



Accompanying note: Model-based Clustering with Missing Not At Random Data

Aude Sportisse, Matthieu Marbac, Fabien Laporte, Gilles Celeux, Claire Boyer, Christophe Biernacki, Julie Josse

► To cite this version:

Aude Sportisse, Matthieu Marbac, Fabien Laporte, Gilles Celeux, Claire Boyer, et al.. Accompanying note: Model-based Clustering with Missing Not At Random Data. 2023. hal-04358192

HAL Id: hal-04358192

<https://hal.science/hal-04358192>

Preprint submitted on 21 Dec 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Accompanying note: Model-based Clustering with Missing Not At Random Data

Aude Sportisse, Matthieu Marbac,
Fabien Laporte, Gilles Celeux
Claire Boyer, Christophe Biernacki, Julie Josse

December 21, 2023

Contents

| | | |
|----------|--|-----------|
| 1 | Organization of this accompanying note | 2 |
| 1.1 | Contributions | 2 |
| 1.2 | Notations | 2 |
| 2 | Variant MNAR modeling | 3 |
| 2.1 | Rationale for MNAR assumptions | 3 |
| 2.2 | Models | 4 |
| 3 | Identifiability of the model parameters | 5 |
| 4 | Estimation of the proposed MNAR models | 7 |
| 5 | Implementation and numerical experiments | 9 |
| A | Appendix 1: Identifiability | 12 |
| A.1 | Continuous and count data | 12 |
| A.2 | On identifiability of the Gaussian mixture | 13 |
| A.3 | On identifiability of the Poisson mixture | 15 |
| A.4 | Categorical data | 16 |
| B | Appendix 2: Details on EM algorithm | 17 |
| C | Appendix 3: Details on SEM algorithm | 20 |
| D | Appendix 4: Complements on the numerical experiment | 26 |

1 Organization of this accompanying note

This document is the accompanying note of the main paper (Sportisse et al., 2023). We assume the data missing not at random (MNAR) values (Rubin, 1976; Ibrahim et al., 2001; Mohan et al., 2018), *i.e.* the effect of missingness depends on the missing values themselves. An example includes clinical data collected in emergency situations, where doctors may choose to treat patients before measuring heart rate: the missingness of heart rate depends on the missing heart rate itself. For such a setting, the observed data are therefore not representative of the population. The main paper focuses on the specific MNAR_z setting, for which the only effect of missingness is on the class membership; in this document, we give some details for other MNAR settings.

1.1 Contributions

In Section 2, we present and illustrate a relevant inventory of distributions for the MNAR missingness process in the context of unsupervised classification based on mixture models for different types of data (continuous, count, categorical and mixed). In Section 3, we provide the identifiability of the mixture model parameters and missingness process parameters under certain conditions (including the data type and the link functions governing the missingness mechanism distribution). This is a real issue in the context of MNAR data, as models often lead to unidentifiable parameters. When all variables are continuous or count, all models lead to identifiable parameters. In the categorical and mixed cases, only the models for which missingness depends uniquely on the class membership have identifiable parameters. These identifiability results represent a substantial extension of the work of Miao et al. (2016) to more complex missing scenario and to the multivariate case. For each model or submodel, an EM or Stochastic EM algorithm is proposed in Section 4, implemented, and made available for reproducibility¹. Finally, we propose a numerical experiment in Section 5 to assess the clustering performance in each proposed MNAR setting.

1.2 Notations

Set the dataset $Y = (\mathbf{y}_1 | \dots | \mathbf{y}_n)^T$ consisting of n individuals, where each observation $\mathbf{y}_i = (y_{i1}, \dots, y_{id})^T$ belongs to a space \mathcal{Y} , depending on the type of data, defined by d features. The pattern of missing data is denoted by $C = (\mathbf{c}_1 | \dots | \mathbf{c}_n)^T \in \{0, 1\}^{n \times d}$, with $\mathbf{c}_i = (\mathbf{c}_{i1}, \dots, \mathbf{c}_{id})^T \in \{0, 1\}^d$: $\mathbf{c}_{ij} = 1$ indicates that the value y_{ij} is missing and $\mathbf{c}_{ij} = 0$ otherwise. The values of the observed (resp. missing) variables for individual i are denoted by $\mathbf{y}_i^{\text{obs}}$ (resp. $\mathbf{y}_i^{\text{mis}}$). The objective of clustering is to estimate an unknown partition $Z = (\mathbf{z}_1 | \dots | \mathbf{z}_n)^T \in \{0, 1\}^{n \times K}$ that groups the full dataset Y into K classes, with $\mathbf{z}_i = (z_{i1}, \dots, z_{iK})^T \in \{0, 1\}^K$

¹The code is available on <https://anonymous.4open.science/r/Clustering-MNAR-7E29>.

and where $z_{ik} = 1$ if \mathbf{y}_i belongs to cluster k , $z_{ik} = 0$ otherwise. Consequently, in a clustering context, the missing data are not only the values $\mathbf{y}_i^{\text{mis}}$ but also the partition labels \mathbf{z}_i .

Mixture models allow for clustering by modeling the distribution of the observed data $(\mathbf{y}_i^{\text{obs}}, \mathbf{c}_i)$. Assuming an underlying mixture model with K components, the probability distribution function (pdf) of the couple $(\mathbf{y}_i, \mathbf{c}_i)$ reads as

$$f(\mathbf{y}_i, \mathbf{c}_i; \theta) = \sum_{k=1}^K \pi_k f_k(\mathbf{y}_i; \lambda_k) f_k(\mathbf{c}_i | \mathbf{y}_i; \psi_k), \quad (1)$$

where $\theta = (\gamma, \psi)$ gathers all the model parameters, $\gamma = (\pi, \lambda)$ groups the parameters related to the marginal distribution of \mathbf{y}_i , $\pi = (\pi_1, \dots, \pi_K)$ is the vector of proportions with $\sum_{k=1}^K \pi_k = 1$ and $\pi_k > 0$ for all $k \in \{1, \dots, K\}$. Given $\lambda = (\lambda_1, \dots, \lambda_K)$, $f_k(\cdot; \lambda_k)$ is the pdf of the k -th component parameterized by λ_k , $\psi = (\psi_1, \dots, \psi_K)$ groups the parameters of the missingness mechanisms and $f_k(\mathbf{c}_i | \mathbf{y}_i; \psi_k)$ is the pdf related to the missingness mechanism under component k (i.e., $f_k(\mathbf{c}_i | \mathbf{y}_i; \psi_k) = f(\mathbf{c}_i | \mathbf{y}_i, z_{ik} = 1; \psi_k)$). In many cases, the parameter ψ is interpreted as a nuisance parameter. However, when the mechanism is not ignorable, we need to consider the whole parameter θ to achieve clustering since the pdf of the observed data is

$$f(\mathbf{y}_i^{\text{obs}}, \mathbf{c}_i; \theta) = \int f(\mathbf{y}_i, \mathbf{c}_i; \theta) d\mathbf{y}_i^{\text{miss}}. \quad (2)$$

Different types of pdf $f_k(\cdot; \lambda_k)$ can be considered, depending on the types of features at hand. Thus, if \mathbf{y}_i is a vector of continuous variables, the pdf of a d -variate Gaussian distribution (McLachlan and Basford, 1988; Banfield and Raftery, 1993) can be considered for $f_k(\mathbf{y}_i; \lambda_k)$ and thus λ_k groups the mean vector and the covariance matrix. Moreover, if some components of \mathbf{y}_i are discrete or categorical, the latent class model (see Geweke et al. (1994); McParland and Gormley (2016)) defining $f_k(\mathbf{y}_i; \lambda_k) = \prod_{j=1}^d f_{kj}(y_{ij}; \lambda_{kj})$ can be used, with $\lambda_k = (\lambda_{k1}, \dots, \lambda_{kd})$. In such case, f_{kj} could be the pdf of a Poisson (resp. multinomial) distribution with parameter λ_{kj} if y_{ij} is an integer (resp. categorical) variable. The next subsection discusses the choice of the modeling for the missingness mechanism (i.e., the distribution $f_k(\mathbf{c}_i | \mathbf{y}_i; \psi_k)$).

2 Variant MNAR modeling

2.1 Rationale for MNAR assumptions

To handle MNAR data in selection models, the distribution of the missing-data pattern given the data and the partition should be specified. We consider the following assumptions:

1. The elements of \mathbf{c}_i are conditionally independent given $(\mathbf{y}_i, \mathbf{z}_i)$.

2. The element c_{ij} is conditionally independent given $(\mathbf{y}_i, \mathbf{z}_i)$ from $y_{ij'}$ for $j \neq j'$.

By the categorical nature of the mask \mathbf{c}_i , the independence assumption 1. is a quite natural hypothesis in the context of clustering (Du Roy De Chaumaray and Marbac, 2020; Chi et al., 2016). The independence assumption 2. of the variables amounts to considering self-masked class-wise MNAR mechanisms for each variable: the missingness of the variable j may depend on its value itself (self-masked) and on the class membership (class-wise). Note that the self-masked feature, apart from limiting the number of parameters to be estimated, is now commonly met in the literature (Mohan, 2018; Sportisse et al., 2020; Le Morvan et al., 2020), and is able to retrieve some existing ad hoc MNAR procedures already used in machine learning community (see more details in (Sportisse et al., 2023, Theorem 1)).

More specifically, the conditional distribution of c_{ij} given $(\mathbf{y}_i, \mathbf{z}_i)$ is assumed to be a (classical) generalized linear model with link function ρ , so that finally

$$f_k(\mathbf{c}_i | \mathbf{y}_i; \psi_k) = \prod_{j=1}^d (\rho(\alpha_{kj} + \beta_{kj}y_{ij}))^{c_{ij}} (1 - \rho(\alpha_{kj} + \beta_{kj}y_{ij}))^{1-c_{ij}}, \quad (3)$$

where $\psi_k = (\alpha_{k1}, \beta_{k1}, \dots, \alpha_{kK}, \beta_{kK})$.

The parameter α_{kj} represents a mean effect of missingness on the k -th class membership for the variable j (note that within a same class k , α_{kj} is not necessarily equal to $\alpha_{kj'}$ for $j \neq j'$). The parameter β_{kj} represents the direct effect of missingness on the variable j which depends on the class k as well. This model is called $\text{MNAR}y^k z^j$ in the following.

Some variations include $f_k(\mathbf{c}_{ij} = 1 | \mathbf{y}_i; \psi_k) = \rho(\alpha_{kj} + \gamma_{kj'}y_{ij'})$, where the j' -th variable is always observed. In such a case, the missing values are missing at random (MAR) as only depending on observed variables. Identifiability guarantees and estimation are still valid in such a setting.

2.2 Models

Simpler models can be derived from (3) by imposing equal parameters either across the class membership, or across the variables likely to be missing. First, we introduce three models, with a lower complexity than (3), that still allow the probability of being missing to depend on both the variable itself and the class membership. For the so-called $\text{MNAR}y z^j$ model, the effect of missingness on a variable is the same regardless of the class (while keeping different mean effects α_{kj} on the class membership), so that

$$\text{MNAR}y z^j: \quad \beta_{1j} = \dots = \beta_{Kj}, \quad \forall j. \quad (4)$$

For the $\text{MNAR}y^k z$ model, the missingness has a same mean effect on class membership shared by all variables (while allowing different self-masked and class-wise parameters β_{kj}):

$$\text{MNAR}y^k z: \quad \alpha_{k1} = \dots = \alpha_{kd}, \quad \forall k. \quad (5)$$

The effects on a particular variable and on the class membership can be respectively the same for all the classes and for all the variables, entailing the so-called MNAR_{yz} model:

$$\text{MNAR}_{yz}: \beta_{1j} = \dots = \beta_{Kj}, \forall j \quad \text{and} \quad \alpha_{k1} = \dots = \alpha_{kd}, \forall k. \quad (6)$$

Secondly, the probability to be missing can also depend only on the variable itself. This is actually a particular case of MNAR mechanisms, widely used in practice (Mohan, 2018), that we call MNAR_y here. The only effect of missingness is thus on the variable j , being the same regardless of the class membership,

$$\text{MNAR}_y: \alpha_{11} = \dots = \alpha_{1d} = \alpha_{21} = \dots = \alpha_{Kd} \quad \text{and} \quad \beta_{1j} = \dots = \beta_{Kj} \quad \forall j. \quad (7)$$

A slightly more general case can be considered by allowing the effect of missingness on the variable j to depend on the class k , as in the following MNAR_y ^{k} model,

$$\text{MNAR}_y^k: \alpha_{11} = \dots = \alpha_{1d} = \alpha_{21} = \dots = \alpha_{Kd}. \quad (8)$$

Thirdly, the probability to be missing can also depend only on the class membership, so that the missingness is class-wise only. In the MNAR_z model, we consider that the only effect of missingness is on the class membership k , being the same for all variables,

$$\text{MNAR}_z: \beta_{kj} = 0, \forall (k, j) \quad \text{and} \quad \alpha_{k1} = \dots = \alpha_{kd}, \forall k. \quad (9)$$

The MNAR_z ^{j} model is a slightly more general case than the MNAR_z model, because the effect of missingness on the class membership k is not the same for all the variables,

$$\text{MNAR}_z^j: \beta_{kj} = 0, \forall (k, j). \quad (10)$$

Finally, the simplest model is the missing completely at random (MCAR) one, characterized by no dependence on variables, neither on class membership, *i.e.*, each variable has the same probability of missing,

$$\text{MCAR}: \beta_{kj} = 0, \forall (k, j) \quad \text{and} \quad \alpha_{1j} = \dots = \alpha_{Kj}, \forall j. \quad (11)$$

3 Identifiability of the model parameters

The generic identifiability (Allman et al., 2009) of parameters for continuous, count, categorical, and mixed data (*i.e.*, when the set of unidentifiable parameters has a zero Lebesgue measure) is ensured by the following theorem. We consider the following assumptions:

- A1.** The parameters (π, λ) of the marginal mixture defined by the density $\sum_{k=1}^K \pi_k f_k(\mathbf{y}_i; \lambda_k)$ are identifiable;

A2. There exists a total ordering \preceq of $\mathcal{F}_j \times \mathcal{R}$, for $j \in \{1, \dots, d\}$ fixed, where \mathcal{F}_j is the family of the data densities $\{f_{1j}, \dots, f_{Kj}\}$ and \mathcal{R} is the family of the mechanism densities $\{\rho_1, \dots, \rho_K\} = \{\rho(\cdot; \psi_1), \dots, \rho(\cdot; \psi_K)\}$, where ρ is the cumulative distribution function of any continuous distribution function and $(\psi_k)_{k \in \{1, \dots, K\}}$ its parameter. The total ordering is such that $\forall k < \ell \in \{1, \dots, K\}, \forall j \in \{1, \dots, d\}, F_{kj} \preceq F_{\ell j}$ (denoting $F_{kj} = \rho_k f_{kj}$ and $F_{\ell j} = \rho_\ell f_{\ell j}$) implies $\lim_{u \rightarrow +\infty} \frac{\rho_\ell(u) f_{\ell j}(u)}{\rho_k(u) f_{kj}(u)} = 0$;

A3. The missing-data distribution ρ is assumed to be strictly monotone.

A4. The feature are independently drawn conditionally to the group membership, *i.e.*,

$$f_k(\cdot; \lambda_k) = \prod_{j=1}^d f_{kj}(\cdot; \lambda_{kj}); \quad (12)$$

A5. The dimension d of the observations is related to the number K of clusters so that

$$d \geq 2 \lceil \log_2 K \rceil + 1,$$

with $\lceil x \rceil$ the least integer greater than or equal to x .

Assumption **A1.** means that the identifiability of the parameters (π, λ, ψ) of the model (2) requires the identifiability of the parameters (π, λ) of the marginal mixture of (Y, Z) (*i.e.*, considering the case without missing values). Some authors have already studied the identifiability of the mixture models, when no missing values in Y occur, especially [Teicher \(1963\)](#) for Gaussian mixtures (continuous variables) and [Yakowitz and Spragins \(1968\)](#) for Poisson mixtures (count variables). Assumption **A2.** is the core ingredient to prove the identifiability of the parameters, requiring that a total ordering of the mixture densities holds. We illustrate it by considering concrete examples in Appendix A. Note that under Assumption **A3.** requires that the link function of the missing data mechanism is strictly monotone, but no assumption about its form (*e.g.* logit, probit) is made. Assumption **A4.** requires the conditional independence of the features given the group membership and Assumption **A5.** links the dimension of the observations and the logarithm of the number of clusters. Both assumptions **A4.** and **A5.** are classical in the categorical case, even without missing values ([Allman et al., 2009](#)).

Theorem 3.1. *Define the conditions:*

- C1 The variables correspond to continuous or count data, **A1.** and **A2.** hold true,*
- C2 All the variables are categorical, **A4.** and **A5.** hold true and the mechanism is stated by (9), (10) or (11),*
- C3 At least one variable is continuous or count data and has a marginal distribution that satisfies **A1.** and **A2.**, **A4.** holds true,*

C4 At least one variable is categorical and its associated mechanism is stated by (9), (10) or (11), A4. and A5. hold true.

Assume that Assumption A3. holds and that at least one of conditions C1-C4 is satisfied, then the parameters of the model in (2) are generically identifiable, up to label swapping.

The proof is given in Appendix A. In the case of continuous and count variables, the proof follows the reasoning used by Teicher (1963, Theorem 2) which proves the identifiability of univariate finite mixtures. For categorical variables, the generic identifiability holds only for the MCAR, MNAR z and MNAR z^j mechanisms. The idea of the proof is to rewrite the observed likelihood as the finite mixture of K multinomial distributions, for which the identifiability is given by Corollary 5 of Allman et al. (2009). For MNAR y mechanisms, the rewriting is impossible, because of the dependency on y of the mechanism. The identifiability of mixed data directly follows from the identifiability of continuous and categorical components.

4 Estimation of the proposed MNAR models

Assuming identifiability, we estimate parameters via likelihood maximization using EM and SEM algorithms specifically designed for Gaussian, Poisson, multinomial and mixed data with MNAR data.

For the MNAR z mechanism, an EM algorithm can be derived (see the main paper (Sportisse et al., 2023)), because the effect of the missingness does not depend on y_i . As the MNAR z^j mechanism has the same property, an EM algorithm can also be derived (see Appendix B). However, the EM algorithm becomes untractable when the missingness depends on variables y , such models being generically denoted by MNAR y in the sequel. In particular, some distributions entail untractable integrals at the E-step (e.g., Gaussian components with MNAR y mechanism defined with logit link, see Appendix B for more details).

The stochastic EM algorithm (Celeux and Diebolt, 1985) can overpass the EM's intractability, by imputing missing values using Gibbs sampling instead of integrating over them. In addition, it has another advantage, unlike the EM algorithm, not to be necessarily trapped by the first encountered local maximum of the likelihood function in play (Celeux and Diebolt, 1985). The principle of the SEM algorithm is to involve a stochastic-E step (SE-step) instead of the traditional E-step of the EM algorithm. Details of the algorithm is given in Appendix C. The iteration $[r]$ then becomes:

SE-step: Draw the missing data $(\mathbf{z}_i^{[r]}, \mathbf{y}_i^{\text{mis}[r]})$ according to their conditional distribution given the observed data $(\mathbf{y}_i^{\text{obs}}, \mathbf{c}_i)$ and the current parameter $\theta^{[r-1]}$. As simulating according to this conditional distribution may be difficult, we simulate instead according to the following two conditional probabilities using a Gibbs sam-

pler, by noting $\mathbf{y}_i^{[r]} = (\mathbf{y}_i^{\text{obs}}, \mathbf{y}_i^{\text{mis}[r]})$,

$$\mathbf{z}_i^{[r]} \sim \mathbf{z}_i \mid \mathbf{y}_i^{[r-1]}, \mathbf{c}_i; \theta^{[r-1]} \quad \text{and} \quad \mathbf{y}_i^{\text{mis}[r]} \sim \mathbf{y}_i^{\text{mis}} \mid \mathbf{y}_i^{\text{obs}}, \mathbf{z}_i^{[r]}, \mathbf{c}_i; \theta^{[r-1]}. \quad (13)$$

M-step: Let $Y^{[r]} = (\mathbf{y}_1^{[r]} \mid \dots \mid \mathbf{y}_n^{[r]})$ be the imputed matrix and let $Z^{[r]} = (\mathbf{z}_1^{[r]} \mid \dots \mid \mathbf{z}_n^{[r]})$ be the current and corresponding partition. The parameter $\theta^{[r]}$ is computed using the maximum likelihood estimate in the complete case. For all $k \in \{1, \dots, K\}$, the parameter $\pi_k^{[r]}$ is the proportion of rows of $Y^{[r]}$ belonging to class k . The parameter $\lambda_k^{[r]}$ is updated in a standard way, depending on the parametric mixture family in play. Finally, the parameter $\psi_k^{[r]}$ is the resulting coefficients of a GLM with a binomial link function, cf Appendix C for details.

In the SE-step, note that the sampling of $\mathbf{z}_i^{[r]}$ is performed by a multinomial distribution. The conditional distribution of $\mathbf{y}_i^{\text{mis}}$ given $(\mathbf{y}_i^{\text{obs}}, \mathbf{z}_{ik}^{[r]} = 1, \mathbf{c}_i)$ parameterized by $\theta^{[r-1]}$ is

$$\begin{aligned} & f_k(\mathbf{y}_i^{\text{mis}} \mid \mathbf{y}_i^{\text{obs}}, \mathbf{c}_i; \theta^{[r-1]}) \\ &= \frac{f_k(\mathbf{y}_i^{\text{mis}} \mid \mathbf{y}_i^{\text{obs}}; \theta^{[r-1]}) f_k(\mathbf{c}_i \mid \mathbf{y}_i^{\text{mis}}, \mathbf{y}_i^{\text{obs}}; \psi^{[r-1]})}{\int f_k(\mathbf{y}_i^{\text{mis}} \mid \mathbf{y}_i^{\text{obs}}; \theta^{[r-1]}) f_k(\mathbf{c}_i \mid \mathbf{y}_i^{\text{mis}}, \mathbf{y}_i^{\text{obs}}; \psi^{[r-1]}) d\mathbf{y}_i^{\text{mis}}}. \end{aligned}$$

This distribution may not be classical in general. For example, for MNAR y^* models, it is not explicit if the components are Gaussian and if the missing data distribution ρ is logistic (since the product of logistic and Gaussian distributions is not a standard law). Therefore, the SEM algorithm cannot be easily applied. However, if ρ is the probit function, we can make the distribution of interest explicit (it is a truncated Gaussian distribution when the variables are Gaussian). For MNAR z and MNAR z^j models, all the computations remain feasible. Table 1 summarizes the cases for which the EM or SEM algorithm is feasible.

| | EM | | | SEM | | |
|------------|----------------|---------------------------|------------------|----------|------------------------------------|------------------|
| | Gaussian | | Categorical | Gaussian | | Categorical |
| MNAR z | ✓ | | ✓ | ✓ | | ✓ |
| MNAR z^j | | | | | | |
| | Probit | Logit | | Probit | Logit | |
| MNAR y^* | no closed form | no closed form, optim. pb | not identifiable | ✓ | require algorithms as SIR (costly) | not identifiable |

Table 1: Summary of the cases for which the EM and the SEM lead to feasible (or not feasible) computations. The symbol ✓ means that the computations are feasible (derived in Appendix B).

5 Implementation and numerical experiments

The SEM algorithm has been implemented for each MNAR setting for Gaussian data and is available on <https://anonymous.4open.science/r/Clustering-MNAR-7E29>.

For each MNAR setting, we assess the clustering performance through the consistency of the partition, by computing the ARI between the true partition Z and the estimated one, given by $\hat{Z}^{\text{MAP}} = \{z_{ik}^{\text{MAP}}(\hat{\theta})\}_{i,k} \in \mathbb{R}^{n \times K}$ as follows

$$\text{with } z_{ik}^{\text{MAP}}(\theta) = \underset{k \in \{1, \dots, K\}}{\operatorname{argmax}} \mathbb{P}(z_{ik} = 1 | \mathbf{y}_i^{\text{obs}}, \mathbf{c}_i; \theta).$$

We consider thus the following methods:

- the EM algorithm (for MCAR (11), MNAR z (9) and MNAR z^j (10)),
- the SEM algorithm (for MNAR y (7), MNAR y^k (8), MNAR $y^k z^j$ (3), MNAR yz (6))

The data are generated using a Gaussian mixture with three components having unequal proportions ($\pi_1 = 0.5$, $\pi_2 = \pi_3 = 0.25$) and independent variables:

$$\forall j \in \{1, \dots, d\}, y_{ij} = \delta \sum_{k=1}^3 \varphi_{kj} z_{ik} + \epsilon_{ij}, \quad (14)$$

with $\epsilon_{ij} \sim \mathcal{N}(0, 1)$ the noise term, $\varphi_k \in \{0, 1\}^d$ and $\delta > 0$. We consider $d = 6$ variables and we vary the number of observations $n = 100, 250, 500$. In Figure 1, as expected, considering the mechanism always gives better results than using the MCAR model, especially for models with many parameters and larger sample sizes (as the MNAR yz , MNAR $y^k z^j$, MNAR $y^k z$, MNAR yz^j settings for $n = 250$ and $n = 500$). Finally, consistency seems satisfactory in each scenario, indicating that our tuning parameters for the algorithm (starting values, stopping rules) are quite suitable.

In the main paper (Sportisse et al., 2023), we compare the MNAR z setting with other ones, and discuss the computational cost of the estimation for the different MNAR settings.

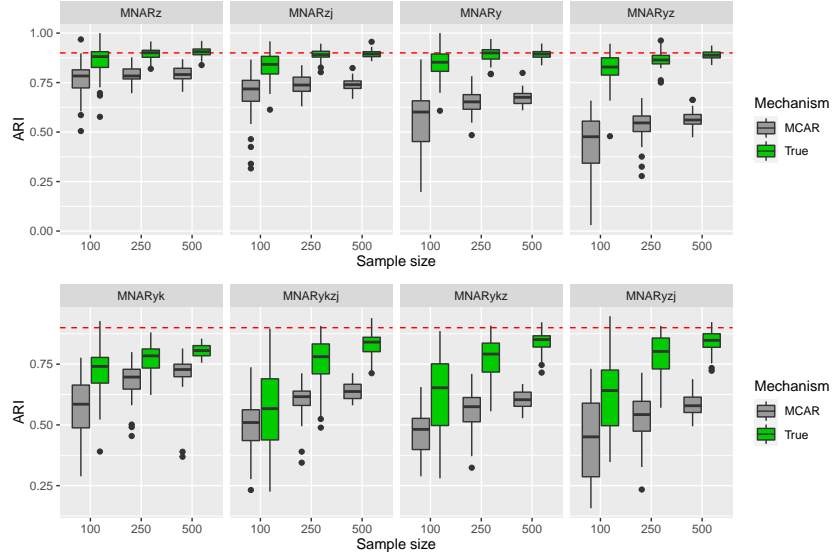


Figure 1: Boxplot of the ARI obtained for 50 samples composed of $d = 6$ variables. The sample size varies by $\{100, 250, 500\}$. The boxplots in green (True) correspond to the performance of the algorithm considering the MNAR setting matching the one that has been used for the missing value generation. The red dashed line indicates the theoretical ARI. This experiment been performed for a theoretical rate of misclassification of 10% and a theoretical missing rate in the whole dataset of 30%.

References

- E. S Allman, C. Matias, J. A Rhodes, et al. Identifiability of parameters in latent structure models with many observed variables. *The Annals of Statistics*, 37 (6A):3099–3132, 2009.
- T. W. Anderson. *An Introduction to Multivariate Statistical Analysis*. Wiley, 2003.
- J. D Banfield and A. E Raftery. Model-based gaussian and non-gaussian clustering. *Biometrics*, pages 803–821, 1993.
- G. Celeux and J. Diebolt. The SEM algorithm: a probabilistic teacher algorithm derived from the EM algorithm for the mixture problem. *Computational Statistics Quarterly*, 1985.
- Jocelyn T Chi, Eric C Chi, and Richard G Baraniuk. k-pod: A method for k-means clustering of missing data. *The American Statistician*, 70(1):91–99, 2016.
- M. Du Roy De Chaumaray and M. Marbac. Clustering data with nonignorable missingness using semi-parametric mixture models. *arXiv preprint*, 2020.

- J. Geweke, M. Keane, and D. Runkle. Alternative computational approaches to inference in the multinomial probit model. *The review of economics and statistics*, 1994.
- J. G Ibrahim, M.-H. Chen, and S. R Lipsitz. Missing responses in generalised linear mixed models when the missing data mechanism is nonignorable. *Biometrika*, 2001.
- Marine Le Morvan, Julie Josse, Thomas Moreau, Erwan Scornet, and Gaël Varoquaux. Neumiss networks: differentiable programming for supervised learning with missing values. *Advances in Neural Information Processing Systems*, 33: 5980–5990, 2020.
- G. J McLachlan and K. E Basford. *Mixture models: Inference and applications to clustering*. M. Dekker New York, 1988.
- D. McParland and Isobel C. Gormley. Model based clustering for mixed data: clustmd. *Advances in Data Analysis and Classification*, 10(2):155–169, 2016.
- W. Miao, P. Ding, and Z. Geng. Identifiability of normal and normal mixture models with nonignorable missing data. *Journal of the American Statistical Association*, 2016.
- K. Mohan. On handling self-masking and other hard missing data problems. 2018.
- K. Mohan, F. Thoemmes, and J. Pearl. Estimation with incomplete data: The linear case. In *IJCAI*, pages 5082–5088, 2018.
- Donald B Rubin. Inference and missing data. *Biometrika*, 63(3):581–592, 1976.
- Aude Sportisse, Claire Boyer, and Julie Josse. Imputation and low-rank estimation with missing not at random data. *Statistics and Computing*, 30(6):1629–1643, 2020.
- Aude Sportisse, Matthieu Marbac, Christophe Biernacki, Claire Boyer, Gilles Celeux, Julie Josse, and Fabien Laporte. Model-based clustering with missing not at random data. *arXiv preprint arXiv:2112.10425*, 2023.
- H. Teicher. Identifiability of finite mixtures. *The annals of Mathematical statistics*, 1963.
- S. J Yakowitz and J. D Spragins. On the identifiability of finite mixtures. *The Annals of Mathematical Statistics*, pages 209–214, 1968.

Supplementary material

This file is a supplementary material. In Appendix A, the proof for Proposition 3.1 is given. Some complements on the EM algorithm are given in Appendix B. The SEM algorithm presented in Section 4 is detailed in Appendix C. Appendix D gives the values of hyperparameters for the numerical experiments on synthetic data.

A Appendix 1: Identifiability

A.1 Continuous and count data

Proof of Proposition 3.1, continuous case. Suppose there exists two sets of parameters $\{\gamma, \psi\}$ and $\{\gamma', \psi'\}$ which have the same observed distribution, i.e., $f(\mathbf{y}_i^{\text{obs}}, \mathbf{c}_i; \gamma, \psi) = f(\mathbf{y}_i^{\text{obs}}, \mathbf{c}_i; \gamma', \psi')$. More precisely, one has

$$\begin{aligned} \forall \mathbf{y}_i \in \mathbb{R}^d, \forall \mathbf{c}_i \in \{0, 1\}^d, \sum_{k=1}^K \int_{\mathcal{Y}_i^{\text{mis}}} \pi_k f_k(\mathbf{y}_i; \lambda_k) \prod_{j=1}^d \rho(\alpha_{kj} + \beta_{kj} y_{ij})^{c_{ij}} [1 - \rho(\alpha_{kj} + \beta_{kj} y_{ij})]^{1-c_{ij}} dy^{\text{mis}} \\ = \sum_{k=1}^{K'} \int_{\mathcal{Y}_i^{\text{mis}}} \pi'_k f_k(\mathbf{y}_i; \lambda'_k) \prod_{j=1}^d \rho((\alpha')_{kj} + (\beta')_{kj} y_{ij})^{c_{ij}} [1 - \rho((\alpha')_{kj} + (\beta')_{kj} y_{ij})]^{1-c_{ij}} dy^{\text{mis}} \end{aligned}$$

Let us consider the case when $c_{ij} = 0$ for all $j = 1, \dots, d$. One has

$$\sum_{k=1}^K \pi_k f_k(\mathbf{y}_i; \lambda_k) \prod_{j=1}^d (1 - \rho(\alpha_{kj} + \beta_{kj} y_{ij})) = \sum_{k=1}^{K'} \pi'_k f_k(\mathbf{y}_i; \lambda'_k) \prod_{j=1}^d (1 - \rho((\alpha')_{kj} + (\beta')_{kj} y_{ij})).$$

By using the identifiability of the marginal mixture, one obtains $\lambda_k = \lambda'_k$ and $\pi_k \prod_{j=1}^d (1 - \rho(\alpha_{kj} + \beta_{kj} y_{ij})) = \pi'_k \prod_{j=1}^d (1 - \rho((\alpha')_{kj} + (\beta')_{kj} y_{ij}))$.

In the sequel, we use the same reasoning of Theorem 2 in (Teicher, 1963). Let us denote $F_k(y_{ij}) = f_{kj}(y_{ij}; \lambda_{kj}) \prod_{j=1}^d (1 - \rho(\alpha_{kj} + \beta_{kj} y_{ij}))$ and $F'_k(y_{ij}) = f_{kj}(y_{ij}; \lambda'_{kj}) \prod_{j=1}^d (1 - \rho((\alpha')_{kj} + (\beta')_{kj} y_{ij}))$. Without loss of generality, assume that $F_k \prec F_l$ and $F'_k \prec F'_l$ for $k < l$. If $F_1 \neq F'_1$, we assume also without loss of generality that $F_1 \preceq F'_1$. Then, $F_1 \prec F'_k$ for $1 \leq k \leq K'$. For $u \in T_1$, where $T_1 = S_{F_1} \cap \{u : F_1(u) \neq 0\}$ is the domain of definition of F_1 such that $f_{1j}(u; \lambda_{1j}) \prod_{j=1}^d (1 - \rho(\alpha_{1j} + \beta_{1j} u)) \neq 0$, one has

$$\pi_1 + \sum_{k=1}^K \pi_k \frac{F_k(u)}{F_1(u)} = \sum_{k=1}^{K'} \pi'_k \frac{F'_k(u)}{F_1(u)}.$$

Letting $u \rightarrow +\infty$, $\pi_1 = 0$ which is in contradiction with the mixture model (where $\pi_k > 0$, $\forall k = 1, \dots, K$). It implies that $F_1 = F'_1$. For any $u \in T_1$, one has

$$\pi_1 + \sum_{k=2}^K \pi_k \frac{F_k(u)}{F_1(u)} = \pi'_1 + \sum_{k=2}^{K'} \pi'_k \frac{F'_k(u)}{F_1(u)}.$$

Letting again $u \rightarrow +\infty$, one obtains $\pi_1 = \pi'_1$ and $\sum_{k=2}^K \pi_k \frac{F_k(u)}{F_1(u)} = \sum_{k=2}^{K'} \pi'_k \frac{F'_k(u)}{F_1(u)}$. We repeat this argument to conclude that $F_k = F'_k$ and $\pi_k = \pi'_k$ for $k = 1, \dots, \min\{K, K'\}$. Finally, if $K \neq K'$, say $K > K'$, $\sum_{k=K'+1}^K \pi_k F_k(u) = 0$ implies $\pi_k = 0$ for $K' + 1 \leq k \leq K$ which is in contradiction with the definition of the mixture model. Thus $K = K'$.

Finally, $F_k = F'_k$ implies that $\prod_{j=1}^d (1 - \rho(\alpha_{kj} + \beta_{kj} y_{ij})) = \prod_{j=1}^d (1 - \rho((\alpha')_{kj} + (\beta')_{kj} y_{ij}))$. By integrating out over all the elements but the j -th element, one has for all $y_{ij} \in \mathbb{R}$, $\rho(\alpha_{kj} + \beta_{kj} y_{ij}) = \rho((\alpha')_{kj} + (\beta')_{kj} y_{ij})$. and $\alpha_{kj} = (\alpha')_{kj}$ and $\beta_{kj} = (\beta')_{kj}$, since ρ is an injective function. Indeed, ρ is assumed to be strictly monotone. \square

A.2 On identifiability of the Gaussian mixture

Finite Gaussian mixtures are identifiable and, for any variable j , there is a total ordering defined by $\sigma_{kj}^2 > \sigma_{(k+1)j}^2$ and $\mu_{kj} > \mu_{(k+1)j}$ if $\sigma_{kj}^2 = \sigma_{(k+1)j}^2$, where μ_{kj} and σ_{kj}^2 are respectively the mean and the variance of variable j under component k . Example A.1 shows that the identifiability holds for Gaussian mixtures when there are missing values and that the distribution of the MNAR mechanism is a probit one.

Example A.1 (Gaussian + Probit). *Let us consider that ρ is the probit function and f_k (respectively f_{k+1}) the Gaussian density with parameters (μ_k, σ_k) (respectively $(\mu_{k+1}, \sigma_{k+1})$). Suppose without loss of generality that $\beta_k \geq \beta_{k+1}$. One want to prove that*

$$\lim_{u \rightarrow +\infty} E_u := \lim_{u \rightarrow +\infty} \frac{\int_{-\infty}^{\alpha_{k+1} + \beta_{k+1} u} e^{-t^2/2} dt}{\int_{-\infty}^{\alpha_k + \beta_k u} e^{-t^2/2} dt} \frac{\sigma_k \exp - \frac{(u - \mu_{k+1})^2}{2\sigma_{k+1}^2}}{\sigma_{k+1} \exp - \frac{(u - \mu_k)^2}{2\sigma_k^2}} = 0$$

Let us denote $\phi(u) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^u e^{-t^2/2} dt$. One has

$$\lim_{u \rightarrow +\infty} \phi(u) = \begin{cases} 1 & \text{if } u > 0 \\ 1/2 & \text{if } u = 0 \\ 0 & \text{if } u < 0 \end{cases} \quad (15)$$

- If $\beta_{k+1} > 0$ (and $\beta_k > 0$):

$$\lim_{u \rightarrow +\infty} E_u = \lim_{u \rightarrow +\infty} \exp - \left(u^2 \left(\frac{1}{2\sigma_{k+1}^2} - \frac{1}{2\sigma_k^2} \right) + u \left(\frac{\mu_k}{\sigma_k} - \frac{\mu_{k+1}}{\sigma_{k+1}} \right) \right) = 0.$$

assuming without loss of generality that $\sigma_k^2 > \sigma_{k+1}^2$ or $\mu_k > \mu_{k+1}$ if $\sigma_k^2 = \sigma_{k+1}^2$.

- If $\beta_{k+1} \leq 0$ (and $\beta_k \geq 0$):

$$\lim_{u \rightarrow +\infty} E_u = 0$$

since

$$\lim_{u \rightarrow +\infty} \exp - \left(u^2 \left(\frac{1}{2\sigma_{k+1}^2} - \frac{1}{2\sigma_k^2} \right) + u \left(\frac{\mu_k}{\sigma_k} - \frac{\mu_{k+1}}{\sigma_{k+1}} \right) \right) = 0$$

and

$$\lim_{u \rightarrow +\infty} \frac{\int_{-\infty}^{\alpha_{k+1} + \beta_{k+1} u} e^{-t^2/2} dt}{\int_{-\infty}^{\alpha_k + \beta_k u} e^{-t^2/2} dt} = \begin{cases} 0 & \text{if } \beta_{k+1} < 0 \\ 1/2 & \text{if } \beta_{k+1} = 0 \text{ and } \beta_k > 0 \\ 1 & \text{if } \beta_{k+1} = 0 \text{ and } \beta_k = 0. \end{cases} \quad (16)$$

- If $\beta_{k+1} < 0$ and $\beta_k < 0$: One uses the upper and lower bounds for the probit function.

$$\frac{1}{-t + \sqrt{t^2 + 4}} < \sqrt{\frac{\pi}{2}} \exp \frac{t^2}{2} \phi(t) < \frac{1}{-t + \sqrt{t^2 + 8/\pi}},$$

i.e., $\phi(t) < \sqrt{\frac{2}{\pi}} \frac{1}{-t + \sqrt{t^2 + 8/\pi}} \exp -\frac{t^2}{2}$ and $\frac{1}{\phi(t)} < (-t + \sqrt{t^2 + 4}) \sqrt{\frac{\pi}{2}} \exp \frac{t^2}{2}$ Thus, noting that $\lim_{u \rightarrow +\infty} \phi(\alpha_{k+1} + \beta_{k+1}u) = \lim_{u \rightarrow +\infty} \phi(\beta_{k+1}u)$,

$$\frac{\int_{-\infty}^{\alpha_{k+1} + \beta_{k+1}u} e^{-t^2/2} dt}{\int_{-\infty}^{\alpha_k + \beta_k u} e^{-t^2/2} dt} \underset{u \rightarrow +\infty}{=} \frac{\phi(\beta_{k+1}u)}{\phi(\beta_k u)} \underset{u \rightarrow +\infty}{<} \exp \left(\left(\frac{\beta_k^2}{2} - \frac{\beta_{k+1}^2}{2} \right) u^2 \right). \quad (17)$$

As $\beta_{k+1} \leq \beta_k < 0$, one has $\beta_k^2/2 - \beta_{k+1}^2/2 < 0$ and it implies

$$\lim_{u \rightarrow +\infty} \frac{\int_{-\infty}^{\alpha_{k+1} + \beta_{k+1}u} e^{-t^2/2} dt}{\int_{-\infty}^{\alpha_k + \beta_k u} e^{-t^2/2} dt} = 0.$$

Given that

$$\lim_{u \rightarrow +\infty} \exp - \left(u^2 \left(\frac{1}{2\sigma_{k+1}^2} - \frac{1}{2\sigma_k^2} \right) + u \left(\frac{\mu_k}{\sigma_k} - \frac{\mu_{k+1}}{\sigma_{k+1}} \right) \right) = 0,$$

assuming without loss of generality that $\sigma_k^2 > \sigma_{k+1}^2$ or $\mu_k > \mu_{k+1}$ if $\sigma_k^2 = \sigma_{k+1}^2$, one has

$$\lim_{u \rightarrow +\infty} E_u = 0.$$

This result has been already stated, in the case of univariate distributions, by [Miao et al. \(2016\)](#). In particular, the identifiability conditions in [Miao et al. \(2016\)](#) (conditions 1 and 2 of their paper) imply the existence of the total ordering defined in Assumption [A2.](#). However, these conditions exclude the case of Gaussian mixture with a logistic missing-data distribution, which is very used in practice. In Corollary [A.2](#), we therefore extend this result to the multivariate case with a logistic missing-data distribution.

Note first that with a logistic distribution, a total ordering cannot be defined. Indeed, for variable j , such an ordering cannot be defined if the two univariate variances are equal (i.e., $\sigma_{kj}^2 = \sigma_{(k+1)j}^2$) and $\mu_{kj} - \beta_{kj} - \mu_{(k+1)j} + \beta_{(k+1)j} = 0$. However, for the specific case of Gaussian mixture where all the univariate variances are different between the components, then conditions of Proposition [3.1](#) hold true with a logistic missing-data distribution and so does its identifiability. In addition, for more parsimonious MNAR models for which the effect on the variable j does not depend on the class membership k (i.e., $\beta_{kj} = \beta_{(k+1)j}$), the conditions of Proposition [3.1](#) hold true with a logistic missing-data distribution. Finally, as stated by Corollary [A.2](#) below, the condition on the covariance matrices (including the case of homoscedastic Gaussian mixture) can be relaxed to obtain the generic identifiability of the model (i.e., all not-identifiable parameter choices lie within a proper submanifold, and thus form a set of Lebesgue zero measure; [Allman et al. \(2009\)](#)).

Corollary A.2. Assume that $\sum_{k=1}^K \pi_k f_k(y_i; \lambda_k)$ is a multivariate Gaussian mixture, ρ is the logistic function and that the missingness scenario is defined by (3), (5) or (8), then, the parameters (π, λ, ψ) of the model given by (2) are generically identifiable up to label swapping, i.e., all not-identifiable parameter choices lie within a proper submanifold, and thus form a set of Lebesgue zero measure.

For the other MNAR models given in (4), (6), (7), (9) and (10), the parameters (π, λ, ψ) of the model given by (2) are identifiable up to label swapping.

Proof of Corollary A.2. We use Proposition [3.1](#). We fix j . By abuse of notation, α_k, β_k, μ_k and σ_k correspond to the parameters $\alpha_{kj}, \beta_{kj}, \mu_{kj}$ and Σ_{kj} of the variable j . Let us first consider the missing scenarios (3), (5) and (8) for which

$\beta_k \neq \beta_{k+1}$. To obtain the total ordering, we need to prove that

$$\lim_{u \rightarrow +\infty} E_u = \frac{(1 + e^{-\alpha_k - \beta_k u}) e^{-\frac{(u - \mu_{k+1})^2}{2\sigma_{k+1}^2}} \sigma_k}{(1 + e^{-\alpha_{k+1} - \beta_{k+1} u}) e^{-\frac{(u - \mu_k)^2}{2\sigma_k^2}} \sigma_{k+1}} = 0.$$

- If $\sigma_k^2 > \sigma_{k+1}^2$, $\lim_{u \rightarrow +\infty} E_u = \lim_{u \rightarrow +\infty} \exp -\frac{1}{2} \left(\frac{1}{\sigma_{k+1}^2} - \frac{1}{\sigma_k^2} \right) u^2 = 0$.
- If $\sigma_k^2 = \sigma_{k+1}^2$, one has $\lim_{u \rightarrow +\infty} E_u = \lim_{u \rightarrow +\infty} \exp((\mu_k - \beta_k) - (\mu_{k+1} - \beta_{k+1}))u = 0$ discarding the case where $(\mu_k - \beta_k) - (\mu_{k+1} - \beta_{k+1}) = 0$ and assuming without loss of generality that $(\mu_k - \beta_k) > (\mu_{k+1} - \beta_{k+1})$. The set of nonidentifiable parameters is $\{\mu_k, \beta_k, \mu_{k+1}, \beta_{k+1} \text{ s.t. } (\mu_k - \beta_k) - (\mu_{k+1} - \beta_{k+1}) = 0\}_{k=1, \dots, K}$ and is of Lebesgue zero measure.

Finally, for the missing scenarios (9) and (10), note that $\beta_k = \beta_{k+1} = 0$. For the missing scenarios (4), (6) and (7), one has $\beta_k = \beta_{k+1}$. Following the same reasoning as above, in the case where $\sigma_{k+1}^2 = \sigma_k^2$, one obtains the set of nonidentifiable parameters such that $\mu_k = \mu_{k+1}$, which is empty since $\mu_k \neq \mu_{k+1}$ if $\sigma_k^2 = \sigma_{k+1}^2$. \square

A.3 On identifiability of the Poisson mixture

Proposition A.2 can also be applied for variables with integer value (*i.e.*, count data), as shown below in Examples A.3 and A.4 for the Poisson mixture with probit or logistic missing-data distributions.

Example A.3 (Poisson + Probit). *Considering that ρ is the probit function and f_k (respectively f_{k+1}) the Poisson distribution with parameters λ_k (respectively λ_{k+1}). Suppose without loss of generality that $\beta_k > \beta_{k+1}$ and $\lambda_k > \lambda_{k+1}$. One want to prove*

$$\lim_{u \rightarrow +\infty} E_u := \lim_{u \rightarrow +\infty} \frac{\int_{-\infty}^{\alpha_{k+1} + \beta_{k+1} u} e^{-t^2/2} dt}{\int_{-\infty}^{\alpha_k + \beta_k u} e^{-t^2/2} dt} \frac{\lambda_{k+1}^u e^{-\lambda_{k+1}}}{\lambda_k^u e^{-\lambda_k}} = 0.$$

- If $\beta_{k+1} > 0$ (and $\beta_k > 0$): using (15), one has

$$\lim_{u \rightarrow +\infty} E_u = \lim_{u \rightarrow +\infty} \exp u \log \frac{\lambda_{k+1}}{\lambda_k} = 0.$$

- If $\beta_{k+1} \leq 0$ (and $\beta_k \geq 0$): one has

$$\lim_{u \rightarrow +\infty} E_u = 0.$$

using

$$\lim_{u \rightarrow +\infty} \exp u \log \frac{\lambda_{k+1}}{\lambda_k} = 0$$

and (16) for the missing distribution part.

- If $\beta_{k+1} < 0$ and $\beta_k < 0$: using (17), one obtains

$$\lim_{u \rightarrow +\infty} E_u < \lim_{u \rightarrow +\infty} \exp \left(\left(\frac{\beta_k^2}{2} - \frac{\beta_{k+1}^2}{2} \right) u^2 \right) \exp u \log \frac{\lambda_{k+1}}{\lambda_k} = 0,$$

because $\beta_k^2/2 - \beta_{k+1}^2/2 < 0$.

Example A.4 (Poisson + Logistic). Considering that ρ is the logistic function and f_k (respectively f_{k+1}) the Poisson distribution with parameters λ_k (respectively λ_{k+1}). One want to prove that

$$\lim_{u \rightarrow +\infty} E_u = \lim_{u \rightarrow +\infty} \frac{1 + e^{-\alpha_k - \beta_k u}}{1 + e^{-\alpha_{k+1} - \beta_{k+1} u}} \exp u \log \frac{\lambda_{k+1}}{\lambda_k} = 0.$$

Assume that $\lambda_k > \lambda_{k+1}$ without loss of generality.

- For the missing scenarios (3), (5) and (8) for which $\beta_k \neq \beta_{k+1}$, one obtains the generic identifiability where the set of non-identifiable parameters is $\{\alpha_k, \beta_k, \lambda_k \text{ s.t. } (\log \lambda_k - \beta_k) - (\log \lambda_{k+1} - \beta_{k+1}) = 0\}_{k=1, \dots, K}$ and is of Lebesgue zero measure.
- For the missing scenarios (9) and (10), note that $\beta_k = \beta_{k+1} = 0$. For the missing scenarios (4), (6) and (7), one has $\beta_k = \beta_{k+1}$. It implies that idenfiability holds since

$$\lim_{u \rightarrow +\infty} E_u = \lim_{u \rightarrow +\infty} \exp u \log \frac{\lambda_{k+1}}{\lambda_k} = 0.$$

A.4 Categorical data

Proposition 3.1 states that generic identifiability holds only for the MNAR z and the MNAR z^j missing scenarios and that the other missing scenarios lead to non-identifiable models. The proof uses Corollary 5 of Allman et al. (2009) which gives the identifiability of finite mixtures of Bernoulli products.

Proof of Proposition 3.1, categorical case. Let us first consider the case where $\beta_{kj} = (0, \dots, 0) \in \mathbb{R}^{\ell_j}, \forall k = 1, \dots, K, \forall j = 1, \dots, d$. Suppose there exists two sets of parameters $\{\gamma, \psi\}$ and $\{\gamma', \psi'\}$ which have the same observed distribution.

$$\begin{aligned} \forall \mathbf{y}_i \in \mathbb{R}^d, \forall \mathbf{c}_i \in \{0, 1\}^d, \sum_{k=1}^K \int_{\mathcal{Y}_i^{\text{mis}}} \pi_k f_k(\mathbf{y}_i; \lambda_k) \prod_{j=1}^d \rho(\alpha_{kj})^{c_{ij}} [1 - \rho(\alpha_{kj})]^{1-c_{ij}} dy^{\text{mis}} \\ = \sum_{k=1}^{K'} \int_{\mathcal{Y}_i^{\text{mis}}} \pi'_k f_k(\mathbf{y}_i; \lambda'_k) \prod_{j=1}^d \rho((\alpha')_{kj})^{c_{ij}} [1 - \rho(\alpha'_{kj})]^{1-c_{ij}} dy^{\text{mis}}. \end{aligned}$$

Let us consider the case where all the elements of \mathbf{y}_i are observed, i.e., $c_{ij} = 0, \forall j = 1, \dots, d$. One has

$$\sum_{k=1}^K \pi_k f_k(\mathbf{y}_i; \lambda_k) \prod_{j=1}^d (1 - \rho(\alpha_{kj})) = \sum_{k=1}^{K'} \pi'_k f_k(\mathbf{y}_i; \lambda'_k) \prod_{j=1}^d (1 - \rho(\alpha'_{kj})),$$

i.e., by independence to the group membership,

$$\begin{aligned} \sum_{k=1}^K \pi_k \prod_{j=1}^d f_{kj}(y_{ij}; \lambda_{kj}) (1 - \rho(\alpha_{kj})) &= \sum_{k=1}^{K'} \pi'_k \prod_{j=1}^d f_{kj}(y_{ij}; \lambda'_{kj}) (1 - \rho(\alpha'_{kj})), \\ \Leftrightarrow \sum_{k=1}^K \pi_k \prod_{j=1}^d (1 - \rho(\alpha_{kj}))^{1-c_{ij}} \prod_{h=1}^{\ell_j} (\lambda_{kj}^h)^{y_{ij}^h} &= \sum_{k=1}^{K'} \pi'_k \prod_{j=1}^d (1 - \rho(\alpha_{kj}))^{1-c_{ij}} \prod_{h=1}^{\ell_j} ((\lambda'_{kj})^h)^{y_{ij}^h}. \end{aligned}$$

We recognize the finite mixture of K multinomial distributions with d components for $w_{ij} = ((1 - c_{ij}), y_{ij})$, $j = 1, \dots, d$ with paramaters $(\lambda_{kj}) = ((1 - \rho(\alpha_{kj})), \lambda_{kj}^1, \dots, \lambda_{kj}^{\ell_j})$, $j = 1, \dots, d$ and proportions π_k . We can thus apply Theorem 4 (Allman et al., 2009) with the model $\mathcal{M}(K; \ell_1, \dots, \ell_d)$ which gives the generic identifiability of the model paramaters up to a label swapping, *i.e.*,

$$\begin{aligned}\forall k, \forall j, \lambda_{kj}^h &= (\lambda'_{kj})^h \\ \forall k, \forall j, \rho(\alpha_{kj}) &= \rho(\alpha'_{kj}) \\ \forall k, \pi_k &= \pi'_k\end{aligned}$$

As the function ρ is strictly monotone, the equality $\rho(\alpha_{kj}) = \rho(\alpha'_{kj})$ implies $\alpha_{kj} = \alpha'_{kj}$. In addition, if $K \neq K'$, say $K > K'$, $\sum_{k=K'+1}^K \pi_k \prod_{j=1}^d (1 - \rho(\alpha_{kj})) \prod_{h=1}^{\ell_j} (\lambda_{kj}^h)^{y_{ij}^h} = 0$ implies $\pi_k = 0$ for $K' + 1 \leq k \leq K$.

We consider now the missing scenarios for which $\beta_{kj} \neq 0$. The identifiability does not hold. We can present a counter-example. The set of parameters $\psi = \{\alpha = (1, \dots, 1), \beta = (1, \dots, 1)\}$ has the same observed distribution than another set of parameters $\psi' = \{\alpha' = (0, \dots, 0), \beta' = (2, \dots, 2)\}$. Indeed, in the case where $y_{ij} = (1, \dots, 1)$, $\rho(\alpha_{kj} + \beta_{kj} y_{ij}) = \rho(\alpha'_{kj} + \beta'_{kj} y_{ij})$. □

B Appendix 2: Details on EM algorithm

The EM algorithm consists on two steps iteratively proceeded: the E-step and M-step. For the E-step, one has

$$Q(\theta; \theta^{[r-1]}) = \mathbb{E}[\ell_{\text{comp}}(\theta; \mathbf{y}, \mathbf{z}, \mathbf{c}) | \mathbf{y}_i^{\text{obs}}, \mathbf{c}_i; \theta^{[r-1]}],$$

where

$$\ell_{\text{comp}}(\theta; \mathbf{Y}, \mathbf{Z}, \mathbf{C}) = \sum_{i=1}^n \sum_{k=1}^K z_{ik} \log (\pi_k f_k(\mathbf{y}_i; \lambda_k) f_k(\mathbf{c}_i | \mathbf{y}_i; \psi_k)).$$

It leads to the decomposition

$$Q(\theta; \theta^{[r-1]}) = \sum_{i=1}^n \sum_{k=1}^K t_{ik}(\theta^{[r-1]}) \left[\log(\pi_k) + \tau_y(\lambda_k; \mathbf{y}_i^{\text{obs}}, \mathbf{c}_i, \theta^{[r-1]}) + \tau_c(\psi_k; \mathbf{y}_i^{\text{obs}}, \mathbf{c}_i, \theta^{[r-1]}) \right],$$

with $t_{ik}(\theta^{[r-1]}) = f(z_{ik} = 1 | \mathbf{y}_i^{\text{obs}}, \mathbf{c}_i; \theta^{[r-1]})$. The terms involved in this decomposition are now detailed.

- (a) the expectation of the data mixture part over the missing values given the available information (*i.e.*, the observed data and the indicator pattern), the class membership and the current value of the parameters:

$$\tau_y(\lambda_k; \mathbf{y}_i^{\text{obs}}, \mathbf{c}_i, \theta^{[r-1]}) = \mathbb{E}_{\theta^{[r-1]}} \left[\log f_k(y_i; \lambda_k) | \mathbf{y}_i^{\text{obs}}, z_{ik} = 1, \mathbf{c}_i \right],$$

- (b) the expectation of the missing mechanism part over the missing values given the available information, the class membership and the current value of the parameters:

$$\tau_c(\psi_k; \mathbf{y}_i^{\text{obs}}, \mathbf{c}_i, \theta^{[r-1]}) = \mathbb{E}_{\theta^{[r-1]}} \left[\log f_k(c_i | y_i; \psi_k) | \mathbf{y}_i^{\text{obs}}, z_{ik} = 1, \mathbf{c}_i \right].$$

(c) the conditional probability for an observation i to belong to the class k given the available information and the current value of the parameters:

$$t_{ik}(\theta^{[r-1]}) = f(z_{ik} = 1 \mid \mathbf{y}_i^{\text{obs}}, \mathbf{c}_i; \theta^{[r-1]}).$$

Terms (a) and (b) require to integrate over the distribution $f(\mathbf{y}_i^{\text{mis}} \mid \mathbf{y}_i^{\text{obs}}, z_{ik} = 1, \mathbf{c}_i; \theta^{[r-1]})$. Term (c) corresponds to the conditional probability for an observation i to arise from the k th mixture component with the current values of the model parameter. More particularly, one has

$$\begin{aligned} t_{ik}(\theta^{[r-1]}) &= \frac{f(z_{ik} = 1, \mathbf{y}_i^{\text{obs}}, \mathbf{c}_i; \theta^{[r-1]})}{f(\mathbf{y}_i^{\text{obs}}, \mathbf{c}_i; \theta^{[r-1]})} \\ &= \frac{\pi_k^{[r-1]} f_k(\mathbf{y}_i^{\text{obs}}; \lambda_k^{[r-1]}) f(\mathbf{c}_i \mid \mathbf{y}_i^{\text{obs}}, z_{ik} = 1; \theta^{[r-1]})}{\sum_{h=1}^K \pi_h^{[r-1]} f_h(\mathbf{y}_i^{\text{obs}}; \lambda_h^{[r-1]}) f(\mathbf{c}_i \mid \mathbf{y}_i^{\text{obs}}, z_{ih} = 1; \theta^{[r-1]})} \end{aligned} \quad (18)$$

The quantity that can cause numerical difficulties is the probability $f(\mathbf{c}_i \mid \mathbf{y}_i^{\text{obs}}, z_{ik} = 1; \theta^{[r-1]})$.

Computations for the MNAR z^j model For the MNAR z and MNAR z^j models, the effect of the missingness is only due to the class membership. The EM algorithm for the MNAR z model is detailed in the main paper (Sportisse et al., 2023). Term (a) is the same for both MNAR z and MNAR z^j models but (b) and (c) differ. For Term (b), $f(\mathbf{c}_i \mid \mathbf{y}_i, z_{ik} = 1; \psi)$ is independent of \mathbf{y}_i , which implies

$$\log(f(\mathbf{c}_i \mid z_{ik} = 1; \psi)) = \sum_{j=1}^d c_{ij} \log \rho(\alpha_{kj}) + (1 - c_{ij}) \log(1 - \rho(\alpha_{kj})). \quad (19)$$

For Term (c), one remark that

$$\mathbb{P}(\mathbf{c}_i \mid \mathbf{y}_i^{\text{obs}}, z_{ik} = 1; \theta^{[r-1]}) = \prod_{j=1}^d \mathbb{P}(c_{ij} = 1 \mid \mathbf{y}_i^{\text{obs}}, z_{ik} = 1; \theta^{[r-1]})^{c_{ij}} \mathbb{P}(c_{ij} = 0 \mid \mathbf{y}_i^{\text{obs}}, z_{ik} = 1; \theta^{[r-1]})^{1-c_{ij}} \quad (\text{MNAR}z^j).$$

and that, by independence of \mathbf{y}_i , one has

$$\mathbb{P}(c_{ij} = 1 \mid \mathbf{y}_i^{\text{obs}}, z_{ik} = 1; \theta^{[r-1]}) = \mathbb{P}(c_{ij} = 1 \mid z_{ik} = 1; \theta^{[r-1]}) = \rho(\alpha_{kj}).$$

Using (18), one obtains

$$t_{ik}^{[r-1]}(\theta^{[r-1]}) = \frac{\pi_k^{[r-1]} \phi(\mathbf{y}_i^{\text{obs}}; (\mu_{ik}^{\text{obs}})^{[r-1]}, (\Sigma_{ik}^{\text{obs,obs}})^{[r-1]}) \prod_{j=1}^d \rho(\alpha_{kj}^{[r-1]})^{c_{ij}} (1 - \rho(\alpha_{kj}^{[r-1]}))^{1-c_{ij}}}{\sum_{h=1}^K \pi_h^{[r-1]} \phi(\mathbf{y}_i^{\text{obs}}; (\mu_{ih}^{\text{obs}})^{[r-1]}, (\Sigma_{ih}^{\text{obs,obs}})^{[r-1]}) \prod_{j=1}^d \rho(\alpha_{hj}^{[r-1]})^{c_{ij}} (1 - \rho(\alpha_{hj}^{[r-1]}))^{1-c_{ij}}} \quad (\text{MNAR}z^j). \quad (20)$$

The E-step is derived in the same way as for the MNAR z model with these terms. The M-step for ψ consists of performing a GLM with a binomial link for the matrices $(\mathcal{J}_{kj}^{\text{MNAR}z^j})_{j=1,\dots,d}^{[r]}$ and by giving $t_{ik}(\theta^{[r-1]})$ as prior weights to fit the process.

$$(\mathcal{J}_{kj}^{\text{MNAR}z^j})^{[r]} = c_{.j} \quad 1 \quad (21)$$

Intractability of the EM algorithm for MNAR y_\star models For missing scenarios which model the effect of the missingness depending on the variable, the computations are more difficult.

- Because of the dependence of y , $f(\mathbf{y}_i^{\text{mis}} | \mathbf{y}_i^{\text{obs}}, z_{ik} = 1, \mathbf{c}_i; \theta^{[r-1]}) = f(\mathbf{y}_i^{\text{mis}} | \mathbf{y}_i^{\text{obs}}, z_{ik} = 1; \theta^{[r-1]})$ does not hold anymore. Here, one has

$$\begin{aligned} & f(\mathbf{y}_i^{\text{mis}} | \mathbf{y}_i^{\text{obs}}, z_{ik} = 1, \mathbf{c}_i; \theta^{[r-1]}) \\ &= \frac{\prod_{h=1}^d \rho(\alpha_{kh}^{[r-1]} + \beta_{kh}^{[r-1]} y_{ih}^{\text{mis}})^{c_{ih}} (1 - \rho(\alpha_{kh}^{[r-1]} + \beta_{kh}^{[r-1]} y_{ih}^{\text{obs}}))^{1-c_{ih}} f(\mathbf{y}_i^{\text{mis}}, \mathbf{y}_i^{\text{obs}}, z_{ik} = 1; \theta^{[r-1]})}{\int_{\mathcal{Y}_i^{\text{mis}}} \prod_{h=1}^d \rho(\alpha_{kh}^{[r-1]} + \beta_{kh}^{[r-1]} y_{ih}^{\text{mis}})^{c_{ih}} (1 - \rho(\alpha_{kh}^{[r-1]} + \beta_{kh}^{[r-1]} y_{ih}^{\text{obs}}))^{1-c_{ih}} f(\mathbf{y}_i^{\text{mis}}, \mathbf{y}_i^{\text{obs}}, z_{ik} = 1; \theta^{[r-1]}) d\mathbf{y}_i^{\text{mis}}} \\ &= \frac{\prod_{h, c_{ih}=1} \rho(\alpha_{kh}^{[r-1]} + \beta_{kh}^{[r-1]} y_{ih}^{\text{mis}}) f(\mathbf{y}_i^{\text{mis}} | \mathbf{y}_i^{\text{obs}}, z_{ik} = 1; \theta^{[r-1]})}{\int_{\mathcal{Y}_i^{\text{mis}}} \prod_{h, c_{ih}=1} \rho(\alpha_{kh}^{[r-1]} + \beta_{kh}^{[r-1]} y_{ih}^{\text{mis}}) f(\mathbf{y}_i^{\text{mis}} | \mathbf{y}_i^{\text{obs}}, z_{ik} = 1; \theta^{[r-1]}) d\mathbf{y}_i^{\text{mis}}}. \end{aligned} \quad (22)$$

which implies that Term (a) requires difficult computations if this distribution is not classical.

- For Term (b), it is the same problem, since $f(\mathbf{c}_i | \mathbf{y}_i, z_{ik} = 1; \psi)$ is no longer independent of \mathbf{y} , then it requires a specific numerical integration. Using (22),

$$\begin{aligned} \tau_c(\psi_k; \mathbf{y}_i^{\text{obs}}, \mathbf{c}_i, \theta^{[r-1]}) &= \mathbb{E} \left[\log \left(\prod_{j=1}^d \rho(\alpha_{kj} + \beta_{kj} y_{ij})^{c_{ij}} (1 - \rho(\alpha_{kj} + \beta_{kj} y_{ij}))^{1-c_{ij}} \right) | \mathbf{y}_i^{\text{obs}}, z_{ik} = 1, \mathbf{c}_i; \theta^{[r-1]} \right] \\ &= \sum_{j=1}^d c_{ij} \int_{\mathcal{Y}_{ij}^{\text{mis}}} \log(\rho(\alpha_{kj} + \beta_{kj} y_{ij}^{\text{mis}})) f(y_{ij}^{\text{mis}} | \mathbf{y}_i^{\text{obs}}, z_{ik} = 1, \mathbf{c}_i; \theta^{[r-1]}) dy_{ij}^{\text{mis}} \\ &\quad + (1 - c_{ij}) \log(1 - \rho(\alpha_{kj} + \beta_{kj} y_{ij}^{\text{obs}})) \end{aligned}$$

where

$$\begin{aligned} & f(y_{ij}^{\text{mis}} | y_i^{\text{obs}}, z_{ik} = 1, \mathbf{c}_i; \theta^{[r-1]}) \\ &= \frac{\rho(\alpha_{kj}^{[r-1]} + \beta_{kj}^{[r-1]} y_{ij}^{\text{mis}})^{c_{ij}} (1 - \rho(\alpha_{kj}^{[r-1]} + \beta_{kj}^{[r-1]} y_{ij}^{\text{obs}}))^{1-c_{ij}} f(y_{ij}^{\text{mis}} | y_i^{\text{obs}}, z_{ik} = 1; \theta^{[r-1]})}{\int_{\mathcal{Y}_{ij}^{\text{mis}}} \rho(\alpha_{kj}^{[r-1]} + \beta_{kj}^{[r-1]} y_{ij}^{\text{mis}})^{c_{ij}} (1 - \rho(\alpha_{kj}^{[r-1]} + \beta_{kj}^{[r-1]} y_{ij}^{\text{obs}}))^{1-c_{ij}} f(y_{ij}^{\text{mis}} | y_i^{\text{obs}}, z_{ik} = 1; \theta^{[r-1]}) dy_{ij}^{\text{mis}}}. \end{aligned}$$

Therefore,

$$\begin{aligned} & \tau_c(\psi_k; \mathbf{y}_i^{\text{obs}}, \mathbf{c}_i, \theta^{[r-1]}) \\ &= \sum_{j=1}^d c_{ij} \int_{\mathcal{Y}_{ij}^{\text{mis}}} \log(\rho(\alpha_{kj} + \beta_{kj} y_{ij}^{\text{mis}})) \frac{\rho(\alpha_{kj}^{[r-1]} + \beta_{kj}^{[r-1]} y_{ij}^{\text{mis}})^{c_{ij}} f(y_{ij}^{\text{mis}} | y_i^{\text{obs}}, z_{ik} = 1; \theta^{[r-1]})}{\int_{\mathcal{Y}_{ij}^{\text{mis}}} \rho(\alpha_{kj}^{[r-1]} + \beta_{kj}^{[r-1]} x)^{c_{ij}} f(x | y_i^{\text{obs}}, z_{ik} = 1; \theta^{[r-1]}) dx} dy_{ij}^{\text{mis}} \\ &\quad + (1 - c_{ij}) \log(1 - \rho(\alpha_{kj} + \beta_{kj} y_{ij}^{\text{obs}})). \end{aligned}$$

- There is no closed-form expression for Term (c).

$$\begin{aligned} & f(c_{ij} | \mathbf{y}_i^{\text{obs}}, z_{ik} = 1; \theta^{[r-1]}) \\ &= \int_{\mathcal{Y}_{ij}^{\text{mis}}} f(c_{ij} | \mathbf{y}_i^{\text{obs}}, y_{ij}^{\text{mis}}, z_{ik} = 1; \psi^{[r-1]}) f(\mathbf{y}_{ij}^{\text{mis}} | \mathbf{y}_i^{\text{obs}}, z_{ik} = 1; \theta^{[r-1]}) dy_{ij}^{\text{mis}} \\ &= c_{ij} \int_{-\infty}^{+\infty} \rho(\alpha_{kj}^{[r-1]} + \beta_{kj}^{[r-1]} y_{ij}^{\text{mis}}) \phi(y_{ij}^{\text{mis}}; (\tilde{\mu}_{ik}^{\text{mis}})_j^{[r-1]}, (\tilde{\Sigma}_{ik}^{\text{mis}})_{jj}^{[r-1]}) dy_{ij}^{\text{mis}} + (1 - c_{ij})(1 - \rho(\alpha_{kj}^{[r-1]} + \beta_{kj}^{[r-1]} y_{ij}^{\text{obs}})). \end{aligned} \quad (23)$$

Using (18), the probabilities $t_{ik}(\theta^{[r-1]})$ can be deduced from Equation (23).

Let us detail the difficulties for two particular cases, if ρ is logistic or probit.

- **ρ is logistic:** Equation (22) leads to none classical distribution because

$$f(\mathbf{y}_i^{\text{mis}} \mid \mathbf{y}_i^{\text{obs}}, z_{ik} = 1, \mathbf{c}_i; \theta^{[r-1]}) \propto \prod_{h, c_{ih}=1} \frac{1}{\exp(-(\alpha_{kh}^{[r-1]} + \beta_{kh}^{[r-1]} y_{ih}^{\text{mis}}))} \phi(y_i^{\text{mis}}; (\tilde{\mu}_{ik}^{\text{mis}})^{[r-1]}, (\tilde{\Sigma}_{ik}^{\text{mis}})^{[r-1]}).$$

Term (b) is

$$\begin{aligned} \tau_c(\psi_k; \mathbf{y}_i^{\text{obs}}, \mathbf{c}_i, \theta^{[r-1]}) \\ \propto \sum_{j=1}^d c_{ij} \int_{\mathcal{Y}_{ij}^{\text{mis}}} -\frac{\log(1 + \exp(-(\alpha_{kj} + \beta_{kj} y_{ij}^{\text{mis}})))}{1 + \exp(-(\alpha_{kj}^{[r-1]} + \beta_{kj}^{[r-1]} y_{ij}^{\text{mis}}))} \phi(y_{ij}^{\text{mis}}; (\tilde{\mu}_{ik}^{\text{mis}})_j^{[r-1]}, (\tilde{\Sigma}_{ik}^{\text{mis}})_{jj}^{[r-1]}) dy_{ij}^{\text{mis}} \\ - (1 - c_{ij}) \log(1 + \exp(\alpha_{kj} + \beta_{kj} y_{ij}^{\text{obs}})), \end{aligned}$$

which amounts to compute the Gaussian moment of $\frac{\log(1+\exp(-u))}{1+\exp(-u)}$, but it has no closed form to our knowledge.

Finally, Equation (23) does not have a closed form either because it requires the computation of

$$\int_{-\infty}^{+\infty} \frac{1}{1 + \exp(-(\alpha_{kj}^{[r-1]} + \beta_{kj}^{[r-1]} y_{ij}^{\text{mis}}))} \phi(y_{ij}^{\text{mis}}; (\tilde{\mu}_{ik}^{\text{mis}})_j^{[r-1]}, (\tilde{\Sigma}_{ik}^{\text{mis}})_{jj}^{[r-1]}) dy_{ij}^{\text{mis}},$$

i.e., the computation of the Gaussian moment of $\frac{1}{1+\exp(-u)}$.

- **ρ is Probit:** One can prove (presented in Appended C) that the conditional distribution $(\mathbf{y}_i^{\text{mis}} \mid \mathbf{y}_i^{\text{obs}}, z_{ik} = 1, \mathbf{c}_i)$ is a truncated Gaussian, which makes possible the computation of Term (a). Term (b) has no closed form to our knowledge

$$\begin{aligned} \tau_c(\psi_k; \mathbf{y}_i^{\text{obs}}, \mathbf{c}_i, \theta^{[r-1]}) \\ \propto \sum_{j=1}^d c_{ij} \int_{\mathcal{Y}_{ij}^{\text{mis}}} \frac{\log\left(\int_{-\infty}^{\alpha_{kj} + \beta_{kj} y_{ij}^{\text{mis}}} e^{-t^2} dt\right)}{1 + \exp(\alpha_{kj}^{[r-1]} + \beta_{kj}^{[r-1]} y_{ij}^{\text{mis}})} \phi(y_{ij}^{\text{mis}}; (\tilde{\mu}_{ik}^{\text{mis}})_j^{[r-1]}, (\tilde{\Sigma}_{ik}^{\text{mis}})_{jj}^{[r-1]}) dy_{ij}^{\text{mis}} \\ - (1 - c_{ij}) \log\left(1 - \int_{-\infty}^{\alpha_{kj} + \beta_{kj} y_{ij}^{\text{obs}}} e^{-t^2} dt\right), \end{aligned}$$

Equation (23) does not have a closed form either because it requires the computation of

$$\int_{-\infty}^{+\infty} \left(\int_{-\infty}^{\alpha_{kj} + \beta_{kj} y_{ij}^{\text{mis}}} e^{-t^2} dt \right) \phi(y_{ij}^{\text{mis}}; (\tilde{\mu}_{ik}^{\text{mis}})_j^{[r-1]}, (\tilde{\Sigma}_{ik}^{\text{mis}})_{jj}^{[r-1]}) dy_{ij}^{\text{mis}}.$$

C Appendix 3: Details on SEM algorithm

The SEM algorithm consists on two steps iteratively proceeded as presented in Section 4. The key issue is to draw the missing data $(\mathbf{y}_i^{\text{mis}})^r$ and \mathbf{z}_i^r according to their current conditional distribution $f(\mathbf{y}_i^{\text{mis}}, \mathbf{z}_i \mid \mathbf{y}_i^{\text{obs}}, \mathbf{c}_i; \pi^{[r-1]}, \lambda^{[r-1]}, \psi^{[r-1]})$. By convenience, we use a Gibbs sampling and simulate two easier probabilities recalled here

$$\mathbf{z}_i^{[r]} \sim f(\cdot \mid \mathbf{y}_i^{[r-1]}, \mathbf{c}_i; \pi^{[r-1]}, \lambda^{[r-1]}, \psi^{[r-1]}) \quad \text{and} \quad (\mathbf{y}_i^{\text{mis}})^{[r]} \sim f(\cdot \mid \mathbf{y}_i^{\text{obs}}, \mathbf{z}_i^r, \mathbf{c}_i; \lambda^{[r-1]}, \psi^{[r-1]}),$$

where $\mathbf{y}_i^{[r-1]} = (\mathbf{y}_i^{\text{obs}}, (\mathbf{y}_i^{\text{mis}})^{[r-1]})$. For the latter distribution, the membership k of $z_i^{[r]}$ is drawn from the multinomial distribution with probabilities $(\mathbb{P}(z_{ik} = 1 \mid \mathbf{y}_i^{[r-1]}, \mathbf{c}_i; \lambda^{[r-1]}, \psi^{[r-1]}))_{k=1, \dots, K}$ detailed as follows

$$\begin{aligned} & \mathbb{P}(z_{ik} = 1 \mid \mathbf{y}_i^{[r-1]}, \mathbf{c}_i; \theta^{[r-1]}) \\ &= \frac{\mathbb{P}(z_{ik} = 1, \mathbf{y}_i^{[r-1]}, \mathbf{c}_i; \pi^{[r-1]}, \lambda^{[r-1]}, \psi^{[r-1]})}{\mathbb{P}(\mathbf{y}_i^{[r-1]}, \mathbf{c}_i; \pi^{[r-1]}, \lambda^{[r-1]}, \psi^{[r-1]})} \end{aligned} \quad (24)$$

$$\begin{aligned} &= \frac{\mathbb{P}(\mathbf{c}_i \mid \mathbf{y}_i^{[r-1]}, z_{ik} = 1; \psi^{[r-1]}) \mathbb{P}(\mathbf{y}_i^{[r-1]} \mid z_{ik} = 1; \lambda^{[r-1]}) \mathbb{P}(z_{ik} = 1; \pi^{[r-1]})}{\sum_{h=1}^K \mathbb{P}(\mathbf{c}_i \mid \mathbf{y}_i^{[r-1]}, z_{ih} = 1; \psi^{[r-1]}) \mathbb{P}(\mathbf{y}_i^{[r-1]} \mid z_{ih} = 1; \lambda^{[r-1]}) \mathbb{P}(z_{ih} = 1; \pi^{[r-1]})} \\ &= \frac{\mathbb{P}(\mathbf{c}_i \mid \mathbf{y}_i^{[r-1]}, z_{ik} = 1; \psi^{[r-1]}) \mathbb{P}(\mathbf{y}_i^{[r-1]} \mid z_{ik} = 1; \lambda^{[r-1]}) \pi_k^{[r-1]}}{\sum_{h=1}^K \mathbb{P}(\mathbf{c}_i \mid \mathbf{y}_i^{[r-1]}, z_{ih} = 1; \psi^{[r-1]}) \mathbb{P}(\mathbf{y}_i^{[r-1]} \mid z_{ih} = 1; \lambda^{[r-1]}) \pi_h^{[r-1]}}. \end{aligned} \quad (25)$$

The conditional distribution of $((\mathbf{y}_i^{\text{mis}})^{[r]} \mid \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r]} = 1, \mathbf{c}_i)$ has already been detailed in Equation (22) and recalled here

$$\begin{aligned} & f(\mathbf{y}_i^{\text{mis}} \mid \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r]} = 1, \mathbf{c}_i; \theta^{[r-1]}) \\ &= \frac{\prod_{j, c_{ij}=1} f(c_{ij} = 1 \mid \mathbf{y}_i^{\text{mis}}, \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r]} = 1; \psi^{[r-1]}) f(\mathbf{y}_i^{\text{mis}} \mid \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r]} = 1; \lambda^{[r-1]})}{\int_{\mathcal{Y}_i^{\text{mis}}} \prod_{j, c_{ij}=1} f(c_{ij} = 1 \mid \mathbf{y}_i^{\text{mis}}, \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r]} = 1; \psi^{[r-1]}) f(\mathbf{y}_i^{\text{mis}} \mid \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r]} = 1; \lambda^{[r-1]}) d\mathbf{y}_i^{\text{mis}}}. \end{aligned} \quad (26)$$

Gaussian mixture for continuous data First note that the probabilities of the multinomial distribution for drawing $z_i^{[r]}$ given in (25) can be easily computed for all cases.

$$\begin{aligned} & \mathbb{P}(z_{ik} = 1 \mid \mathbf{y}_i^{[r-1]}, \mathbf{c}_i; \theta^{[r-1]}) \\ &= \frac{\prod_{j=1}^d f(c_{ij} = 1 \mid \mathbf{y}_i^{[r-1]}, z_{ik}^{[r-1]} = 1; \psi^{[r-1]})^{c_{ij}} f(c_{ij} = 0 \mid \mathbf{y}_i^{[r-1]}, z_{ik}^{[r-1]} = 1; \psi^{[r-1]})^{1-c_{ij}} \phi(\mathbf{y}_i^{[r-1]}; \lambda_k^{[r-1]}) \pi_k^{[r-1]}}{\sum_{h=1}^K \prod_{j=1}^d f(c_{ij} = 1 \mid \mathbf{y}_i^{[r-1]}, z_{ih}^{[r-1]} = 1; \psi^{[r-1]})^{c_{ij}} f(c_{ij} = 0 \mid \mathbf{y}_i^{[r-1]}, z_{ih}^{[r-1]} = 1; \psi^{[r-1]})^{1-c_{ij}} \phi(\mathbf{y}_i^{[r-1]}; \lambda_h^{[r-1]}) \pi_h^{[r-1]}}, \end{aligned}$$

where $\phi(\mathbf{y}_i; \lambda_k) = \phi(\mathbf{y}_i; \mu_k, \Sigma_k)$ is assumed to be a Gaussian distribution with mean vector μ_k and covariance matrix Σ_k , and $f(c_{ij} = 1 \mid \mathbf{y}_i^{[r-1]}, z_{ih}^{[r-1]} = 1; \psi^{[r-1]})$ is specified depending the MNAR model and the distribution ρ . The only difficulty of the SE-step is thus to draw from the distribution $(\mathbf{y}_i^{\text{mis}} \mid \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r]} = 1, \mathbf{c}_i)$.

In the sequel, we detail the distribution $(\mathbf{y}_i^{\text{mis}} \mid \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r]} = 1, \mathbf{c}_i)$ and the M-step for ψ depending the MNAR model.

For MNAR y_\star models, the conditional distribution $(y_i^{\text{mis}} \mid y_i^{\text{obs}}, z_{ik}^{[r]} = 1, c_i)$ depends on the distribution ρ at hand. For the MNAR y_\star models, we will consider two classical distributions for ρ : the logistic function and probit one.

Logistic distribution: For the logistic function, the distribution given in (26) is not classical and drawing y_i^{mis} from it seems complicated. Indeed, one has

$$\begin{aligned} & f(\mathbf{y}_i^{\text{mis}} \mid \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r]} = 1, \mathbf{c}_i; \theta^{[r-1]}) \\ & \propto \prod_{j=1, c_{ij}=1} \frac{1}{1 + \exp(\alpha_{kj}^{[r-1]} + \beta_{kj}^{[r-1]} y_{ij}^{\text{mis}})} \phi(\mathbf{y}_i^{\text{mis}}; (\tilde{\mu}_{ik}^{\text{mis}})^{[r-1]}, (\tilde{\Sigma}_{ik}^{\text{mis}})^{[r-1]}), \end{aligned}$$

where $(\tilde{\mu}_{ik}^{\text{mis}})^{[r-1]}$ and $(\tilde{\Sigma}_{ik}^{\text{mis}})^{[r-1]}$ are given in (33) and (34). We could use the Sampling Importance Resampling (SIR) algorithm which simulates a realization of $(\mathbf{y}_i^{\text{mis}} \mid \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r]} = 1, \mathbf{c}_i)$ with a known instrumental distribution (for example: $(\mathbf{y}_i^{\text{mis}} \mid \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r]} = 1)$) and includes a re-sampling step. However, this algorithm may be computationnaly costly.

Probit distribution: For the probit function, the distribution in (26) can be made explicit by using a latent variable \mathbf{L}_i .

More particularly, let \mathbf{L}_i such that $\mathbf{L}_i = \alpha_k^{[r-1]} + \beta_k^{[r-1]} \mathbf{y}_i + \epsilon_i$, with $\epsilon_i \sim \mathcal{N}(0_d, I_{d \times d})$. Then, \mathbf{c}_i can be viewed as an indicator for whether this latent variable is positive, *i.e.*, for all $j = 1, \dots, d$,

$$c_{ij} = \begin{cases} 1 & \text{if } L_{ij} > 0 \\ 0 & \text{otherwise.} \end{cases} \quad (27)$$

Thus, indeed to draw $(\mathbf{y}_i^{\text{mis}})^{[r]}$ and $\mathbf{z}_i^{[r]}$ according to $f(\mathbf{y}_i^{\text{mis}}, \mathbf{z}_i \mid \mathbf{y}_i^{\text{obs}}, \mathbf{c}_i; \theta^{[r-1]})$, we draw $\mathbf{L}_i^{[r]}$, $(\mathbf{y}_i^{\text{mis}})^{[r]}$ and $\mathbf{z}_i^{[r]}$ according to $f(\mathbf{L}_i, \mathbf{y}_i^{\text{mis}}, \mathbf{z}_i \mid \mathbf{y}_i^{\text{obs}}, \mathbf{c}_i; \theta^{[r-1]})$ by using a Gibbs sampling.

First, we have to draw $\mathbf{L}_i^{[r]}$ according to $f(\cdot \mid \mathbf{y}_i^{[r-1]}, z_{ik}^{[r-1]} = 1, \mathbf{c}_i; \psi^{[r-1]})$. One has

$$\begin{aligned} f(\mathbf{L}_i \mid \mathbf{y}_i^{[r-1]}, z_{ik}^{[r-1]} = 1, \mathbf{c}_i) & \propto f(\mathbf{L}_i, \mathbf{c}_i \mid \mathbf{y}_i^{[r-1]}, z_{ik}^{[r-1]} = 1; \psi^{[r-1]}) \\ & \propto f(\mathbf{c}_i \mid \mathbf{L}_i^{[r]}, \mathbf{y}_i^{[r-1]}, z_{ik}^{[r-1]} = 1; \psi^{[r-1]}) f(\mathbf{L}_i^{[r]} \mid \mathbf{y}_i^{\text{mis}}, \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r-1]} = 1; \psi^{[r-1]}) \\ & \stackrel{(i)}{\propto} f(\mathbf{c}_i \mid \mathbf{L}_i^{[r]}; \psi^{[r-1]}) f(\mathbf{L}_i^{[r]} \mid \mathbf{y}_i^{[r-1]}, z_{ik}^{[r-1]} = 1; \psi^{[r-1]}) \\ & \stackrel{(ii)}{=} \mathbb{1}_{\{\mathbf{c}_i = 1\} \cap \{\mathbf{L}_i^{[r]} > 0\}} f(\mathbf{L}_i^{[r]} \mid \mathbf{y}_i^{\text{mis}}, \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r-1]} = 1; \psi^{[r-1]}) \end{aligned}$$

where we use that $\mathbf{L}_i^{[r]}$ is a function of $\mathbf{y}_i^{\text{mis}}, \mathbf{y}_i^{\text{obs}}, z_{ik} = 1$ in step (i). Step (ii) is obtained by using (27). By abuse of notation, $\{\mathbf{c}_i = 1\} \cap \{\mathbf{L}_i^{[r]} > 0\}$ means that for all $j = 1, \dots, d$, $\{c_{ij} = 1\} \cap \{L_{ij}^{[r]} > 0\}$. Finally the conditional distribution $(\mathbf{L}_i \mid \mathbf{y}_i^{[r-1]}, z_{ik}^{[r-1]} = 1, \mathbf{c}_i)$ is a multivariate truncated Gaussian distribution denoted as \mathcal{N}_t , as detailed here

$$(\mathbf{L}_i \mid \mathbf{y}_i^{[r-1]}, z_{ik}^{[r-1]} = 1, \mathbf{c}_i) \sim \mathcal{N}_t(\alpha_k^{[r-1]} + \beta_k^{[r-1]} \mathbf{y}_i, I_{d \times d}; a, b), \quad (28)$$

with $a \in \mathbb{R}^d$ and $b \in \mathbb{R}^d$ the lower and upper bounds such that for all $j = 1, \dots, d$,

$$\begin{aligned} a_j &= \begin{cases} 0 & \text{if } c_{ij} = 1, \\ -\infty & \text{otherwise.} \end{cases} \\ b_j &= \begin{cases} +\infty & \text{if } c_{ij} = 1, \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

Secondly, we draw the membership k of $\mathbf{z}_i^{[r]}$ from the multinomial distribution with probabilities, for all $k = 1, \dots, K$ detailed as follows

$$\begin{aligned} \mathbb{P}(z_{ik} = 1 \mid \mathbf{L}_i^{[r]}, \mathbf{y}_i^{[r-1]}, \mathbf{c}_i; \theta^{[r-1]}) &= \frac{\mathbb{P}(z_{ik} = 1, \mathbf{L}_i^{[r]}, \mathbf{y}_i^{[r-1]}, \mathbf{c}_i; \theta^{[r-1]})}{\sum_{k=1}^K \mathbb{P}(z_{ik} = 1, \mathbf{L}_i^{[r]}, \mathbf{y}_i^{[r-1]}, \mathbf{c}_i; \theta^{[r-1]})} \\ &= \frac{f(\mathbf{L}_i^{[r]} \mid z_{ik} = 1, \mathbf{y}_i^{[r-1]}, \mathbf{c}_i; \psi^{[r-1]}) f(z_{ik} = 1, \mathbf{y}_i^{[r-1]}, \mathbf{c}_i; \theta^{[r-1]})}{\sum_{k=1}^K f(\mathbf{L}_i^{[r]} \mid z_{ik} = 1, \mathbf{y}_i^{[r-1]}, \mathbf{c}_i; \psi^{[r-1]}) f(z_{ik} = 1, \mathbf{y}_i^{[r-1]}, \mathbf{c}_i; \theta^{[r-1]})}. \end{aligned} \quad (29)$$

The part involving $f(z_{ik} = 1, \mathbf{y}_i^{[r-1]}, \mathbf{c}_i; \theta^{[r-1]})$ is given in (25) and $f(\mathbf{L}_i^{[r]} \mid z_{ik} = 1, \mathbf{y}_i^{[r-1]}, \mathbf{c}_i; \psi^{[r-1]})$ is only the density of the multivariate truncated Gaussian distribution described in (28) evaluated in $\mathbf{L}_i^{[r]}$.

Finally, $\mathbf{y}_i^{[r]}$ is drawn according to $f(\cdot | \mathbf{L}_i^{[r]}, z_{ik}^{[r]} = 1, \mathbf{y}_i^{\text{obs}}, \mathbf{c}_i; \theta^{[r-1]})$. One has

$$\begin{aligned}
& f(\mathbf{y}_i^{\text{mis}} | \mathbf{L}_i^{[r]}, z_{ik}^{[r]} = 1, \mathbf{y}_i^{\text{obs}}, \mathbf{c}_i; \theta^{[r-1]}) \\
& \propto f(\mathbf{c}_i, \mathbf{L}_i^{[r]} | \mathbf{y}_i^{\text{mis}}, \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r]} = 1; \psi^{[r-1]}) f(\mathbf{y}_i^{\text{mis}} | \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r]} = 1; \theta^{[r-1]}) \\
& \propto f(\mathbf{c}_i | \mathbf{L}_i^{[r]}, \mathbf{y}_i^{\text{mis}}, \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r]} = 1; \psi^{[r-1]}) f(\mathbf{L}_i^{[r]} | \mathbf{y}_i^{\text{mis}}, \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r]} = 1; \psi^{[r-1]}) f(\mathbf{y}_i^{\text{mis}} | \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r]} = 1; \theta^{[r-1]}) \\
& \propto f(\mathbf{c}_i | \mathbf{L}_i^{[r]}; \psi^{[r-1]}) f(\mathbf{L}_i^{[r]} | \mathbf{y}_i^{\text{mis}}, \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r]} = 1; \psi^{[r-1]}) f(\mathbf{y}_i^{\text{mis}} | \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r]} = 1; \theta^{[r-1]}) \\
& \propto f(\mathbf{L}_i^{[r]} | \mathbf{y}_i^{\text{mis}}, \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r]} = 1; \psi^{[r-1]}) f(\mathbf{y}_i^{\text{mis}} | \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r]} = 1; \theta^{[r-1]}).
\end{aligned}$$

Yet, one has

$$f(\mathbf{L}_i^{[r]} | \mathbf{y}_i^{\text{mis}}, \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r]} = 1; \psi^{[r-1]}) \propto \exp \left(-\frac{1}{2} \left[(\mathbf{L}_i^{[r]} - (\alpha_k^{[r-1]} + \beta_k^{[r-1]} \mathbf{y}_i))^T (\mathbf{L}_i^{[r]} - (\alpha_k^{[r-1]} + \beta_k^{[r-1]} \mathbf{y}_i)) \right] \right) \quad (30)$$

$$f(\mathbf{y}_i^{\text{mis}} | \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r]} = 1; \theta^{[r-1]}) \propto \exp \left(-\frac{1}{2} \left[(\mathbf{y}_i^{\text{mis}} - (\tilde{\mu}_{ik}^{\text{mis}})^{[r]})^T ((\tilde{\Sigma}_{ik}^{\text{mis}})^{[r-1]})^{-1} (\mathbf{y}_i^{\text{mis}} - (\tilde{\mu}_{ik}^{\text{mis}})^{[r]}) \right] \right), \quad (31)$$

with $(\tilde{\mu}_{ik}^{\text{mis}})^{[r-1]}$ and $(\tilde{\Sigma}_{ik}^{\text{mis}})^{[r-1]}$ the standard expression of the mean vector and covariance matrix of a conditional Gaussian distribution (see for instance [Anderson \(2003\)](#)). In particular, one has in this case:

$$(\mathbf{y}_i^{\text{mis}} | \mathbf{y}_i^{\text{obs}}, z_{ik} = 1; \lambda^{[r-1]}) \sim \mathcal{N} \left((\tilde{\mu}_{ik}^{\text{mis}})^{[r-1]}, (\tilde{\Sigma}_{ik}^{\text{mis}})^{[r-1]} \right). \quad (32)$$

with:

$$(\tilde{\mu}_{ik}^{\text{mis}})^{[r-1]} = (\mu_{ik}^{\text{mis}})^{[r-1]} + (\Sigma_{ik}^{\text{mis,obs}})^{[r-1]} \left((\Sigma_{ik}^{\text{obs,obs}})^{[r-1]} \right)^{-1} \left(\mathbf{y}_i^{\text{obs}} - (\mu_{ik}^{\text{obs}})^{[r-1]} \right), \quad (33)$$

$$(\tilde{\Sigma}_{ik}^{\text{mis}})^{[r-1]} = (\Sigma_{ik}^{\text{mis,mis}})^{[r-1]} - (\Sigma_{ik}^{\text{mis,obs}})^{[r-1]} \left((\Sigma_{ik}^{\text{obs,obs}})^{[r-1]} \right)^{-1} (\Sigma_{ik}^{\text{obs,mis}})^{[r-1]}. \quad (34)$$

Finally combining the two equations (30) and (31) one obtains

$$(\mathbf{y}_i^{\text{mis}} | \mathbf{L}_i^{[r]}, z_{ik}^{[r]} = 1, \mathbf{y}_i^{\text{obs}}, \mathbf{c}_i) \sim \mathcal{N} (\mu_{ik}^{\text{SEM}}, \Sigma_{ik}^{\text{SEM}}), \quad (35)$$

where

$$\Sigma_{ik}^{\text{SEM}} = \left(((\tilde{\Sigma}_{ik}^{\text{mis}})^{[r-1]})^{-1} + ((\beta_k^{\text{mis}})^{[r-1]})^T (\beta_k^{\text{mis}})^{[r-1]} \right)^{-1},$$

$$\mu_{ik}^{\text{SEM}} = \Sigma_{ik}^{\text{SEM}} \left[((\tilde{\Sigma}_{ik}^{\text{mis}})^{[r-1]})^{-1} \tilde{\mu}_{ik}^{\text{mis}} + ((\beta_k^{\text{mis}})^{[r-1]})^T (\mathbf{L}_i^{\text{mis}})^{[r]} - ((\beta_k^{\text{mis}})^{[r-1]})^T (\alpha_k^{\text{mis}})^{[r-1]} \right],$$

with $(\beta_k^{\text{mis}})^{[r-1]}$ (resp. $(\mathbf{L}_i^{\text{mis}})^{[r]}$ and $(\alpha_k^{\text{mis}})^{[r-1]}$) the vector β_k (resp. $(\mathbf{L}_i)^{[r]}$ and $(\alpha_k)^{[r-1]}$) restricted to the coordinates $j \in \mathcal{Y}_i^{\text{mis}}$.

Finally, for fully describing the SEM-algorithm, in the M-step, $\psi^{[r-1]}$ is computed using a GLM with a binomial link function for a matrix depending on the MNAR model. In particular,

- For MNAR_y, the coefficient obtained with a GLM for the matrix $(\mathcal{H}_j^{\text{MNAR}_y})^{[r]}$ are α_0 and $\beta_1^{[r]}, \dots, \beta_d^{[r]}$, with

$$(\mathcal{H}^{\text{MNAR}_y})^{[r]} = \begin{array}{c|ccccc} c.1 & 1 & y_{.1}^{[r]} & 0 & \dots & 0 \\ c.2 & 1 & 0 & y_{.2}^{[r]} & \dots & 0 \\ \vdots & \vdots & & \ddots & \ddots & \\ c.d & 1 & 0 & 0 & \dots & y_{.d}^{[r]} \end{array}. \quad (36)$$

- For $\text{MNAR}y^k$, the coefficient obtained with a GLM for the matrix $(\mathcal{H}_k^{\text{MNAR}y^k})^{[r]}$ is α_0 and $\beta_{11}^{[r]}, \dots, \beta_{K1}^{[r]}, \dots, \beta_{Kd}^{[r]}$ with

$$(\mathcal{H}_k^{\text{MNAR}y^k})^{[r]} = \begin{pmatrix} (c_{u1})_{u,z_{u1}^{[r]}=1} \\ \vdots \\ (c_{u1})_{u,z_{uK}^{[r]}=1} \\ \vdots \\ (c_{ud})_{u,z_{uK}^{[r]}=1} \end{pmatrix} \begin{vmatrix} 1 & (y_{u1}^{[r]})_{u,z_{u1}^{[r]}=1} & 0 & \dots & 0 \\ \vdots & \vdots & & \ddots & \vdots \\ 1 & 0 & (y_{u1}^{[r]})_{u,z_{uK}^{[r]}=1} & & 0 \\ \vdots & \vdots & & \ddots & \\ 1 & 0 & 0 & & (y_{ud}^{[r]})_{u,z_{uK}^{[r]}=1} \end{vmatrix}. \quad (37)$$

- For $\text{MNAR}yz$, the coefficients obtained with a GLM for the matrix $(\mathcal{H}^{\text{MNAR}yz})^{[r]}$ are $\beta_1^{[r]}, \dots, \beta_d^{[r]}$ and $\alpha_1^{[r]}, \dots, \alpha_K^{[r]}$, with

$$(\mathcal{H}^{\text{MNAR}yz})^{[r]} = \begin{pmatrix} c.1 \\ c.2 \\ \vdots \\ c.d \end{pmatrix} \begin{vmatrix} y_{.1}^{[r]} & 0 & \dots & 0 & z_{.1}^{[r]} & \dots & z_{.K}^{[r]} \\ 0 & y_{.2}^{[r]} & \dots & 0 & z_{.1}^{[r]} & \dots & z_{.K}^{[r]} \\ & \ddots & \ddots & & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & y_{.d}^{[r]} & z_{.1}^{[r]} & \dots & z_{.K}^{[r]} \end{vmatrix}. \quad (38)$$

- For $\text{MNAR}yz^j$, the coefficients obtained with a GLM for the matrix $(\mathcal{H}_j^{\text{MNAR}yz^j})^{[r]}$ are $\beta_j^{[r]}, \alpha_{1j}^{[r]}, \dots, \alpha_{Kj}^{[r]}$, with

$$(\mathcal{H}_j^{\text{MNAR}yz^j})^{[r]} = \begin{pmatrix} c.j \end{pmatrix} \begin{vmatrix} y_{.j}^{[r]} & z_{.1}^{[r]} & \dots & z_{.K}^{[r]} \end{vmatrix}. \quad (39)$$

- For $\text{MNAR}y^k z$, the coefficients obtained with a GLM for the matrix $(\mathcal{H}_k^{\text{MNAR}y^k z})^{[r]}$ are $\beta_{k1}^{[r]}, \dots, \beta_{kd}^{[r]}, \alpha_k^{[r]}$, with

$$(\mathcal{H}_k^{\text{MNAR}y^k z})^{[r]} = \begin{pmatrix} (c_{u1})_{u,z_{uk}^{[r]}=1} \\ (c_{u2})_{u,z_{uk}^{[r]}=1} \\ \vdots \\ (c_{ud})_{u,z_{uk}^{[r]}=1} \end{pmatrix} \begin{vmatrix} (y_{u1}^{[r]})_{u,z_{uk}^{[r]}=1} & 0 & \dots & 0 & 1 \\ 0 & (y_{u2}^{[r]})_{u,z_{uk}^{[r]}=1} & \dots & 0 & 1 \\ & \ddots & \ddots & & 1 \\ 0 & 0 & \dots & (y_{ud}^{[r]})_{u,z_{uk}^{[r]}=1} & 1 \end{vmatrix}. \quad (40)$$

- For $\text{MNAR}y^k z^j$, the coefficients obtained with a GLM for the matrix $(\mathcal{H}_{kj}^{\text{MNAR}y^k z^j})^{[r]}$ are $\beta_{kj}^{[r]}, \alpha_{kj}^{[r]}$, with

$$(\mathcal{H}_{kj}^{\text{MNAR}y^k z^j})^{[r]} = \begin{pmatrix} (c_{uj})_{u,z_{uk}^{[r]}=1} \end{pmatrix} \begin{vmatrix} (y_{uj}^{[r]})_{u,z_{uk}^{[r]}=1} & 1 \end{vmatrix} \quad (41)$$

When ρ is the probit function, the SEM algorithm can be derived, see Algorithm 1.

Remark C.1 (SEM for $\text{MNAR}z$ and $\text{MNAR}z^j$ mechanisms). *A SEM algorithm can also be derived for these two mechanisms. For continuous data, we can prove that*

$$f(\mathbf{y}_i^{\text{mis}} \mid \mathbf{y}_i^{\text{obs}}, \mathbf{z}_i^{[r-1]}, \mathbf{c}_i; \theta^{[r-1]}) = f(\mathbf{y}_i^{\text{mis}} \mid \mathbf{y}_i^{\text{obs}}, \mathbf{z}_i^{[r-1]}; \lambda^{[r-1]}),$$

and that this conditional distribution is Gaussian (the SE step is then just a draw from this law). The M-step for ψ are given in the following.

Algorithm 1 SEM algorithm for Gaussian mixture, MNAR y_{\star} models, ρ is probit

Input: $Y \in \mathbb{R}^{n \times d}$ (matrix containing missing values), $K \geq 1, r_{\max}$.

Initialize $Z^0, \pi_k^0, \mu_k^0, \Sigma_k^0$ and ψ_k^0 , for $k \in \{1, \dots, K\}$.

for $r = 0$ **to** r_{\max} **do**

SE-step:

for $i = 1$ **to** n **do**

Draw $(\mathbf{L}_i)^{[r]}$ from the multivariate truncated Gaussian distribution given in (28).

Draw $\mathbf{z}_i^{[r]}$ from the multinomial distribution with probabilities detailed in (29).

Draw $(\mathbf{y}_i^{\text{mis}})^{[r]}$ from the multivariate Gaussian distribution given in (35).

end for

Let $Y^{[r]} = (\mathbf{y}_1^{[r]} | \dots | \mathbf{y}_n^{[r]})$ be the imputed matrix.

Let $Z^{[r]} = (\mathbf{z}_1^{[r]} | \dots | \mathbf{z}_n^{[r]})$ be the partition.

M-step:

for $k = 1$ **to** K **do**

Let $\pi_k^{[r]}$ be the proportion of rows of $Y^{[r]}$ belonging class k .

Let $\mu_k^{[r]}, \Sigma_k^{[r]}$ be the mean and covariance matrix of rows of $Y^{[r]}$ belonging to class k.

Let $\psi_k^{[r]}$ be the resulted coefficients of a GLM with a binomial link function, *i.e.*, the optimization problem is $\forall j \in \{1, \dots, d\}$,

$$\mathcal{M}_{kj}\psi_k^{[r]} = \log\left(\frac{1 - \mathbb{E}[\mathbf{c}_{\cdot j}|\mathcal{M}_{kj}]}{\mathbb{E}[\mathbf{c}_{\cdot j}|\mathcal{M}_{kj}]}\right),$$

for a matrix \mathcal{M}_{kj} depending on the MNAR model (see (36), (37), (38), (43), (41) and (42)) and $\mathbf{c}_{\cdot j}$ the missing data pattern for the variable j .

end for**end for**

- For $MNARz$, the coefficients obtained with a GLM for the matrix $(\mathcal{H}^{MNARz})^{[r]}$ are $\alpha_1, \dots, \alpha_K$, with

$$(\mathcal{H}^{MNARz})^{[r]} = \begin{array}{c|ccc} c_{.1} & z_{.1} & \dots & z_{.K} \\ \vdots & \vdots & \vdots & \vdots \\ c_{.d} & z_{.1} & \dots & z_{.K} \end{array} = \begin{array}{c|ccc} c_{11} & z_{11}^{[r]} & \dots & z_{1K}^{[r]} \\ \vdots & \vdots & \vdots & \vdots \\ c_{n1} & z_{n1}^{[r]} & \dots & z_{nK}^{[r]} \\ \vdots & \vdots & \vdots & \vdots \\ c_{1d} & z_{11}^{[r]} & \dots & z_{1K}^{[r]} \\ \vdots & \vdots & \vdots & \vdots \\ c_{nd} & z_{n1}^{[r]} & \dots & z_{nK}^{[r]} \end{array}. \quad (42)$$

- For $MNARz^j$, the coefficients obtained with a GLM for the matrix $(\mathcal{H}_j^{MNARz^j})^{[r]}$ are $\alpha_{1j}, \dots, \alpha_{Kj}$, with

$$(\mathcal{H}_j^{MNARz^j})^{[r]} = \begin{array}{c|ccc} c_{.j} & z_{.1}^{[r]} & \dots & z_{.K}^{[r]} \end{array} \quad (43)$$

For categorical data, we have $\phi(\mathbf{y}_i; \lambda_k) = \prod_{j=1}^d \phi(y_{ij}; \lambda_{kj}) = \prod_{j=1}^d \prod_{\ell=1}^{\ell_j} (\lambda_{kj}^\ell)^{y_{ij}^\ell}$. For drawing from the conditional distribution $(\mathbf{y}_i^{\text{mis}} \mid \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r]} = 1)$, by independence of the features conditionally to the membership, we can draw for $j = 1, \dots, d$ $y_{ij}^{\text{mis}} = ((y_{ij}^{\text{mis}})^1, \dots, (y_{ij}^{\text{mis}})^{\ell_j})$ from the conditional distribution $(y_{ij}^{\text{mis}} \mid \mathbf{y}_i^{\text{obs}}, z_{ik}^{[r]} = 1)$. This latter is a multinomial distribution with probabilities $(\lambda_{kj}^\ell)_{\ell=1, \dots, \ell_j}$.

D Appendix 4: Complements on the numerical experiment

This section gives the values of δ (see (14)) ψ (see (3)) and φ (see (14)) used during the numerical experiments. For most MNAR models, the values of these parameters are available in the main paper (Sportisse et al., 2023).

| d | δ | α | β |
|-----|----------|----------|--|
| 6 | 1.92 | -0.75 | $\begin{pmatrix} -3 & 0.3 & -3 & -3 & -2 & 1 \\ 0.5 & -2 & 1 & 1 & 1 & 0.5 \\ 1 & 1 & 0.5 & 0.5 & 0.5 & 2 \end{pmatrix}$ |

Table 2: Choice of the values of δ , α and β for all the experiments of Section 5 for the $MNARy^k$ mechanism.

| d | δ | α | β |
|-----|----------|------------------------------|--|
| 6 | 1.91 | $(-0.9 \quad -0.15 \quad 0)$ | $\begin{pmatrix} -3 & 0.3 & -3 & -3 & -2 & 1 \\ 0.3 & -3 & 0.3 & -0.3 & -2 & 0.2 \\ -3 & 0.3 & -3 & -3 & -2 & 1 \end{pmatrix}$ |

Table 3: Choice of the values of δ , α and β for all the experiments of Section 5 for the $MNARy^k z$ mechanism.

| d | δ | α | β |
|-----|----------|--|---------------------------------|
| 6 | 2.15 | $\begin{pmatrix} -1.4 & -1.4 & -1.2 & -1.1 & -1 & -0.9 \\ -0.6 & 0.4 & 0.4 & 0.3 & 0.1 & 0.1 \\ -0.8 & -0.8 & 0.8 & -0.8 & -0.8 & 0.8 \end{pmatrix}$ | $(-3 \ 0.3 \ -3 \ -3 \ -2 \ 1)$ |

Table 4: Choice of the values of δ , α and β for all the experiments of Section 5 for the MNAR_{yz^j} mechanism.