to m,\ell,\mathcal{H}[k']} \text{Cov}(Y_{.k}, Y_{.k'}) + o_{\text{covmiss}}(\sigma^2), \quad (51)$$

where

$$o_{\text{covmiss}}(\sigma^2) = -2\sigma^2 \sum_{(k,k') \in \{m,\ell\} \cup \mathcal{H}} \mathcal{B}^c_{j \to m,\ell,\mathcal{H}[k]} \mathcal{B}^c_{j \to m,\ell,\mathcal{H}[k']} \sum_{\ell' \in \{m,\ell\} \cup \mathcal{H}} (\Gamma_Z)_{k'\ell'} \text{Cov}(Y_{.k}, Y_{.\ell'})$$

$$+ \sigma^4 \sum_{k \in \{m,\ell\} \cup \mathcal{H}} (\mathcal{B}^c_{j \to m,\ell,\mathcal{H}[k]})^2 \left(\sum_{(k' < \ell') \in \{m,\ell\} \cup \mathcal{H}} (\Gamma_Z)^2_{kk'} \text{Var}(Y_{.k'}) - 2(\Gamma_Z)_{kk'}(\Gamma_Z)_{k\ell'} \text{Cov}(Y_{.k'}, Y_{.\ell'})\right)$$

$$- 2\sigma^4 \sum_{(k < k') \in \{m,\ell\} \cup \mathcal{H}} \mathcal{B}^c_{j \to m,\ell,\mathcal{H}[k]} \mathcal{B}^c_{j \to m,\ell,\mathcal{H}[k']} \sum_{(k'',\ell') \in \{m,\ell\} \cup \mathcal{H}} (\Gamma_Z)_{kk''}(\Gamma_Z)_{k'\ell'} \text{Cov}(Y_{.k''}, Y_{.\ell'}) \quad (52)$$

Combining (49), (50) and (51), one get the following formula for $\text{Cov}(Y_{.m}, Y_{.\ell})$,

$$2\mathcal{B}^c_{j \to m,\ell,\mathcal{H}[m]} \mathcal{B}^c_{j \to m,\mathcal{H}[\ell]} \text{Cov}(Y_{.m}, Y_{.\ell}) = \text{Var}(Y_{.j}) - q^c - \sum_{k \in \{m,\ell\} \cup \mathcal{H}} (\mathcal{B}^c_{j \to m,\ell,\mathcal{H}[k]})^2 \text{Var}(Y_{.k})$$

$$- \sum_{k < k', k \in \{m,\ell\} \cup \mathcal{H}, k' \in \mathcal{H}} 2\mathcal{B}^c_{j \to m,\ell,\mathcal{H}[k]} \mathcal{B}^c_{j \to m,\ell,\mathcal{H}[k']} \text{Cov}(Y_{.k}, Y_{.k'}) - o_{\text{covmiss}}(\sigma^2)$$

An estimator of $\text{Cov}(Y_{.m}, Y_{.\ell})$ is then derived as in (48), given that $\sigma^2$ tends to zero and $K = \mathcal{B}^c_{j \to m,\ell,\mathcal{H}[m]} \mathcal{B}^c_{j \to m,\ell,\mathcal{H}[\ell]}$ is non zero.

We use the consistent estimators defined in Proposition 5 for estimating $\text{Var}(Y_{.m})$ and $\text{Cov}(Y_{.m}, Y_{.k})_{k \in \mathcal{H}}$. If $Y_{.\ell}$ is also a MNAR variable, Proposition 5 is applied for estimating $\text{Var}(Y_{.\ell})$ and $\text{Cov}(Y_{.\ell}, Y_{.k})_{k \in \mathcal{H}}$. Otherwise, if $Y_{.\ell}$ is not a MNAR variable, we use **A9.**.

Eventually, **A3.** and **A8.** lead to the consistency of $\widehat{\text{Cov}}(Y_{.m}, Y_{.\ell})$. $\qquad \square$

## B.5   Extension to more general mechanisms for the not MNAR variables

The results of Proposition 4, 5 and 8 can be extended to a more general setting than the one presented in Section 2. The pivot variables may be assumed to be MCAR (or observed). The variables which are neither MNAR nor pivot may be observed or satisfying

$$\forall \ell \in \bar{\mathcal{J}} \backslash \mathcal{M}, \forall i \in \{1, \ldots, n\}, \quad \mathbb{P}(\Omega_{i\ell} = 1|Y_{i.}) = \mathbb{P}(\Omega_{i\ell} = 1|(Y_{ik})_{k \in \bar{\mathcal{J}} \backslash \{\ell\} \cup \mathcal{M}}), \quad (53)$$

*i.e.* they are MCAR or MAR but their missing-data mechanisms may not depend on the pivot variables.

The proofs are similar and not presented here for the sake of brevity.

Note that the main difference is that the complete case has to be extended. For instance, for $j \in \mathcal{J}$ and $k \in \mathcal{J}_{-j}$, the coefficients standing respectively for the intercept and the effects of $Y_{.j}$ on $(Y_{.m}, (Y_{.j'})_{j' \in \mathcal{J}_{-j}})$ in the complete case, *i.e.* when $\Omega_{.m} = 1, (\Omega_{.j} = 1)_{j \in \mathcal{J}}$ are in this general setting defined as follows

$$(Y_{.j})_{|\Omega_{.m}=1,(\Omega_j=1)_{j \in \mathcal{J}}} := \mathcal{B}^c_{j \to m,\mathcal{J}_{-j}[0]} + \sum_{j' \in \mathcal{J}_{-j}} \mathcal{B}^c_{j \to m,\mathcal{J}_{-j}[j']} Y_{.j'} + \mathcal{B}^c_{j \to m,\mathcal{J}_{-j}[m]} Y_{.m} + \zeta^c,$$

with $\zeta^c = -\sum_{j' \in \mathcal{J}_{-j}} \mathcal{B}^c_{j \to m,\mathcal{J}_{-j}[j']} \epsilon_{.j'} - \mathcal{B}^c_{j \to m,\mathcal{J}_{-j}[m]} \epsilon_{.m} + \epsilon_{.j}$.

# C Other numerical experiments

**Robustness to noise.** Considering the same setting as in Section 4.1 ($n = 1000$, $p = 10$, $r = 2$ and seven self-masked MNAR variables), the methods are tried for different noise levels $\sigma^2 \in \{0.1, 0.3, 0.5, 0.7, 1\}$. The results are presented for one missing variable and for all the other ones, the results are similar. In Figure 7, Algorithm 1 is the only method that does not give a biased estimate of the mean and the variance regardless of the noise level. In Figure 8, despite a larger bias in the estimation of the covariance between a missing variable and a pivot one as the noise level increases, Algorithm 1 outperforms all the other methods, regarding the estimation of the covariance between two missing variables. Note that the formula for the estimate of the covariance between two missing variables relies on the one for the estimate of the variance, but both differ from the one used for the covariance estimation between a missing variable and a pivot one. As expected, in Figure 9, estimation deteriorates as the data gets noisier and then the loading matrix estimation and the imputation error get closer to the results of mean imputation. In term of imputation error, the proposed method yet remains competitive in regards of the approaches (ii) and (iii). Overall, when the noise level increases, the exogeneity will be worse and that ignoring it in practice can be made to the detriment of performance.



Figure 7: Mean estimation (left graphic) and variance estimation (right graphic) of one missing variable for different values of the level of noise when $r = 2$, $n = 1000$, $p = 10$ and seven variables are MNAR. True values to be estimated are indicated by red lines.



Figure 8: Covariance estimation beetween a missing variable and a pivot one (left graphic) and two missing variables (right graphic) for different values of the level of noise when $r = 2$, $n = 1000$, $p = 10$ and seven variables are MNAR. True values to be estimated are indicated by red lines.
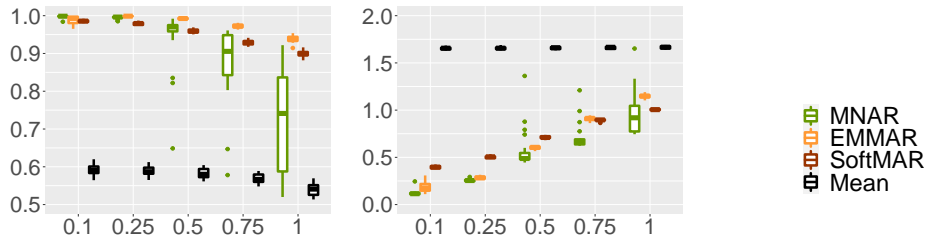


Figure 9: RV coefficients for the loading matrix (left graphic) and imputation error (right graphic) for different values of the level of noise when $r = 2$, $n = 1000$, $p = 10$ and seven variables are MNAR.

**Varying the percentage of missing values.** Considering the same setting as in Section 4.1 ($n = 1000$, $p = 10$, $r = 2$, $\sigma = 0.1$ and seven self-masked MNAR variables), the methods are tried for different percentages of missing values (10%, 30%, 50%). The results are presented in Figure 10. As expected, all the methods deteriorate with an increasing percentage of missing values but our method is stable.
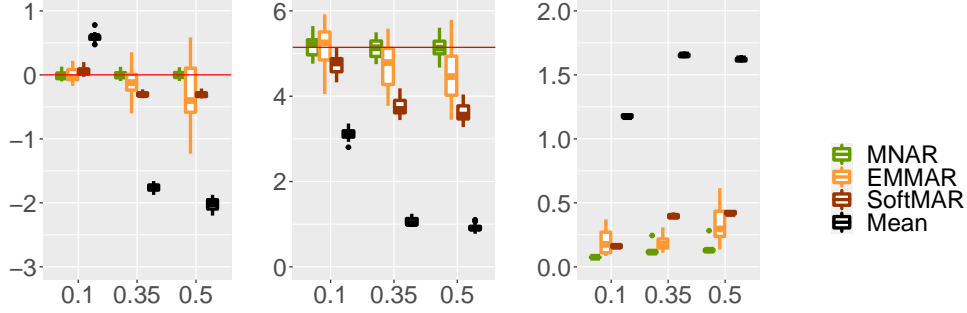


Figure 10: Mean estimation (left graphic), variance estimation (middle graphic) and imputation error (right graphic) for different percentages of missing values when $r = 2$, $n = 1000$, $p = 10$ and seven variables are MNAR.

**Misspecification to the rank.** The misspecification to the parameter $r$ has been evaluated: under a model generated with $r = 3$ latent variables ($n = 1000$, $p = 20$, $\sigma = 0.8$ and ten MNAR self-masked variables), the rank is either underestimated, well estimated or overestimated by giving to Algorithm 1 the information that $r = 2$, $r = 3$ or $r = 4$. Both estimation of the loading matrix and imputation error are shown in Figure 11. The results for an underestimated ($r = 2$) or overestimated ($r = 4$) rank are comparable to the case where the accurate rank is considered instead ($r = 3$), showing a stability of Algorithm 1 to rank misspecification.



Figure 11: RV coefficients for the loading matrix (left) and imputation error (right) when $r = 3$, $n = 1000$, $p = 20$ and ten variables are MNAR for different cases where the rank is either underestimated, well estimated or overestimated.

**General MNAR mechanism.** We consider the setting $n = 1000$, $p = 20$, $r = 3$ and $\sigma = 0.8$. Here, missing values are introduced on ten variables $(Y_{.k})_{k \in [1:10]}$ using a more general MNAR mechanism (see (3)) than the self-masked one. In particular, the MNAR mechanism we consider is defined as follows,

$$\forall m \in [1:10], \forall i \in \{1, \ldots, n\}, \ \mathbb{P}(\Omega_{im} = 1 | Y_{i.}) = \mathbb{P}(\Omega_{im} = 1 | Y_{im}, Y_{ik}, Y_{i\ell}), \qquad (54)$$

where $k$ and $\ell$ are indexes of MNAR variables randomly chosen such that $k \neq \ell \in [1:10] \backslash \{m\}$. In Figure 12, Algorithm 1 provides the best estimators of the mean and the variance (in term of bias) and the smallest imputation error.

**Higher dimension and variation of the rank.** The performance of the different methods for higher dimension is assessed. A data matrix of size $n = 1000$ and $p = 50$ is generated from two
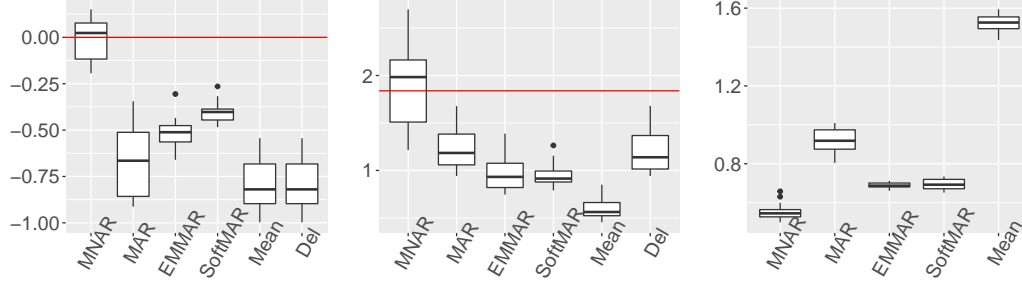
28

Figure 12: Mean estimation (left), variance estimation (middle) of one missing variable and imputation error (right) when $r = 3$, $n = 1000$, $p = 20$ and ten variables are MNAR as in (54). True values are indicated in red lines.
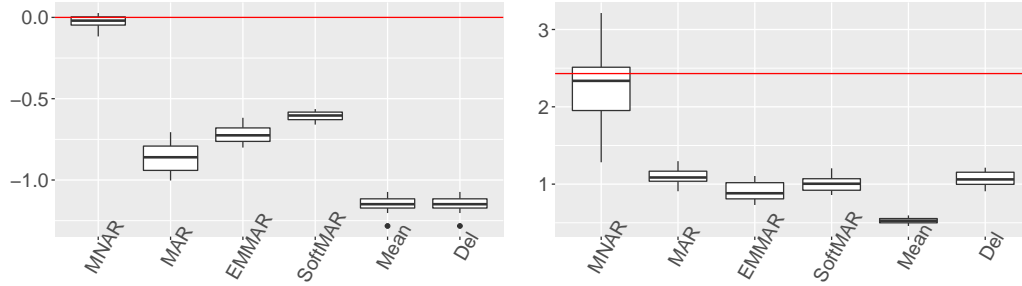


Figure 13: Mean estimation (left) and variance estimation (right) of one missing variable when $r = 2$, $n = 1000$, $p = 50$ and twenty variables are MNAR. True values to be estimated are indicated by red lines.

latent variables ($r = 2$) and with a noise level $\sigma = 1$. Missing values are introduced on twenty variables according to a self-masked MNAR mechanism, leading to 20% of missing values in total. Without loss of generality, the results are presented for one missing variable. Method (iv) has been discarded, as its computational time is too high for this setting.

In Figure 13, as for the estimated mean and variance, Methods (i), (ii) and (iii) suffer from a large bias, while Algorithm 1 gives unbiased estimators. The same comment can be done for the estimation of the covariance between two missing values in Figure 14. As for the covariance estimation between a missing variable and a pivot one Figure 14, Algorithm 1 suffers from a variability, which can be due to the fact that in this higher dimension setting, not all the possible combinations of pivot variables are considered. Indeed, instead of taking the set of pivot variables of all the not MNAR variables *i.e.* $\mathcal{J} = \overline{\mathcal{M}}$, we choose $\mathcal{J} \subset \overline{\mathcal{M}}$ such that $|\mathcal{J}| = 10$. For the mean, 270 combinations of the pivot variables are aggregated over 870 possible combinations if $\mathcal{J} = \overline{\mathcal{M}}$.
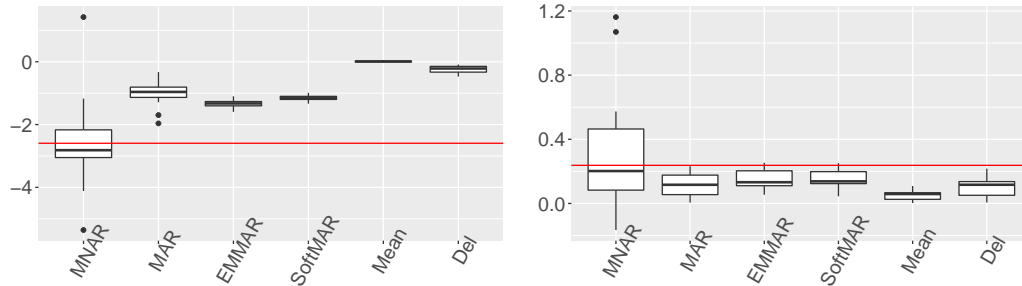


Figure 14: Covariance estimation beetween two missing variable (left) and a missing variable and a pivot one (right) when $r = 2$, $n = 200$, $p = 10$ and seven variables are MNAR. True values to be estimated are indicated by red lines.
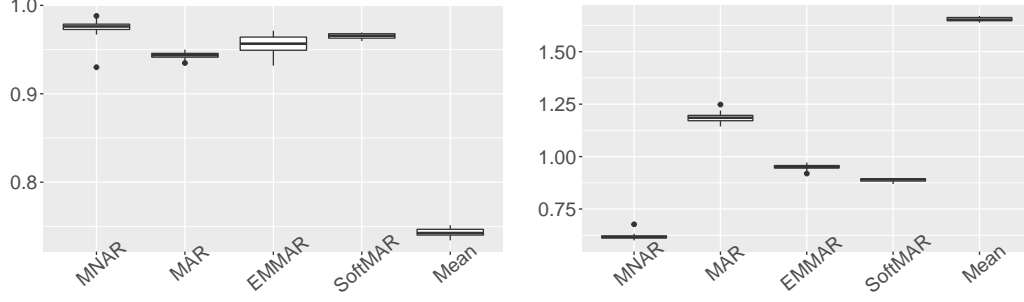
29

Figure 15: RV coefficients for the loading matrix (left) and imputation error (right) when $r = 2$, $n = 1000$, $p = 50$ and twenty variables are MNAR.

Despite this dispersed estimator of the covariance between a MNAR variable and a pivot one, Algorithm 1 gives in Figure 15 a high RV coefficient, by improving Methods (i), (iii) and (ii). Concerning the imputation performance, Algorithm 1 strongly improves Methods (ii) and (iii).

For the same dimension setting ($n = 1000$, $p = 50$) and the same noise level ($\sigma = 1$), we vary the rank to $r = 5$. Similarly as before, missing values are introduced on twenty variables according to a self-masked MNAR mechanism, leading to 20% of missing values in total. In Figure 16, for the mean and the variable estimations, Algorithm 1 gives unbiased estimators. In Figure 17, the covariance between a missing variable and a pivot estimated by Algorithm 1 is biased but still less than the other methods. In addition, the covariance between two missing variables is unbiased but suffers from a high variability. Note that once again we have chosen $\mathcal{J} \subset \mathcal{M}$ such that $|\mathcal{J}| = 10$. For the mean, 1260 combinations of the pivot variables are aggregated over 712530 possible combinations if $\mathcal{J} = \bar{\mathcal{M}}$. In Figure 18, despite such results for the covariance estimators, Algorithm 1 gives a similar RV coefficient than Methods (ii) and (iii) but strongly improves all the methods in term of imputation error.
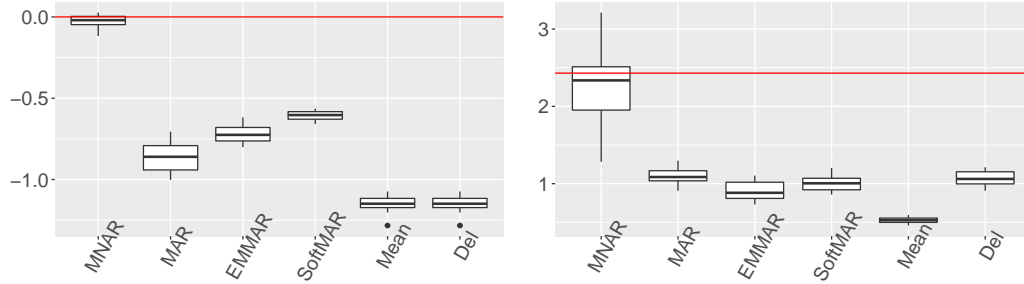


Figure 16: Mean estimation (left) and variance estimation (right) of one missing variable when $r = 5$, $n = 1000$, $p = 50$ and twenty variables are MNAR. True values to be estimated are indicated by red lines.
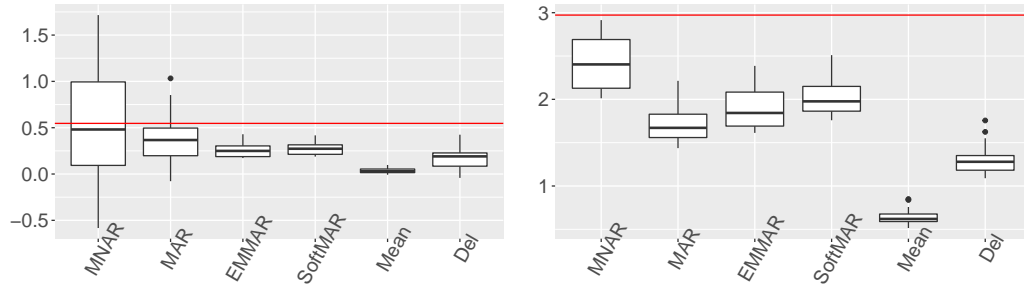


Figure 17: Covariance estimation beetween two missing variable (left) and a missing variable and a pivot one (right) when $r = 5$, $n = 1000$, $p = 50$ and twenty variables are MNAR. True values to be estimated are indicated by red lines.
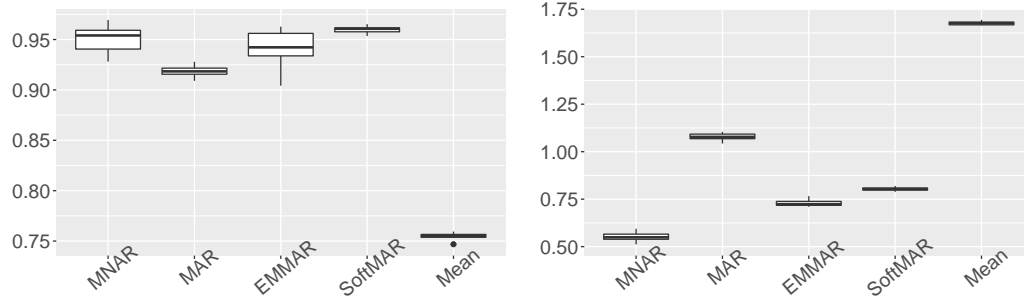
Figure 18: RV coefficients for the loading matrix (left) and imputation error (right) when $r = 5$, $n = 1000$, $p = 50$ and twenty variables are MNAR.

**Efficiency of the *aggregation* approach in the selection of the pivot variables.** As described in Section 3.4, Algorithm 1 requires the selection of $r$ pivot variables (considered M(C)AR) on which the regressions will be performed. To reduce the error committed by the selection pivot variables, we propose to select a bigger set of pivot variables (with a cardinal superior to $r$) and the final estimator will be computed with the median of the estimators over all possible combinations of $r$ pivot variables (this is called the *aggregation* approach). In Figure 19, we consider the same setting as in Section 4.1 ($n = 1000$, $p = 10$, $r = 2$ and seven self-masked MNAR variables) and we perform Algorithm 1 by using the *aggregation* (`MNARagg`) method or not (`MNARnoagg`). By discarding outliers, this *aggregation* approach is more robust than selecting only $r$ pivot variables.
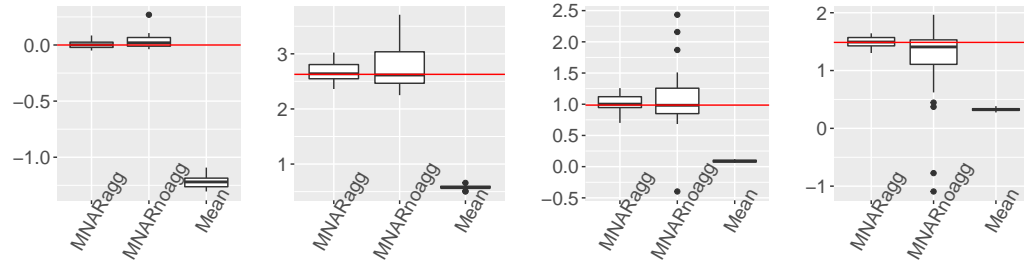


Figure 19: Mean (left) and variance (middle left) estimations of $Y_{.1}$ and covariances estimations of $\mathrm{Cov}(Y_{.1}, Y_{.2})$ (between two missing variables) (middle right) and of $\mathrm{Cov}(Y_{.1}, Y_{.8})$ (between one missing variable and one pivot variable) (right). True values are indicated in red lines.

# D Computation time

Table 1 gathers computation times of the different methods, for both settings considered in Sections 4 and C.

| Method | $r = 2, p = 10, n = 1000$ 35% MNAR values in 7 variables | $r = 5, p = 50, n = 1000$ 20% MNAR values in 20 variables |
|---|---|---|
| MNAR algebraic | 0,1 s | 11 min 48 s (1260 aggregations) |
| SoftMAR | 5,5 s | 28 s |
| EMMAR | 50,8 s | 2 min 9 s |
| Param | 5 h 15 min | not evaluated |

Table 1: Computation time for simulations in Sections 4 and Appendix C. The process time is obtained for a computer with a processor Intel Core i5 of 2,3 GHz.

# E   Additional information on the TraumaBase® dataset

## E.1   Description of the variables

A description of the variables which are used in Section 4.2 is given. The indications given in parentheses ph (pre-hospital) and h (hospital) mean that the measures have been taken before the arrival at the hospital and at the hospital.

- *SBP.ph*, *DBP.ph*, *HR.ph*: systolic and diastolic arterial pressure and heart rate during pre-hospital phase. (ph)

- *HemoCue.init*: prehospital capillary hemoglobin concentration. (ph)

- *SpO2.min*: peripheral oxygen saturation, measured by pulse oxymetry, to estimate oxygen content in the blood. (ph)

- *Cristalloid.volume*: total amount of prehospital administered cristalloid fluid resuscitation (volume expansion). (ph)

- *Shock.index.ph*: ratio of heart rate and systolic arterial pressure during pre-hospital phase. (ph)

- *Delta.shock.index*: Difference of shock index between arrival at the hospital and arrival on the scene. (h)

- *Delta.hemoCue*: Difference of hemoglobin level between arrival at the hospital and arrival on the scene. (h)

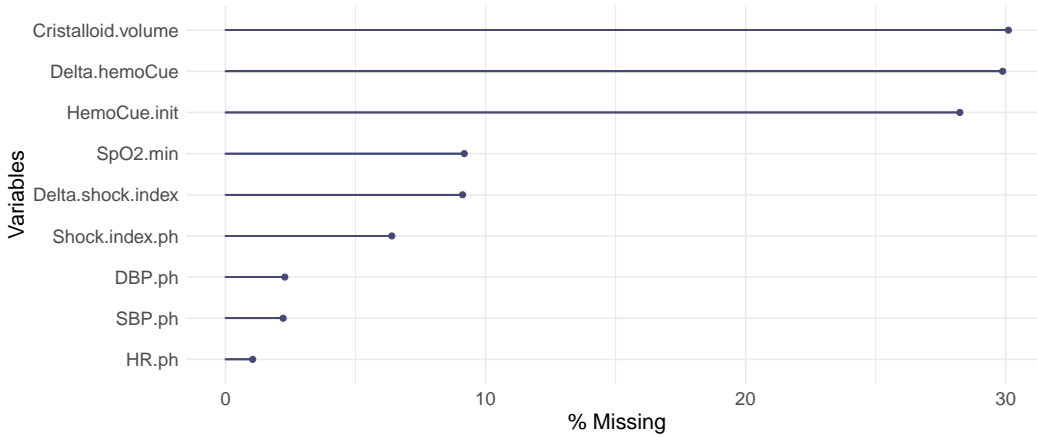The percentage of missing values in each variable is given in Figure 20.



Figure 20:  Percentage of missing values in each variable for the TraumaBase data.

## E.2   Supervised learning task

To predict the administration or not of the tranexomic acid (binary variable), we impute explanatory variables before proceeding to the classification task. In Table 2, Algorithm 1 gives the smallest prediction error.

| | |
|---|---|
| MNAR | 5.06% |
| EMMAR | 5.82% |
| SoftMAR | 5.45% |
| MNARparam | 5.39% |
| Mean | 5.27% |

Table 2: Mean of prediction error over 10 repetitions.

# F    Graphical approach

## F.1    Preliminaries

Lemmas of Mohan et al. [18] are used to construct some estimators of the mean, variance and covariances for a MNAR variable based on a graphical approach.

**Lemma 9** (Lemma 2 [18]). *Let us consider the m-graph $G$. The coefficient of the linear regression of $Y_{.j}$ on $Y_{.k}, k \neq j$, denoted as $\beta_{j \to k, k \neq j}$ is recoverable (i.e. they are consistent in the complete-case analysis) if $Y_{.j} \perp\!\!\!\perp \Omega | Y_{.k}, k \neq j$ and one has*

$$\beta_{j \to k, k \neq j} = \beta^c_{j \to k, k \neq j}.$$

**Lemma 10** (Lemma 1). *[18]](Graphical approach for computing the covariance) Let $G$ be a m-graph with $k$ unblocked paths $p_1, \ldots, p_k$ between two variables $Y_{.\tau}$ and $Y_{.\delta}$. Let $A_{p_i}$ be the ancestor of all nodes on path $p_i$. Let the number of nodes on $p_i$ be $n_{p_i}$. One can derive that*

$$\mathrm{Cov}(Y_{.\tau}, Y_{.\delta}) = \sum_{i=1}^{k} \mathrm{Var}(A_{p_i}) \prod_{j=1}^{n_{p_i}-1} \alpha_j^{p_i},$$

*where $\prod_{j=1}^{n_{p_i}-1} \alpha_j^{p_i}$ is the product of all causal parameters on path $p_i$.*

In addition, let us recall the basic formula,

$$\beta_{Y \to X} = \frac{\mathrm{Cov}(X, Y)}{\mathrm{Var}(X)}, \tag{55}$$

where $Y$ and $X$ are two variables of a linear model.

## F.2    Estimation of the mean, variance and covariances of the MNAR variables

The graphical approach to construct an estimator of $\alpha_1$ is based on the transformation illustrated in Figure 1 of the graphical model of PPCA as structural causal graphs, whose context is introduced in [21]. This latter framework allows to directly apply the results of Mohan et al. [18] who consider the associated (linear) structural causal equations under the exogeneity assumption with MNAR missing values for one variable.

For the sake of brevity, the results are presented for the toy example in Section 3.1 where $p = 3$, $r = 2$, $Y_{.1}$ is self-masked MNAR and the other variables are observed.

Then, one can associate to Figure 1 (bottom right graph) the structural equation model detailled in the following lemma.

**Lemma 11.** *Assuming $\mathbb{E}[\epsilon_{.2}|Y_{.1}, Y_{.3}] = 0$, the structural equation model associated with the bottom right graph in Figure 1 is*

$$Y_{.2} = \beta_{2 \to 1,3[0]} + \beta_{2 \to 1,3[1]} Y_{.1} + \beta_{2 \to 1,3[3]} Y_{.3} + \epsilon_{.2}, \tag{56}$$

*where $\beta_{2 \to 1,3[0]}$, $\beta_{2 \to 1,3[1]}$ and $\beta_{2 \to 1,3[3]}$ are the intercept and the coefficients of the linear regression of $Y_{.2}$ on $Y_{.1}$ and $Y_{.3}$.*

Using Equation (56) and Lemma 9, we apply the results of Mohan et al. [18] to get an estimator for the mean of the MNAR variable.

**Proposition 12** (Mean estimator for the graphical approach). *Under Equation (56), assuming A1. and $\beta^c_{2 \to 1,3[1]} \neq 0$, one can construct an estimator of the mean $\alpha_1$ of the MNAR variable $Y_{.1}$ as follows*

$$\hat{\alpha}_1 := \frac{\hat{\alpha}_2 - \hat{\beta}^c_{2 \to 1,3[0]} - \hat{\beta}^c_{2 \to 1,3[3]} \hat{\alpha}_3}{\hat{\beta}^c_{2 \to 1,3[1]}}, \tag{57}$$

*where $\hat{\beta}^c_{2 \to 1,3[0]}$, $\hat{\beta}^c_{2 \to 1,3[1]}$ and $\hat{\beta}^c_{2 \to 1,3[3]}$ denote some estimators of $\beta^c_{2 \to 1,3[0]}$, $\beta^c_{2 \to 1,3[1]}$ and $\beta^c_{2 \to 1,3[3]}$ given in Lemma 11. This estimator is consistent under additional Assumption A4..*

*Proof.* To derive some estimator of the mean, we want to obtain the following formula

$$\alpha_1 = \frac{\alpha_2 - \beta^c_{2\to1,3[0]} - \beta^c_{2\to1,3[3]}\alpha_3}{\beta^c_{2\to1,3[1]}}.\tag{58}$$

Indeed, one has:

$$\begin{aligned}
\mathbb{E}[Y_{.2}] &= \mathbb{E}[\mathbb{E}[Y_{.2}|Y_{.1}, Y_{.3}]] \\
&= \mathbb{E}[\mathbb{E}[Y_{.2}|Y_{.1}, Y_{.3}, \Omega_{.1}=1]] && \text{(by using \textbf{A1.})} \\
&= \mathbb{E}[\mathbb{E}[\beta^c_{2\to1,3[0]} + \beta^c_{2\to1,3[1]}Y_{.1} + \beta^c_{2\to3,1[3]}Y_{.3} + \epsilon_{.2}|Y_{.1}, Y_{.3}]] \\
&= \beta^c_{2\to1,3[0]} + \beta^c_{2\to1,3[1]}\mathbb{E}[Y_{.1}] + \beta^c_{2\to3,1[3]}\mathbb{E}[Y_{.3}],
\end{aligned}$$

which leads to the desired Equation (58), provided that $\beta^c_{2\to1,3[1]} \neq 0$. A natural estimator fo $\alpha_1$ is then given by (57). It is consistent given that all the quantities involved are consistent, by using **A4.** (for the consistency of $\hat{\alpha}_2$ and $\hat{\alpha}_3$) and Lemma 9 (for the consistency of the coefficients $\hat{\beta}^c_{2\to1,3[0]}$, $\hat{\beta}^c_{2\to1,3[1]}$ and $\hat{\beta}^c_{2\to1,3[3]}$). $\qquad\square$

**Remark 13** (Mean estimation: algebraic vs. graphical approach)**.** *In both approaches, the PPCA model is translated into a linear model. However, both estimators in Equations (9) and (57) theoretically differ. The exogeneity assumption and approximation is not made at the same step. In the algebraic approach, the results are first derived without using any approximation. It gives linear models that do not comply with the standard exogeneity assumption. Consequently, an approximation is done at the estimation step since the parameters $\hat{\mathcal{B}}^c_{2\to1,3[0]}$, $\hat{\mathcal{B}}^c_{2\to1,3[1]}$ and $\hat{\mathcal{B}}^c_{2\to1,3[3]}$ are estimated with the standard linear regression coefficients. In the graphical approach, an approximation is made at the first step when a structural equation model is associated with the graphical model by assuming the exogeneity, i.e. $\mathbb{E}[\epsilon_{.2}|Y_{.1}, Y_{.3}] = 0$. In practice, for both approaches, the same coefficients are naturally computed, i.e. $\hat{\beta}^c_{j\to k,\ell} = \hat{\mathcal{B}}^c_{j\to k,\ell}$, which leads to the same computed estimators for the mean of $Y_{.1}$.*

While only one simplified graphical model between $Y_{.1}$, $Y_{.2}$ and $Y_{.3}$, displayed in the bottom right graph of Figure 1, was required to construct an estimator of the mean of $Y_{.1}$, the variance and covariance estimations rely on Equation (56) and the following one (associating to the bottom left graph of Figure 1),

$$Y_{.3} = \beta_{3\to1,2[0]} + \beta_{3\to1,2[1]}Y_{.1} + \beta_{3\to1,2[2]}Y_{.2} + \epsilon_{.3},\tag{59}$$

assuming $\mathbb{E}[\epsilon_{.3}|Y_{.1}, Y_{.2}] = 0$ and where $\beta_{3\to1,2[0]}$, $\beta_{3\to1,2[1]}$ and $\beta_{3\to1,2[2]}$ are the intercept and the coefficients of the linear regression of $Y_{.3}$ on $Y_{.1}$ and $Y_{.2}$.

Using Equations (56) and (59) and Lemmas 9, 10, one can derive some estimators for the variance and the covariances of $Y_1$.

**Proposition 14** (Variance and covariances formulae resulting from the graphical approach when $p = 3$ and $r = 2$)**.** *Under the two equations (56) and (59), assuming \textbf{A1.} and also $\beta^c_{3\to1} \neq 0$, $\beta^c_{2\to1,3[1]} \neq 0$ and $\mathrm{Var}(Y_{.3}) \neq 0$, one can construct an estimator of the variance of the MNAR variable $Y_{.1}$ and its covariances as follows*

$$\widehat{\mathrm{Var}}(Y_{.1}) := \frac{\widehat{\mathrm{Var}}(Y_{.3})}{\hat{\beta}^c_{3\to1}}\frac{1}{\hat{\beta}^c_{2\to1,3[1]}}\left(\frac{\widehat{\mathrm{Cov}}(Y_{.2}, Y_{.3})}{\widehat{\mathrm{Var}}(Y_{.3})} - \hat{\beta}^c_{2\to1,3[3]}\right),\tag{60}$$

$$\widehat{\mathrm{Cov}}(Y_{.1}, Y_{.2}) := \frac{1}{\hat{\beta}^c_{3\to1,2[1]}}\left(\frac{\widehat{\mathrm{Cov}}(Y_{.2}, Y_{.3})}{\widehat{\mathrm{Var}}(Y_{.2})} - \hat{\beta}^c_{3\to1,2[2]}\right)\widehat{\mathrm{Var}}(Y_{.2}),\tag{61}$$

$$\widehat{\mathrm{Cov}}(Y_{.1}, Y_{.3}) := \frac{1}{\hat{\beta}^c_{2\to1,3[1]}}\left(\frac{\widehat{\mathrm{Cov}}(Y_{.2}, Y_{.3})}{\widehat{\mathrm{Var}}(Y_{.3})} - \hat{\beta}^c_{2\to1,3[3]}\right)\widehat{\mathrm{Var}}(Y_{.3}),\tag{62}$$

*where $\hat{\beta}^c_{3\to1,2[1]}$, $\hat{\beta}^c_{3\to1,2[2]}$ and $\hat{\beta}^c_{3\to1}$ are some estimators of $\beta^c_{3\to1,2[1]}$, $\beta^c_{3\to1,2[2]}$ and $\beta^c_{3\to1}$ given in (59).*

*These estimators are consistent under additional Assumption \textbf{A4.}.*

*Proof.* To derive some estimators of the variance and covariances of the MNAR variable $Y_{.1}$, one want to obtain the following formulae:

$$\mathrm{Var}(Y_{.1}) = \frac{\mathrm{Var}(Y_{.3})}{\beta^c_{3\to1}} \frac{1}{\beta^c_{2\to1,3[1]}} \left( \frac{\mathrm{Cov}(Y_{.2},Y_{.3})}{\mathrm{Var}(Y_{.3})} - \beta^c_{2\to1,3[3]} \right), \tag{63}$$

$$\mathrm{Cov}(Y_{.1},Y_{.2}) = \frac{1}{\beta^c_{3\to1,2[1]}} \left( \frac{\mathrm{Cov}(Y_{.2},Y_{.3})}{\mathrm{Var}(Y_{.2})} - \beta^c_{3\to1,2[2]} \right) \mathrm{Var}(Y_{.2}), \tag{64}$$

$$\mathrm{Cov}(Y_{.1},Y_{.3}) = \frac{1}{\beta^c_{2\to1,3[1]}} \left( \frac{\mathrm{Cov}(Y_{.2},Y_{.3})}{\mathrm{Var}(Y_{.3})} - \beta^c_{2\to1,3[3]} \right) \mathrm{Var}(Y_{.3}). \tag{65}$$

Using Equation (55), one has

$$\mathrm{Cov}(Y_{.1},Y_{.3}) = \mathrm{Var}(Y_{.1})\beta_{3\to1},$$
$$\mathrm{Cov}(Y_{.3},Y_{.1}) = \mathrm{Var}(Y_{.3})\beta_{1\to3},$$

so

$$\mathrm{Var}(Y_{.1}) = \frac{\mathrm{Var}(Y_{.3})\beta_{1\to3}}{\beta_{3\to1}}.$$

Considering the graphical model in the bottom left graph of Figure 1,

$$\mathrm{Cov}(Y_{.2},Y_{.3}) = \beta_{2\to1,3[1]}\beta_{1\to3}\mathrm{Var}(Y_{.3}) + \beta_{2\to1,3[3]}\mathrm{Var}(Y_{.3}) \qquad \text{(by Lemma 10)}$$

$$\Rightarrow \beta_{1\to3} = \frac{1}{\beta_{2\to1,3[1]}} \left( \frac{\mathrm{Cov}(Y_{.2},Y_{.3})}{\mathrm{Var}(Y_{.3})} - \beta_{2\to1,3[3]} \right)$$

$$\Rightarrow \beta_{1\to3} = \frac{1}{\beta^c_{2\to1,3[1]}} \left( \frac{\mathrm{Cov}(Y_{.2},Y_{.3})}{\mathrm{Var}(Y_{.3})} - \beta^c_{2\to1,3[3]} \right) \tag{66}$$

where the last implication is given by Lemma 9 and Assumption **A1.**, giving also

$$\beta_{3\to1} = \beta^c_{3\to1},$$

which leads to Equation (63).

By (55), the covariances can be expressed in two different ways,

$$\mathrm{Cov}(Y_{.1},Y_{.2}) = \beta_{2\to1}\mathrm{Var}(Y_{.1}) \quad \text{and} \quad \mathrm{Cov}(Y_{.1},Y_{.3}) = \beta_{3\to1}\mathrm{Var}(Y_{.1}), \tag{67}$$
$$\mathrm{Cov}(Y_{.1},Y_{.2}) = \beta_{1\to2}\mathrm{Var}(Y_{.2}) \quad \text{and} \quad \mathrm{Cov}(Y_{.1},Y_{.3}) = \beta_{1\to3}\mathrm{Var}(Y_{.3}). \tag{68}$$

In (67), the coefficients $\beta_{2\to1}$ and $\beta_{3\to1}$ can be estimated on the complete case using Lemma 9, but the variance of $Y_{.1}$ has still to be taken care of. Instead of potentially propagate error from (63), we propose to favor the expressions given in (68) to evaluate the covariances.

Focusing on (68), the coefficient $\beta_{1\to3}$ is given in (66) and $\beta_{1\to2}$ can be obtained using the same method, based on the reduced graphical model in the bottom right graph of Figure 1 (by Assumption **A1.**), so that

$$\beta_{1\to2} = \frac{1}{\beta^c_{3\to1,2[1]}} \left( \frac{\mathrm{Cov}(Y_{.2},Y_{.3})}{\mathrm{Var}(Y_{.2})} - \beta^c_{3\to1,2[2]} \right).$$

Therefore, by plugging it in (68), Equations (64) and (65) are obtained.

The natural estimators for $\mathrm{Var}(Y_{.1})$, $\mathrm{Cov}(Y_{.1},Y_{.2})$ and $\mathrm{Cov}(Y_{.1},Y_{.3})$ are then given by (60), (61) and (62). They are consistent given that all the quantites involved are consistent, by using **A4.** (for the consistency of $\widehat{\mathrm{Var}}(Y_{.2})$, $\widehat{\mathrm{Var}}(Y_{.3})$ and $\widehat{\mathrm{Cov}}(Y_{.2},Y_{.3})$) and Lemma 9 (for the consistency of $\hat{\beta}^c_{j\to k,\ell}$). $\qquad\square$

**Remark 15** (Var-covariance estimation: algebraic vs. graphical approach)**.** *As for the mean, the exogeneity assumption is required in the last step of the algebraic approach to estimate coefficients and in the first step of the graphical approach to obtain structural equation models. However, contrary to the estimator suggested for the mean, the estimators in both graphical and algebraic approaches here differ (compare* (10) *with* (60)*,* (61) *and* (62)*). Indeed, the algebraic approach is based on the use of conditionality, while the graphical one relies on graphical results standing for the linear models when exogeneity holds.*

# G  PPCA with MAR data

The following proposition is an adaptation of our method to handle MAR data, called **MAR** in Section 4.1, inspired by [18, Theorems 1, 2, 3]. In this case, the missing variables are assumed to be MAR indexed by $\mathcal{M}$. We assume the following:

**A1$_{MAR}$.** $(B_{.j'})_{j'\in\mathcal{J}}$ is invertible.

**A2$_{MAR}$.** $\forall m \in \mathcal{M}, Y_{.m} \perp\!\!\!\perp \Omega_{.m}|(Y_k)_{k\in\overline{\{m\}}}$

**A3$_{MAR}$.** $\forall m \in \mathcal{M}$, the complete-case coefficients $\mathcal{B}^c_{m\to\mathcal{J}[0]}$ and $\mathcal{B}^c_{m\to\mathcal{J}[k]}, k \in \mathcal{J}$ can be consistently estimated.

**A5$_{MAR}$.** $\forall \ell \in \bar{\mathcal{J}}$, for all set $\mathcal{H} \subset \mathcal{J}_{-j}$ such that $|\mathcal{H}| = r - 1$, $(B_{.\ell} \quad (B_{.j'})_{j'\in\mathcal{H}})$ is invertible,

**A6$_{MAR}$.** $\forall m \in \mathcal{M}, \forall \ell \in \bar{\mathcal{J}}\backslash\mathcal{M}, \forall j \in \mathcal{J}, \; Y_{.m} \perp\!\!\!\perp \Omega_{.\ell}|(Y_{.k})_{k\in\overline{\{m\}}}$.

**A8$_{MAR}$.** $\forall m \in \mathcal{M}, \forall \ell \in \overline{\{m\}}\backslash\mathcal{J}$, for all set $\mathcal{H} \subset \mathcal{J}$ such that $|\mathcal{H}| = r - 1$, the complete-case coefficients $\mathcal{B}^c_{m\to\ell,\mathcal{H}[0]}$ and $\mathcal{B}^c_{m\to\ell,\mathcal{H}[k]}, k \in \{\ell\} \cup \mathcal{H}$ can be consistently estimated.

**Proposition 16** (Expectation, variance and covariances formulae for a MAR variable when $p = 3$ and $r = 2$). *Consider the PPCA model* (1). *Under Assumptions A1$_{MAR}$. and A2$_{MAR}$., one can construct the estimators of the mean, the variance and the covariances with a pivot variable for any MAR variable $Y_{.m}, m \in \mathcal{M}$, as follows*

– *the mean of the missing variable*

$$\hat{\alpha}_m = \hat{\mathcal{B}}^c_{m\to\mathcal{J}[0]} + \sum_{j\in\mathcal{J}} \hat{\mathcal{B}}^c_{m\to\mathcal{J}[j]}\hat{\alpha}_j,$$

*with $\mathcal{J}$ the pivot variables set,*

– *the variance of the missing variable*

$$\widehat{\mathrm{Var}}(Y_{.m}) = \hat{Q}^c_{\mathrm{MAR}} + \sum_{j\in\mathcal{J}} (\hat{\mathcal{B}}^c_{m\to\mathcal{J}[j]})^2\widehat{\mathrm{Var}}(Y_{.j})$$

$$+ 2\sum_{(j<k)\in\mathcal{J}} \hat{\mathcal{B}}^c_{m\to\mathcal{J}[j]}\hat{\mathcal{B}}^c_{m\to\mathcal{J}[k]}\widehat{\mathrm{Cov}}(Y_{.j}, Y_{.k}),$$

*with*

$$\hat{Q}^c_{\mathrm{MAR}} = \left(\widehat{\mathrm{Var}}(Y_{.m})|\Omega_{.m} = 1\right)$$

$$- \left(\widehat{\mathrm{Cov}}((Y_{.j})_{j\in\overline{\{m\}}}, Y_{.m})\widehat{\mathrm{Var}}((Y_{.j})_{j\in\overline{\{m\}}})^{-1}\widehat{\mathrm{Cov}}((Y_{.j})_{j\in\overline{\{m\}}}, Y_{.m})^T|\Omega_{.m} = 1\right).$$

– *the covariances between the missing variable and a pivot variable, for all $\ell \in \mathcal{J}$,*

$$\widehat{\mathrm{Cov}}(Y_{.m}, Y_{.\ell}) = \hat{\mathcal{B}}^c_{m\to\mathcal{J}[0]}\hat{\alpha}_\ell + \hat{\mathcal{B}}^c_{m\to\mathcal{J}[\ell]}(\widehat{\mathrm{Var}}(Y_{.\ell}) + \hat{\alpha}_\ell^2)$$

$$+ \sum_{k\in\mathcal{J}_{-\ell}} \hat{\mathcal{B}}^c_{m\to\mathcal{J}[k]}(\widehat{\mathrm{Cov}}(Y_{.\ell}, Y_{.k}) + \hat{\alpha}_\ell\hat{\alpha}_k) - \hat{\alpha}_m\hat{\alpha}_\ell$$

*Under Assumption A3$_{MAR}$. and A4., these estimators are consistent.*

*In addition, under Assumption A5$_{MAR}$., A6$_{MAR}$. and A7., one can construct the estimator of the covariance between a MAR variable $Y_{.m}$ for $m \in \mathcal{M}$ and any not pivot variable as follows*

– *the covariances between the missing variable and any not pivot variable, for all $\ell \in \overline{\{m\}}\backslash\mathcal{J}$, choose $r - 1$ variable indexes in $\mathcal{J}$ to form the set $\mathcal{H} \cup \mathcal{J}$ such that $|\mathcal{H}| = r - 1$*

$$\widehat{\mathrm{Cov}}(Y_{.m}, Y_{.\ell}) = \mathcal{B}^c_{m\to\ell,\mathcal{H}[0]}\hat{\alpha}_\ell + \hat{\mathcal{B}}^c_{m\to\ell,\mathcal{H}[\ell]}(\widehat{\mathrm{Var}}(Y_{.\ell}) + \hat{\alpha}_\ell^2)$$

$$+ \sum_{k\in\mathcal{H}} \hat{\mathcal{B}}^c_{m\to\ell,\mathcal{H}[k]}(\widehat{\mathrm{Cov}}(Y_{.\ell}, Y_{.k}) + \hat{\alpha}_\ell\hat{\alpha}_k) - \hat{\alpha}_m\hat{\alpha}_\ell$$

*Under the additional Assumptions $\mathbf{A8_{MAR}}$. and $\mathbf{A9}$. this estimator is consistent.*

*Proof.* The proof follows exactly the same direction than in Proposition 4, 5 and 8. The only difference is that the regressions used are not the same.

For the sake of clarity, consider the same toy example as in Section 3.1 where $p = 3, r = 2$, in which only one variable can be missing (at random), and fix $\mathcal{M} = \{1\}$ and $\mathcal{J} = \{2, 3\}$. Note that here the MAR mechanism leads to $\mathbb{P}(\Omega_{.1} = 0 | Y_{.1}, Y_{.2}, Y_{.3}) = \mathbb{P}(\Omega_{.1} = 0 | Y_{.2}, Y_{.3})$.. The goal is to estimate the mean of $Y_{.1}$, without specifying the distribution of the missing-data mechanism and using only the observed data.

Assumption $\mathbf{A1_{MAR}}$. allows to obtain linear link between the MAR variable $Y_{.1}$ and the pivot variables $(Y_{.2}, Y_{.3})$. In particular, one has

$$Y_{.1} = \beta_{1\to2,3[0]} + \beta_{1\to2,3[2]}Y_{.2} + \beta_{1\to2,3[3]}Y_{.3} + \zeta,$$

with $\beta_{1\to2,3[0]}, \beta_{1\to2,3[2]}$ and $\beta_{1\to2,3[3]}$ the intercept and coefficients standing for the effects of $Y_{.1}$ on $Y_{.2}$ and $Y_{.3}$, and with

$$\zeta = -\mathcal{B}_{1\to2,3[2]}\epsilon_{.2} - \mathcal{B}_{1\to2,3[3]}\epsilon_{.3} + \epsilon_{.1}$$

Assumption $\mathbf{A2_{MAR}}$., *i.e.* $Y_{.1} \perp\!\!\!\perp \Omega_{.1} | Y_{.2}, Y_{.3}$, is required to obtain identifiable and consistent parameters of the distribution of $Y_{.1}$ given $Y_{.2}, Y_{.3}$ in the complete-case when $\Omega_{.1} = 1$, denoted as $\beta^c_{1\to2,3[0]}$, $\beta^c_{1\to2,3[2]}$ and $\beta^c_{1\to2,3[3]}$,

$$(Y_{.1})_{|\Omega_{.1}=1} = \beta^c_{1\to2,3[0]} + \beta^c_{1\to2,3[2]}Y_{.2} + \beta^c_{1\to2,3[3]}Y_{.3} + \zeta^c,$$

with

$$\zeta^c = -\mathcal{B}^c_{1\to2,3[2]}\epsilon_{.2} - -\mathcal{B}^c_{1\to2,3[3]}\epsilon_{.3} + \epsilon_{.1}$$

(In the MNAR case, the regression of $Y_{.1}$ on $(Y_{.2}, Y_{.3})$ is prohibited, as $\mathbf{A2_{MAR}}$. does not hold. That is why we used the regression of $Y_{.2}$ on $Y_{.1}$ and $Y_{.3}$.);

Using again $\mathbf{A2_{MAR}}$., one has

$$\mathbb{E}\left[Y_{.1}|Y_{.2}, Y_{.3}, \Omega_{.1} = 1\right] = \mathbb{E}\left[\beta^c_{1\to2,3[0]} + \beta^c_{1\to2,3[2]}Y_{.2} + \beta^c_{1\to2,3[3]}Y_{.3}|Y_{.2}, Y_{.3}\right] + \mathbb{E}[\zeta^c|Y_{.2}, Y_{.3}],$$

and taking the expectation leads to

$$\mathbb{E}\left[Y_{.1}\right] = \beta^c_{1\to2,3[0]} + \beta^c_{1\to2,3[2]}\mathbb{E}\left[Y_{.2}\right] + \beta^c_{1\to2,3[3]}\mathbb{E}\left[Y_{.3}\right],$$

given that $\mathbb{E}[\epsilon_{.k}] = 0, \ \forall k \in \{1, 2, 3\}$.

One obtains

$$\alpha_1 = \beta^c_{1\to2,3[0]} + \beta^c_{1\to2,3[2]}\alpha_2 + \beta^c_{1\to2,3[3]}\alpha_3$$

A natural estimator for $\alpha_1$ is

$$\hat{\alpha}_1 = \hat{\beta}^c_{1\to2,3[0]} + \hat{\beta}^c_{1\to2,3[2]}\hat{\alpha}_2 + \hat{\beta}^c_{1\to2,3[3]}\hat{\alpha}_3,$$

which is consistent using Assumption $\mathbf{A3_{MAR}}$. and $\mathbf{A4}$.. $\qquad\square$