

Model-based clustering with missing not at random data. Missing mechanism

Fabien Laporte, Christophe Biernacki, Gilles Celeux, Julie Josse

► To cite this version:

Fabien Laporte, Christophe Biernacki, Gilles Celeux, Julie Josse. Model-based clustering with missing not at random data. Missing mechanism. Working Group on Model-Based Clustering Summer Session, Jul 2019, Vienne, Austria. hal-02398987

HAL Id: hal-02398987 https://hal.science/hal-02398987v1

Submitted on 8 Dec 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.

Model-based clustering with missing not at random data.



Fabien Laporte^{1,2}, Christophe Biernacki², Gilles Celeux,³ and Julie Josse,^{1,4}
 ¹CMAP, UMR 0320 / UMR 87641 Centre de Mathématiques Appliquées, 91128 Palaiseau, France
 ²MODAL project team, INRIA Lille, 59650 Villeneuve d'Ascq, France
 ³INRIA Saclay - Île-de-France, 91120 Palaiseau, France
 ⁴ XPOP project team, INRIA Saclay - Île-de-France, 91120 Palaiseau, France



Summary

Since the 90s, model-based clustering is largely used to classify data. Nowadays, with the increase of available data, missing values are more frequent. We defend the need to embed the missingness mechanism directly within the clustering modeling step. There exist three types of missing data: missing completely at random (MCAR), missing at random (MAR) and missing not at random (MNAR). In all situations, logistic regression is proposed as a natural and flexible candidate model. In this unified context, standard model selection criteria can be used to select between such different missing data mechanisms, simultaneously with the number of clusters. Practical interest of our proposal is illustrated on data derived from medical studies suffering from many missing data.

Model

Model-based clustering:

$$P(y_i;\theta) = \sum_{k=1}^{K} \pi_k \phi_k(y_i)$$

where $\star \pi_k = P(z_i^k = 1),$

 $\star \phi_k$ are gaussian or multinomial distributions, with parameters μ_k and Σ_k or p_k .

Dataset Example

 $\mathbf{y} = \begin{pmatrix} 18 & ? & blue & 175 \\ 25 & 63 & green & 162 \\ ? & 56 & brown & 164 \\ 22 & 79 & ? & 182 \end{pmatrix}$

 $\star \theta = (\pi_1, \dots, \pi_K, (\mu_1, \Sigma_1)/p_1, \dots, (\mu_K, \Sigma_K)/p_K)$

Missing mechanism models:

 $logit(P(c_i^j | y_i, z_i; \psi)) = \alpha_0, \ \psi = \alpha_0 \ (MCAR)$ $logit(P(c_i^j | y_i, z_i = 1; \psi)) = \alpha_0 + \beta z_i, \ \psi = (\alpha_0, \beta_1, \dots, \beta_K) \ (MNAR\mathbf{z})$ $logit(P(c_i^j | y_i, z_i; \psi)) = \alpha_0 + \alpha_j y_i^j + \beta z_i, \ \psi = (\alpha_0, \dots, \alpha_d, \beta_1, \dots, \beta_K) \ (MNAR\mathbf{yz})$

Inference: EM-like algorithm

Separability of the update:

$$Q_{Y,C,Z}(\theta,\psi|\theta^{(t)},\psi^{(t)}) = Q_{Y,Z}(\theta|\theta^{(t)},\psi^{(t)}) + Q_C(\psi|\theta^{(t)},\psi^{(t)})$$

 \implies Maximize on θ and ψ separatly.

E step simple for MCAR and MNARz:

 $\tau_i^{k,(t)} \propto \pi_k \phi_k(y_i^{o_i}) P(c_i | y_i^{o_i}, z_i; \psi^{(t)})$

Simpler way: SEM algorithm

 $\mathbf{c} = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}$ $\mathbf{o}_{1} = (1, 3, 4)$ $\mathbf{y}_{1}^{\mathbf{o}_{1}} = (18, blue, 175)$ $\mathbf{m}_{1} = (2)$ $\mathbf{y}^{\mathbf{o}} = \{\mathbf{y}_{1}^{\mathbf{o}_{1}}, \dots, \mathbf{y}_{4}^{\mathbf{o}_{4}}\}$ $\mathbf{y}^{\mathbf{m}} = \{\mathbf{y}_{1}^{\mathbf{m}_{1}}, \dots, \mathbf{y}_{4}^{\mathbf{m}_{4}}\}$

Missing mechanism (1)

Missing Completely at random (MCAR): $P(c_i | \mathbf{Y}) = P(c_i)$

Missing at random (MAR): $P(c_i | \mathbf{Y}) = P(c_i | \mathbf{Y}^o)$

Missing not at random (MNAR): Not MCAR nor MAR

Link with Classical Methods

All available cases:

Conditional independence of the variables knowing the cluster \Longleftrightarrow AAC to update θ in each step

Concatenation: MBC on Y with a MNARz mechanism \iff MBC on (Y|C) with a MAR mechanism

Application

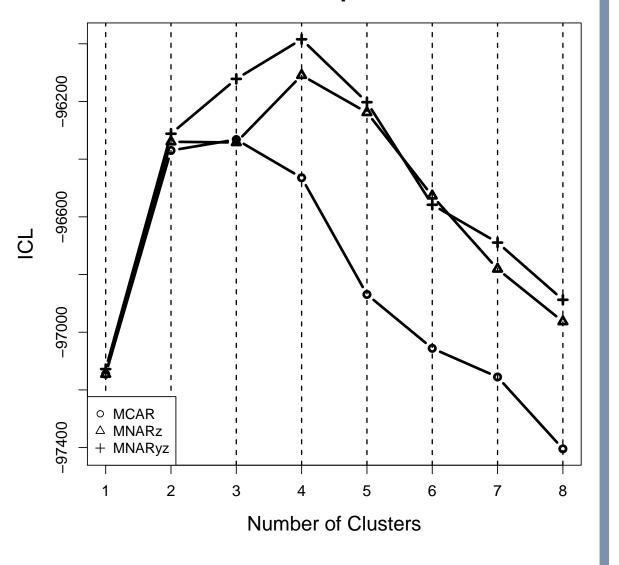
Dataset:

Number of individuals: n = 5 146 Number of feature: d = 7Percentage of missing data: 6.4%

Results:

MCAR and MNARz are equivalent until K = 3

ICL Comparison



References
 Little, R.J.A. and Rubin, D.B. Statis- tical analysis with missing data. 1986. John Wiley & Sons, Inc. New York, NY, USA, 1986

MNARz and MNARyz clearly indicate presence of an additional cluster (K = 4)

Conclusion

The method is implemented in the **R** Package "Relatedness", that performs inference on phased/unphased genotypic data, and accounts for a potential population structure. Core functions are developed in C++ and parallelized to speed-up inference on large panels. More information is available in *F. Laporte, C. Charcosset and T. Mary-Huard, Estimation of the relatedness coefficients from biallelic markers, application in plant mating designs, Biometrics in press.*