



Dealing with missing data in model-based clustering through a MNAR model

Christophe Biernacki, Gilles Celeux, Julie Josse, Fabien Laporte

► To cite this version:

Christophe Biernacki, Gilles Celeux, Julie Josse, Fabien Laporte. Dealing with missing data in model-based clustering through a MNAR model. CRoNos & MDA 2019 - Meeting and Workshop on Multivariate Data Analysis and Software, Apr 2019, Limassol, Cyprus. hal-02103347

HAL Id: hal-02103347

<https://inria.hal.science/hal-02103347>

Submitted on 18 Apr 2019

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.

Dealing with missing data in model-based clustering through a MNAR model

Christophe Biernacki, Gilles Celeux, Julie Josse, Fabien Laporte

Final CRoNoS meeting and Workshop on Multivariate Data Analysis and Software
14-16 April 2019, Limassol, Cyprus



Take home message

- 1 The missing data **pattern** may convey some information on clustering
- 2 **Embed the missingness mechanism** directly within the clustering modeling step

Outline

1 Introduction

2 A model-based MNAR clustering approach

3 Inference procedures

4 Medical study illustration

5 Concluding remarks

Missing data: an inevitable event

The larger the datasets, the more missing data may appear. . .

Two traditional solutions (for obtaining a filled dataset)

- **Discard** individuals with missing data: expect to add **variance** into analysis
- **Impute** missing data: expect to add **bias** modeling into analysis

General guidelines

- Obtaining a non-missing dataset is **not** the final goal
- Missing data management should **take into account the initial analysis target**

Our analysis target: model-based clustering

Embed missing data management into this paradigm. . .

Missing data: notations

- $\mathbf{y} = \{\mathbf{y}_1, \dots, \mathbf{y}_n\}$: full dataset with n individuals
- $\mathbf{y}_i = (y_i^1, \dots, y_i^d) \in \mathbb{R}^d$: full individual $i \in \{1, \dots, n\}$
- $\mathbf{c} = \{\mathbf{c}_1, \dots, \mathbf{c}_n\}$: pattern of missing data for the full dataset
- $\mathbf{c}_i = (c_i^1, \dots, c_i^d) \in \{0, 1\}^d$: pattern of missing data for individual $i \in \{1, \dots, n\}$

$$c_i^j = 1 \Leftrightarrow y_i^j \text{ is missing}$$

- $\mathbf{o}_i = \{j : c_i^j = 0\}$: the observed variables indexes for individual i
- $\mathbf{y}_i^{\mathbf{o}_i}$: the observed variables values for individual i
- $\mathbf{y}^{\mathbf{o}} = \{\mathbf{y}_1^{\mathbf{o}_1}, \dots, \mathbf{y}_n^{\mathbf{o}_n}\}$: the observed values in \mathbf{y}
- $\mathbf{m}_i = \{j : c_i^j = 1\}$: the missing variables indexes for individual i
- $\mathbf{y}_i^{\mathbf{m}_i}$: the missing variables values for individual i
- $\mathbf{y}^{\mathbf{m}} = \{\mathbf{y}_1^{\mathbf{m}_1}, \dots, \mathbf{y}_n^{\mathbf{m}_n}\}$: the missing values in \mathbf{y}

$\mathbf{y} = \{\mathbf{y}^{\mathbf{o}}, \mathbf{y}^{\mathbf{m}}\}$ is the full dataset with its observed and missing parts

Missing data: typology of the missing mechanisms

- Missing completely at random (MCAR):

$$P(\mathbf{c}|\mathbf{y}; \psi) = P(\mathbf{c}; \psi) \quad \forall \mathbf{y}$$

- Missing at random (MAR):

$$P(\mathbf{c}|\mathbf{y}; \psi) = P(\mathbf{c}|\mathbf{y}^o; \psi) \quad \forall \mathbf{y}^m$$

- Missing not at random (MNAR): the mechanism is not MCAR nor MAR

Clustering: model-based approach

- **Partition with K clusters:** $\mathbf{z} = (\mathbf{z}_1, \dots, \mathbf{z}_n)$ where
 - $\mathbf{z}_i = (z_i^1, \dots, z_i^K) \in \{0, 1\}^K$
 - $z_i^k = 1$ if \mathbf{y}_i belongs to cluster k , $z_i^k = 0$ otherwise
- **Gaussian mixture:** $\mathbf{y}_1, \dots, \mathbf{y}_n$ are i.i.d. from the mixture

$$f(\mathbf{y}_i; \boldsymbol{\theta}) = \sum_{k=1}^K \pi_k \phi_k(\mathbf{y}_i; \theta_k)$$

where:

- $\pi_k = P(z_i^k = 1)$
- $\phi_k(\cdot; \theta_k)$: d -variate Gaussian pdf with mean vector and covariance matrix
 $\theta_k = (\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$, d -multinomial pdf with $\theta_k = \mathbf{p}_k$ probabilities vector or mixed pdf.
- $\boldsymbol{\theta} = (\pi_1, \dots, \pi_K, \theta_1, \dots, \theta_K)$

Question we address in this work

Which distribution $P(\mathbf{c}|\mathbf{y}, \mathbf{z}; \boldsymbol{\psi})$ to propose in this clustering context?

Outline

- 1 Introduction
- 2 A model-based MNAR clustering approach**
- 3 Inference procedures
- 4 Medical study illustration
- 5 Concluding remarks

Logistic model: a natural and flexible candidate

$$P(\mathbf{c}|\mathbf{y}, \mathbf{z}; \boldsymbol{\psi}) = \prod_{i=1}^n \prod_{j=1}^d P(c_i^j | \mathbf{y}, \mathbf{z}; \boldsymbol{\psi})$$

- MCAR, with $\boldsymbol{\psi} = \alpha_0$

$$\text{logit}(P(c_i^j = 1 | \mathbf{y}, \mathbf{z}; \boldsymbol{\psi})) = \alpha_0$$

- MNARz (MNARz^j), with $\boldsymbol{\psi} = (\alpha_0, \beta_1^{1\dots d}, \dots, \beta_K^{1\dots d})$

$$\text{logit}(P(c_i^j = 1 | \mathbf{y}, \mathbf{z}; \boldsymbol{\psi})) = \alpha_0 + \sum_{k=1}^K \beta_k^j z_i^k$$

- MNARY, with $\boldsymbol{\psi} = (\alpha_0, \alpha_1, \dots, \alpha_d)$

$$\text{logit}(P(c_i^j = 1 | \mathbf{y}, \mathbf{z}; \boldsymbol{\psi})) = \alpha_0 + \alpha_j y_i^j$$

- MNARYz, with $\boldsymbol{\psi} = (\alpha_0, \alpha_1, \dots, \alpha_d, \beta_1, \dots, \beta_K)$

$$\text{logit}(P(c_i^j = 1 | \mathbf{y}, \mathbf{z}; \boldsymbol{\psi})) = \alpha_0 + \alpha_j y_i^j + \sum_{k=1}^K \beta_k z_i^k$$

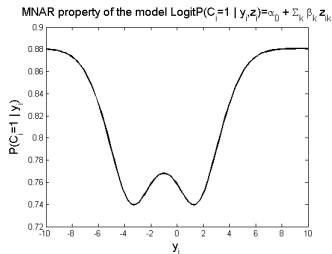
MNARz analysis: it depends on \mathbf{y} through \mathbf{z} !

$$P(c_i^j = 1 | \mathbf{y}; \boldsymbol{\theta}, \boldsymbol{\psi}) = \sum_{k=1}^K P(c_i^j = 1 | \mathbf{y}, \mathbf{z}; \boldsymbol{\psi}) P(\mathbf{z} | \mathbf{y}; \boldsymbol{\theta})$$

Example of a univariate Gaussian model with the three components

$$0.2N(\cdot; 0, 1) + 0.3N(\cdot; 1, 2) + 0.5N(\cdot; 2, 3)$$

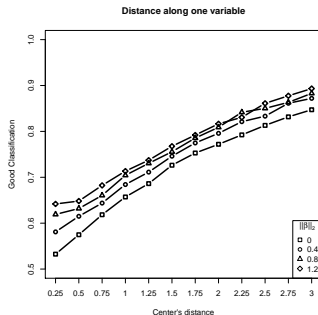
and with parameters of the logit expression: $\alpha_0 = 1, \beta_1 = 1, \beta_2 = -1, \beta_3 = 1$



MNARz analysis: pattern **c** gives information on partition **z**!

Draw Bayes error of a MNARz model with two components and 20% of missing data

$$\pi_k = 0.5, \|\mu_2 - \mu_1\| \text{ varies}, \Sigma_1 = \Sigma_2 = \mathbf{I}, |\beta_2 - \beta_1| \text{ varies}$$



Both μ_k and β_k act on the Bayes error

Outline

- 1 Introduction
- 2 A model-based MNAR clustering approach
- 3 Inference procedures**
- 4 Medical study illustration
- 5 Concluding remarks

Ignorable vs. non ignorable model

A missing mechanism is ignorable if likelihoods can be decomposed as

$$L(\theta, \psi; \underbrace{\mathbf{y}^o, \mathbf{c}}_{\text{observed data}}) = L(\psi; \mathbf{c} | \mathbf{y}^o) \times L(\theta; \mathbf{y}^o)$$

Some simple algebra show that this occurs when missing mechanism is not MNAR

Inference of θ

"If the missing mechanism is **ignorable** then likelihood-based inferences for θ from $L(\theta; \mathbf{y}^o)$ will be the same as likelihood based inference for θ from $L(\theta, \psi; \mathbf{y}^o, \mathbf{c})$." ([Little and Rubin, 2002] Section 6.2)

- MCAR is ignorable
- MNARz, MNARy and MNARyz are non ignorable

EM algorithm: looks simple

Decomposition of $Q(\theta, \psi; \hat{\theta}, \hat{\psi})$

The expected complete log-likelihood conditional related to observed data is:

$$E \left[L(\theta, \psi; \mathbf{c}, \mathbf{y}, \mathbf{z}); \hat{\theta}, \hat{\psi} | \mathbf{y}^o, \mathbf{c} \right] = Q_y(\theta; \hat{\theta}, \hat{\psi}) + Q_c(\psi; \hat{\theta}, \hat{\psi})$$

$$Q_y(\theta; \hat{\theta}, \hat{\psi}) = \sum_{i=1}^n \sum_{k=1}^K \tau_i^k E \left[\log(\pi_k \phi_k(\mathbf{y}_i; \theta_k)) | \mathbf{y}_i^{oi}, \mathbf{c}_i; \hat{\theta}, \hat{\psi} \right]$$

$$Q_c(\psi; \hat{\theta}, \hat{\psi}) = \sum_{i=1}^n \sum_{k=1}^K \tau_i^k E \left[\log(P(\mathbf{c}_i | z_i^k = 1, \mathbf{y}_i; \theta, \psi)) | \mathbf{y}_i^{oi}, \mathbf{c}_i; \hat{\theta}, \hat{\psi} \right]$$

$$\tau_i^k = P(z_i^k = 1 | \mathbf{c}_i, \mathbf{y}_i^{oi}; \hat{\theta}, \hat{\psi}) = \frac{\hat{\pi}_k \phi_k(\mathbf{y}_i^{oi}; \theta_k^{oi}) P(\mathbf{c}_i | z_i^k = 1, \mathbf{y}_i^{oi}; \hat{\psi})}{\sum_{h=1}^K \hat{\pi}_h \phi_h(\mathbf{y}_i^{oi}) P(\mathbf{c}_i | z_i^h = 1, \mathbf{y}_i^{oi}; \hat{\psi})}$$

EM and/or SEM algorithms

- **MCAR** (and also **MAR**...): classical formula! (EM , SEM)

$$\tau_i^k \propto \hat{\pi}_k \phi_k(\mathbf{y}_i^{oi}; \theta_k^{oi})$$

- **MNARz**: needs some new calculus but still simple (EM , SEM)

$$\tau_i^k \propto \hat{\pi}_k \phi_k(\mathbf{y}_i^{oi}; \theta_k^{oi}) \prod_{j=1}^d (1 + \exp(-r_i^j \hat{\beta}_k))^{-1} \text{ where } r_i^j = \begin{cases} 1 & \text{if } c_i^j = 1 \\ -1 & \text{otherwise} \end{cases}$$

- **MNARy**: needs approximations (EM , SEM)

$$P(c_i^j | \mathbf{y}_i^{oi}, z_i^k = 1; \psi) = \begin{cases} \int_{-\infty}^{+\infty} \frac{1}{1 + \exp(-(\alpha_j y_i^j))} \phi_k(y_i^j; \theta_k^j) dy_i^j & \text{if } c_i^j = 1 \\ \frac{1}{1 + \exp(\alpha_j y_i^j)} & \text{otherwise} \end{cases}$$

- In the Gaussian case, there is **no closed form** [Pirjol, 2013] (same for MNARyz)
- But **SEM is still simple** in that case thanks to random drawing instead of expectation

Link with some usual procedures!

Concatenation [Jones, 1996]: model equivalence

$$\text{MNARz}^j(\mathbf{y}^o, \mathbf{c}) \iff \text{MCAR}(\mathbf{y}^o | \mathbf{c})$$

$$\text{MNARz}(\mathbf{y}^o, \mathbf{c}) \iff \text{MCAR} \left(\mathbf{y}^o \middle| \left(\sum_{j=1}^d \mathbf{c}^j \right) \right)$$

“All Available Cases” [Little and Rubin, 2002]: estimation equivalence

In case of **conditional independence** between variables, whatever MCAR or MNAR*:

$$\text{Classical (S)EM} \iff (\text{S)EM without estimating missing } \mathbf{y}^m$$

... an opportunity to reduce the computing time

What about model selection?

Can select between MCAR and MNAR* with any information criterion (BIC, ICL)

Even if the missing mechanism is ignorable for MCAR...

... need to model \mathbf{c} to compare a MCAR and a MNAR model

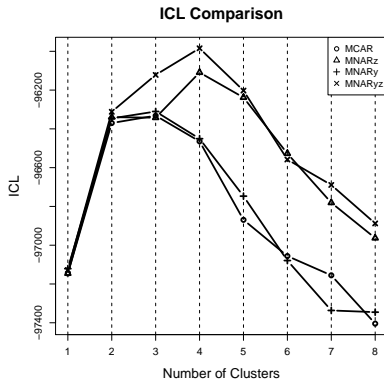
Outline

- 1 Introduction
- 2 A model-based MNAR clustering approach
- 3 Inference procedures
- 4 Medical study illustration**
- 5 Concluding remarks

Hospital Data

- Number of patients: $n = 5\,146$
- Number of features: $d = 7$
 - Age
 - Size
 - Weight
 - Cardiac frequency
 - Hemoglobin concentration
 - Temperature
 - Minimum Diastolic and Systolic Blood Pressure
- Percentage of missing data: 6.4%

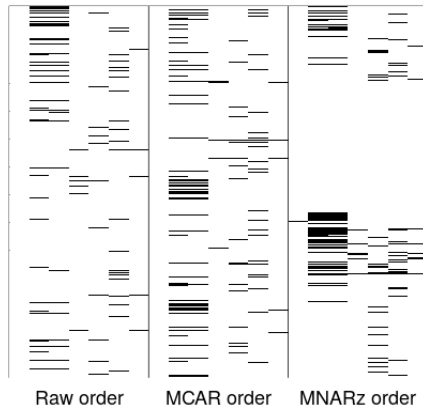
ICL comparison



- MCAR, MNARy and MNARz are equivalent until $K = 3$
- MNARz and MNARyz clearly indicate presence of an additional cluster ($K = 4$)

It seems to be an illustration of the effect of c through MNARz and MNARyz

Missing Pattern



It seems that MNARz modelling leads to a missing free cluster

Outline

- 1 Introduction
- 2 A model-based MNAR clustering approach
- 3 Inference procedures
- 4 Medical study illustration
- 5 Concluding remarks**

Summary

- Interest to put a model on c
- Interest of the simple but meaningful model MNAR $_z$
- Link between our models and usual methods

Ongoing works

- Deeper analysis of the previous results with doctors. . .
- Implement the proposed models/algo. in the Mixmod software^a
- Use **mixed data** algorithms for medical study with **the same MNAR* models**

^a<http://www.mixmod.org>

References



Jones, M. P. (1996).

Indicator and stratification methods for missing explanatory variables in multiple linear regression.

Journal of the American statistical association, 91(433):222–230.



Little, R. J. and Rubin, D. B. (2002).

Statistical Analysis with Missing Data.

Wiley.



Pirjol, D. (2013).

The logistic-normal integral and its generalizations.

Journal of Computational and Applied Mathematics, 237(1):460–469.