

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

SIMPAS Group Meeting

Aude Sportisse ^{1,3} Claire Boyer ^{1,2} Julie Josse ^{3,4,5}

¹Laboratoire de Probabilités Statistique et Modélisation, Sorbonne Université

²Département de Mathématiques et applications, Ecole Normale Supérieure

³Centre de Mathématiques Appliquées, Ecole Polytechnique

⁴XPOP, INRIA

⁵Visiting Researcher Google Brain

18th June 2020

Missing values are everywhere

- ▶ for many reasons: unanswered questions in a survey, lost data, damaged plants, machines that fail...

Traumabase: 15 000 patients/ 250 var/ 15 hospitals

Center	Age	Sex	Weight	Height	Heart rate	Lactates
Beaujon	54	m	85	NA	NA	NA
Lille	33	m	80	1.8	180	4.8
Pitie	26	m	NA	NA	NA	3.9
Beaujon	63	m	80	1.8	190	1.66
Pitie	30	w	NA	NA	NA	NA

NA: Not Available.

Classical definitions

- ★ $Y \in \mathbb{R}^{n \times p}$ the data matrix,
- ★ $\Omega \in \mathbb{R}^{n \times p}$ the missing-data pattern:

$$\Omega_{ij} = \begin{cases} 1 & \text{if } Y_{ij} \text{ is observed,} \\ 0 & \text{otherwise,} \end{cases}$$

- ★ **Issue: Cause of the missingness ?**
[Rubin, 1976], [Little and Rubin, 2014]

Missing-data mechanism

ϕ : the unknown parameters of the missingness.

MCAR mechanism

The missingness does not depend on the variables.

$$p(\Omega|Y; \phi) = p(\Omega; \phi), \quad \forall Y, \phi.$$

MAR mechanism

The missingness depends on the observed variables.

$$p(\Omega|Y; \phi) = p(\Omega|Y_{\text{obs}}; \phi), \quad \forall Y_{\text{mis}}, \phi.$$

MNAR mechanism

Other case, i.e.

$$p(\Omega|Y; \phi) = p(\Omega|Y_{\text{obs}}, Y_{\text{mis}}; \phi), \quad \forall \phi.$$

Self-masked setting when the missingness of a variable depends on the variable itself:

$$p(\Omega_{.j}|Y; \phi) = p(\Omega_{.j}|Y_{.j}; \phi), \quad \forall \phi.$$

MNAR data are very frequent in practice...

Traumabase: 15 000 patients/ 250 var/ 15 hospitals

Center	Age	Sex	Weight	Height	Heart rate	Lactates
Beaujon	54	m	85	NA	NA	NA
Lille	33	m	80	1.8	180	4.8
Pitie	26	m	NA	NA	NA	3.9
Beaujon	63	m	80	1.8	190	1.66
Pitie	30	w	NA	NA	NA	NA

- MNAR case extremely frequent, such as the heart rate (HR).
patient's condition critical
→ HR high or low
→ doctors would rather provide emergency care than measure HR.

...but hard to handle

- ✓ In the **MAR** setting, one can **ignore the mechanism**. Statistical inference is possible without modelling the missing-data mechanism distribution.
- ✗ In the **MNAR** setting, the observed variables are not representative of the population. One should **consider the mechanism**.

...but hard to handle

- ✓ In the **MAR** setting, one can **ignore the mechanism**. Statistical inference is possible without modelling the missing-data mechanism distribution.
- ✗ In the **MNAR** setting, the observed variables are not representative of the population. One should **consider the mechanism**.
 - ▶ Estimation in the MNAR setting: we should **take into account the mechanism** explicitly (by modelling it) or implicitly.
 - ▶ Identifiability in the MNAR setting: the law is identifiable only if the mechanism is identifiable.

Identifiability issue in the MNAR case

$$Y^{\text{NA}} = [1, \text{NA}, 0, 1, \text{NA}, 0].$$

- ▶ $\mathbb{P}(Y)$ is not identifiable without knowing $\mathbb{P}(\Omega|Y)$.
 - Y missing only if $Y = 1$. Thus, $Y = [1, 1, 0, 1, 1, 0]$ and $\mathbb{P}(Y) = 2/3$.
 - Y missing only if $Y = 0$. Thus, $Y = [1, 0, 0, 1, 0, 0]$ and $\mathbb{P}(Y) = 1/3$.
- ▶ 2 mechanisms yield to different data distribution.

Existing works for handling MNAR data (1)

How to estimate parameters of the data distribution in presence of MNAR data?

- ▶ By modeling the MNAR mechanism via a **Logistic Model**

$\forall i \in [1, n]$, $\phi_j = (\phi_{1j}, \phi_{2j})$ denoting a parameter vector:

$$p(\Omega_{ij} | y_{ij}; \phi) = [(1 + e^{-\phi_{1j}(y_{ij} - \phi_{2j})})^{-1}]^{(1 - \Omega_{ij})} [1 - (1 + e^{-\phi_{1j}(y_{ij} - \phi_{2j})})^{-1}]^{\Omega_{ij}}$$

and **using an EM algorithm** to estimate both parameters of the data and mechanism distributions.

- in linear models [Ibrahim et al., 1999],
- in low-rank models [S., Boyer, Josse 2018].

✓ Handling MNAR data.

✗ Often restricted to a limited number of MNAR variables.

✗ Parametric assumption for the mechanism distribution.

✗ Computationally costly.

Existing works for handling MNAR data (2)

How to estimate parameters of the data distribution in presence of MNAR data ?

- ▶ Without modeling the mechanism and by only using all available observed cells
 - for multivariate regression [Miao and Tchetgen, 2018, Tang et al., 2003],
 - in linear models, method based on graphical models [Mohan et al., 2018].
- ✓ Identifiability guarantees.
- ✓ No modelling for the mechanism.
- ✗ Restricted to simple models.
- ✗ The self-masked assumption may be strong.

Existing works for handling MNAR data (2)

How to estimate parameters of the data distribution in presence of MNAR data ?

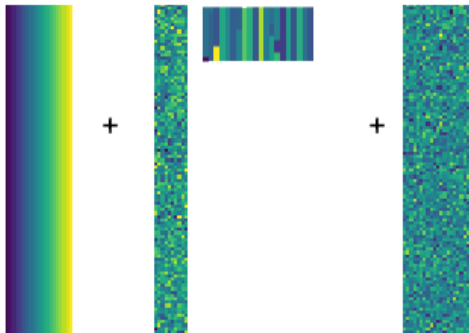
- ▶ Without modeling the mechanism and by only using all available observed cells
 - for multivariate regression [Miao and Tchetgen, 2018, Tang et al., 2003],
 - in linear models, method based on graphical models [Mohan et al., 2018].
- ✓ Identifiability guarantees.
- ✓ No modelling for the mechanism.
- ✗ Restricted to simple models.
- ✗ The self-masked assumption may be strong.

⇒ **Proposal: Handling the MNAR data in probabilistic PCA model.**

Probabilistic Principal Component Analysis (PPCA) model

$Y \in \mathbb{R}^{n \times p}$ noisy realisation of the factorization of the **coefficients matrix**
 $B \in \mathbb{R}^{r \times p}$ and r **latent variables** grouped in the matrix $W \in \mathbb{R}^{n \times r}$:

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



Probabilistic Principal Component Analysis (PPCA) model

$Y \in \mathbb{R}^{n \times p}$ noisy realisation of the factorization of the **coefficients matrix** $B \in \mathbb{R}^{r \times p}$ and r **latent variables** grouped in the matrix $W \in \mathbb{R}^{n \times r}$:

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

where

$$\left\{ \begin{array}{l} W = (W_1 | \dots | W_n)^T, \text{ with } W_i \sim \mathcal{N}(0_r, \text{Id}_{r \times r}), \\ B \text{ with rank } r < \min\{n, p\}, \\ \alpha \in \mathbb{R}^p \text{ and } \mathbf{1} = (1 \dots 1)^T \in \mathbb{R}^n, \\ \epsilon = (\epsilon_1 | \dots | \epsilon_n)^T, \text{ with } \epsilon_i \sim \mathcal{N}(0_p, \sigma^2 \text{Id}_{p \times p}). \end{array} \right.$$

$$\forall i \in \{1, \dots, n\}, \quad Y_i \sim \mathcal{N}(\alpha, \Sigma), \quad \Sigma = B^T B + \sigma^2 \text{Id}_{p \times p}$$

Probabilistic Principal Component Analysis (PPCA) model

$Y \in \mathbb{R}^{n \times p}$ noisy realisation of the factorization of the **coefficients matrix** $B \in \mathbb{R}^{r \times p}$ and r **latent variables** grouped in the matrix $W \in \mathbb{R}^{n \times r}$:

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

--> Access only to the missing-data matrix $Y \odot \Omega + \text{NA} \odot (1 - \Omega)$

$$Y = \begin{pmatrix} Y_{.1} & Y_{.2} & Y_{.3} & Y_{.4} & \dots & Y_{.p} \\ 12 & 28 & 31 & 3 & \dots & \text{NA} \\ \text{NA} & 23 & 89 & 2 & \dots & 85 \\ 32 & 6 & 24 & \text{NA} & \dots & \text{NA} \\ \vdots & \vdots & \vdots & \vdots & \dots & \vdots \\ \text{NA} & 3 & 7 & \text{NA} & \dots & 11 \end{pmatrix}$$

Proposal

Under the PPCA model $\forall i \in \{1, \dots, n\}$, $Y_i \sim \mathcal{N}(\alpha, \Sigma)$, $\Sigma = B^T B + \sigma^2 \text{Id}_{p \times p}$,

- **Identifiability** of the PPCA parameters assuming self-masked MNAR,

Proposition: identifiability of the PPCA parameters

Consider that d variables are self-masked MNAR and $p - d$ variables are MCAR.

- ▶ The parameters (α, Σ) of the PPCA model and the mechanism parameter are identifiable.
 - ▶ Assuming that the noise level σ^2 is known, the parameter B is identifiable up to a row permutation.
- Mean and covariance matrix estimation **without explicitly modeling the MNAR mechanism and by only using all available observed cells**
Consistency results assuming general MNAR mechanism,
 - Estimation of the loading matrix B ,
 - Imputation of the missing values in the data matrix Y .

Toy example

- ▶ $p = 3, r = 2$.
- ▶ $Y_{.1}$ is MNAR (self-masked in this case).
- ▶ As $r = 2$, it requires two **pivot variables**, say $Y_{.2}$ and $Y_{.3}$ **which are independent of the missing-data pattern** $\Omega_{.1}$.

$$(Y_{.1} \ Y_{.2} \ Y_{.3}) = \mathbf{1} (\alpha_1 \ \alpha_2 \ \alpha_3) + (W_{.1} \ W_{.2}) B + \epsilon$$

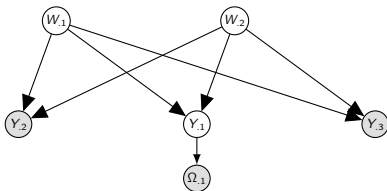
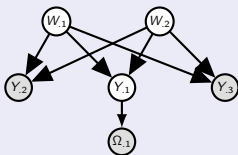


Figure: Graphical model for "fully-connected" PPCA model when $p = 3, r = 2$ and one variable $Y_{.1}$ is missing.

Mean estimation

Exploiting the linear links

As any variable is generated by all the latent variables, linear links can be established.



Lemma exploiting the linear links between variables

$$\mathbf{Y}_{.2} = \mathcal{B}_{2 \rightarrow 1,3[0]} + \mathcal{B}_{2 \rightarrow 1,3[1]} \mathbf{Y}_{.1} + \mathcal{B}_{2 \rightarrow 1,3[3]} \mathbf{Y}_{.3} + \zeta,$$

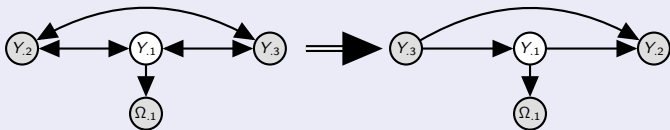
where

- $(\mathcal{B}_{2 \rightarrow 1,3[k]})_{k \in \{0,1,3\}}$ denotes the coefficients standing for the effects of $Y_{.2}$ on $Y_{.1}$ and $Y_{.3}$.
- $(\mathcal{B}_{2 \rightarrow 1,3[k]})_{k \in \{0,1,3\}}$ depends on the unknown matrix B and ζ the noise
- Note that $\mathbb{E}[\zeta | Y_{.1}, Y_{.3}] \neq 0$ (no exogeneity)

Mean estimation

Exploiting the linear links

As any variable is generated by all the latent variables, linear links can be established.

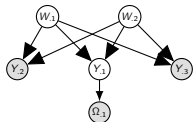


Lemma exploiting the linear links between variables

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

where

- $(B_{2 \rightarrow 1,3[k]})_{k \in \{0,1,3\}}$ denotes the coefficients standing for the effects of $Y_{.2}$ on $Y_{.1}$ and $Y_{.3}$.
- $(B_{2 \rightarrow 1,3[k]})_{k \in \{0,1,3\}}$ depends on the unknown matrix B and ζ the noise
- Note that $\mathbb{E}[\zeta | Y_{.1}, Y_{.3}] \neq 0$ (no exogeneity)



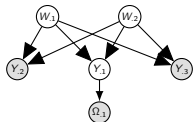
Mean estimation

Estimating in the complete-case analysis

Effects of $Y_{.2}$ on $Y_{.1}$ and $Y_{.3}$ in the complete case when $\Omega_{.1} = 1$:

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

$$Y = \begin{pmatrix} Y_{.1} & Y_{.2} & Y_{.3} \\ 12 & 28 & 31 \\ \text{NA} & 23 & 89 \\ 32 & 6 & 24 \\ \vdots & \vdots & \vdots \\ \text{NA} & 3 & 7 \end{pmatrix}$$



Mean estimation

Estimating in the complete-case analysis

Effects of Y_2 on Y_1 and Y_3 in the complete case when $\Omega_1 = 1$:

$$(Y_2)_{|\Omega_1=1} := \mathcal{B}_{2 \rightarrow 1,3[0]}^c + \mathcal{B}_{2 \rightarrow 1,3[1]}^c Y_1 + \mathcal{B}_{2 \rightarrow 1,3[2]}^c Y_3 + \zeta^c,$$

As $Y_2 \perp\!\!\!\perp \Omega_1 | Y_1, Y_3$, one has

$$\mathbb{E}[Y_2 | Y_1, Y_3, \Omega_1 = 1] = \mathbb{E}[\mathcal{B}_{2 \rightarrow 1,3[0]}^c + \mathcal{B}_{2 \rightarrow 1,3[1]}^c Y_1 + \mathcal{B}_{2 \rightarrow 1,3[3]}^c Y_3 | Y_1, Y_3].$$

Taking the expectation,

$$\mathbb{E}[Y_2] = \mathcal{B}_{2 \rightarrow 1,3[0]}^c + \mathcal{B}_{2 \rightarrow 1,3[1]}^c \mathbb{E}[Y_1] + \mathcal{B}_{2 \rightarrow 1,3[3]}^c \mathbb{E}[Y_3].$$

Mean formula

$$\alpha_1 = \frac{\alpha_2 - \mathcal{B}_{2 \rightarrow 1,3[0]}^c - \mathcal{B}_{2 \rightarrow 1,3[3]}^c \alpha_3}{\mathcal{B}_{2 \rightarrow 1,3[1]}^c},$$

given that $\mathcal{B}_{2 \rightarrow 1,3[1]}^c \neq 0$.

Consistency results

Definition of a mean estimator:

$$\hat{\alpha}_1 := \frac{\hat{\alpha}_2 - \hat{\mathcal{B}}_{2 \rightarrow 1,3[0]}^c - \hat{\mathcal{B}}_{2 \rightarrow 1,3[3]}^c \hat{\alpha}_3}{\hat{\mathcal{B}}_{2 \rightarrow 1,3[1]}^c}.$$

Consistency for the missing variable mean

Assume that:

- ▶ There exist consistent estimators for α_2 and α_3 .
- ▶ There exist consistent estimators for $\mathcal{B}_{2 \rightarrow 1,3[0]}^c$, $\mathcal{B}_{2 \rightarrow 1,3[1]}^c$ and $\mathcal{B}_{2 \rightarrow 1,3[3]}^c$.

Then, the estimator $\hat{\alpha}_1$ is consistent.

Estimation in practice

$$\hat{\alpha}_1 := \frac{\hat{\alpha}_2 - \hat{\beta}_{2 \rightarrow 1,3[0]}^c - \hat{\beta}_{2 \rightarrow 1,3[3]}^c \hat{\alpha}_3}{\hat{\beta}_{2 \rightarrow 1,3[1]}^c}.$$

- ▶ $\hat{\alpha}_2$ and $\hat{\alpha}_3$ are computed as empirical quantities.

$$\triangleright \hat{\alpha}_2 = \bar{Y}_{.2}$$

$$\triangleright \hat{\alpha}_3 = \bar{Y}_{.3}$$

$$Y = \begin{array}{c} Y_{.1} \quad Y_{.2} \quad Y_{.3} \\ \left(\begin{array}{ccc} 12 & 28 & 31 \\ \text{NA} & 23 & 89 \\ 32 & 6 & 24 \\ \vdots & \vdots & \vdots \\ \text{NA} & 3 & 7 \end{array} \right) \end{array}$$

Estimation in practice

$$\hat{\alpha}_1 := \frac{\hat{\alpha}_2 - \hat{\beta}_{2 \rightarrow 1,3[0]}^c - \hat{\beta}_{2 \rightarrow 1,3[3]}^c \hat{\alpha}_3}{\hat{\beta}_{2 \rightarrow 1,3[1]}^c}.$$

- ▶ $\hat{\alpha}_2$ and $\hat{\alpha}_3$ are computed as empirical quantities.
 - ▷ $\hat{\alpha}_2 = \bar{Y}_{.2}$
 - ▷ $\hat{\alpha}_3 = \bar{Y}_{.3}$
- ▶ $(\beta_{2 \rightarrow 1,3[k]}^c)_{k \in \{0,1,3\}}$ estimated by the coefficients of the linear regression of $Y_{.2}$ on $Y_{.1}$ and $Y_{.3}$ using the rows where $Y_{.1}$ is observed.

$$Y = \begin{pmatrix} Y_{.1} & Y_{.2} & Y_{.3} \\ 12 & 28 & 31 \\ \text{NA} & 23 & 89 \\ 32 & 6 & 24 \\ \vdots & \vdots & \vdots \\ \text{NA} & 3 & 7 \end{pmatrix}$$

$$Y = \begin{pmatrix} Y_{.1} & Y_{.2} & Y_{.3} \\ 12 & 28 & 31 \\ \text{NA} & 23 & 89 \\ 32 & 6 & 24 \\ \vdots & \vdots & \vdots \\ \text{NA} & 3 & 7 \end{pmatrix}$$

Estimation of the loading matrix B

- ▶ Same methodology for the variance and covariances.
- ▶ Estimators obtained from the formulae:

$$\hat{\Sigma} = \begin{pmatrix} \widehat{\text{Var}}(Y_1) & \widehat{\text{Cov}}(Y_1, Y_2) & \widehat{\text{Cov}}(Y_1, Y_3) \\ \widehat{\text{Cov}}(Y_2, Y_1) & \widehat{\text{Var}}(Y_2) & \widehat{\text{Cov}}(Y_2, Y_3) \\ \widehat{\text{Cov}}(Y_3, Y_1) & \widehat{\text{Cov}}(Y_3, Y_2) & \widehat{\text{Var}}(Y_3) \end{pmatrix}$$

Estimation of the loading matrix B

- ▶ Same methodology for the variance and covariances.
- ▶ Estimators obtained from the formulae:

$$\hat{\Sigma} = \begin{pmatrix} \widehat{\text{Var}}(Y_{.1}) & \widehat{\text{Cov}}(Y_{.1}, Y_{.2}) & \widehat{\text{Cov}}(Y_{.1}, Y_{.3}) \\ \widehat{\text{Cov}}(Y_{.2}, Y_{.1}) & \widehat{\text{Var}}(Y_{.2}) & \widehat{\text{Cov}}(Y_{.2}, Y_{.3}) \\ \widehat{\text{Cov}}(Y_{.3}, Y_{.1}) & \widehat{\text{Cov}}(Y_{.3}, Y_{.2}) & \widehat{\text{Var}}(Y_{.3}) \end{pmatrix}$$

- ▶ Assuming that σ^2 is known,

$$Y \sim \mathcal{N} \left(\begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{pmatrix}, B^T B + \sigma^2 \text{Id} \right) \Rightarrow \hat{\Sigma} - \sigma^2 \text{Id}_{3 \times 3} \text{ estimates } B^T B.$$

- ▶ Singular value decomposition:

$$\hat{\Sigma} - \sigma^2 \text{Id}_{3 \times 3} =: \hat{U} \hat{D} \hat{U}^T, \text{ with } \hat{U} = (\hat{u}_1 | \hat{u}_2 | \hat{u}_3).$$

- ▶ Assuming that $r = 2$,

$$\hat{B} = \hat{D}_{|2}^{1/2} \hat{U}_{|2}^T = \begin{pmatrix} \sqrt{\hat{d}_1} & 0 \\ 0 & \sqrt{\hat{d}_2} \end{pmatrix} \begin{pmatrix} \hat{u}_1^T \\ \hat{u}_2^T \end{pmatrix}.$$

Imputation of the missing values in Y

- Impute the missing values Y_{i1} for $i \in \{1, \dots, n\}$ such that $M_{i1} = 0$ using the conditional expectation of (Y_{i1}) given Y_{i2} and Y_{i3} .

$$Y = \begin{pmatrix} Y_{.1} & Y_{.2} & Y_{.3} \\ 12 & 28 & 31 \\ \text{NA} & 23 & 89 \\ 32 & 6 & 24 \\ \vdots & \vdots & \vdots \\ \text{NA} & 3 & 7 \end{pmatrix} \rightarrow Y = \begin{pmatrix} Y_{.1} & Y_{.2} & Y_{.3} \\ 12 & 28 & 31 \\ 16 & 23 & 89 \\ 32 & 6 & 24 \\ \vdots & \vdots & \vdots \\ 21 & 3 & 7 \end{pmatrix}$$

The methodology is extended to the general case, for any data with p covariates, r latent variables and d missing variables.

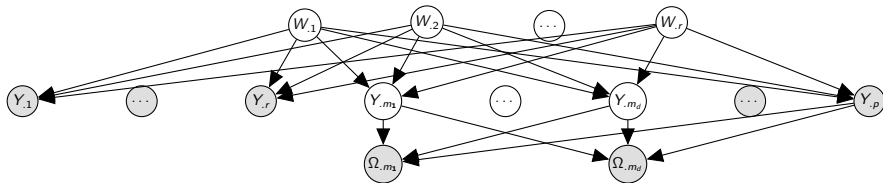
$$Y = \begin{pmatrix} Y_{.1} & Y_{.2} & Y_{.3} & \dots & Y_{.p} \\ \text{NA} & 23 & \text{NA} & \dots & 45 \\ 32 & \text{NA} & 24 & \dots & \text{NA} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \text{NA} & 3 & \text{NA} & \dots & 46 \end{pmatrix} \rightarrow Y = \begin{pmatrix} Y_{.1} & Y_{.2} & Y_{.3} & \dots & Y_{.p} \\ 16 & 23 & 89 & \dots & 45 \\ 32 & 6 & 24 & \dots & 22 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 21 & 3 & 7 & \dots & 46 \end{pmatrix}$$

Imputation of the missing values in Y

- Impute the missing values Y_{i1} for $i \in \{1, \dots, n\}$ such that $M_{i1} = 0$ using the conditional expectation of (Y_{i1}) given Y_{i2} and Y_{i3} .

$$Y = \begin{pmatrix} Y_{.1} & Y_{.2} & Y_{.3} \\ 12 & 28 & 31 \\ \text{NA} & 23 & 89 \\ 32 & 6 & 24 \\ \vdots & \vdots & \vdots \\ \text{NA} & 3 & 7 \end{pmatrix} \rightarrow Y = \begin{pmatrix} Y_{.1} & Y_{.2} & Y_{.3} \\ 12 & 28 & 31 \\ 16 & 23 & 89 \\ 32 & 6 & 24 \\ \vdots & \vdots & \vdots \\ 21 & 3 & 7 \end{pmatrix}$$

The methodology is extended to the general case, for any data with p covariates, r latent variables and d missing variables.



Numerical experiments

Measuring the performance

- estimation of B : RV coefficient (cosine between two subspaces).
- imputation of Y : $\|(\hat{Y} - Y) \odot (1 - \Omega)\|_F^2 / \|Y \odot (1 - \Omega)\|_F^2$.

Other methods

- **MAR** our method which has been adapted to handle MAR data (inspired by [Mohan et al., 2018] in linear models);
- **EMMAR**: EM algorithm to perform PPCA with MAR values [Ilin and Raiko, 2010];
- **SoftMAR**: matrix completion using iterative soft-thresholding singular value decomposition algorithm [Mazumder et al., 2010] relevant only for M(C)AR values;
- **MNARparam**: matrix completion technique modeling the MNAR mechanism with a parametric logistic model [Sportisse et al., 2018];
- **Del**: the naive listwise deletion method;
- **Mean**: the imputation by the mean.

Numerical experiments

Synthetic data

$n = 1000$, $p = 20$, $r = 3$, $\sigma = 0.8$

10 MNAR variables with

$$\forall m \in [1 : 10], \mathbb{P}(\Omega_{.m} = 1 | Y) = \mathbb{P}(\Omega_{.m} = 1 | Y_{.m}, Y_{.k}, Y_{.l}),$$

where k and l are indexes of MNAR variables randomly chosen such that $k \neq l \in [1 : 10] \setminus \{m\}$.

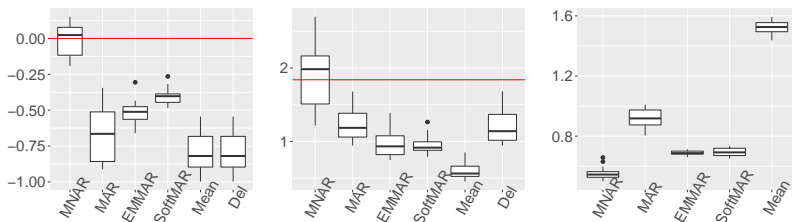


Figure: Mean estimation (left graphic), variance estimation (middle graphic) of one missing variable and prediction error (right graphic).

Numerical experiments

Real data

- ▶ Introduction additional MNAR values in the variable *HR.ph* using a logistic self-masked mechanism. Other variables considered as M(C)AR.

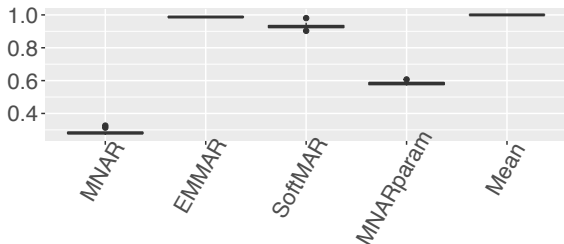
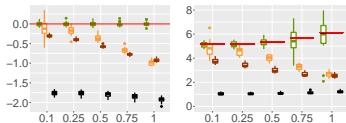
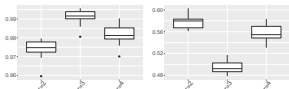
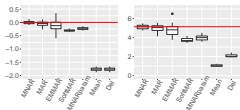


Figure: Comparison of the prediction error for the TraumaBase data.

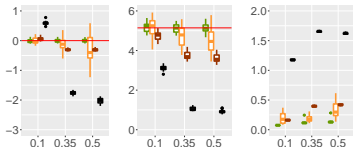
Numerical experiments

Others

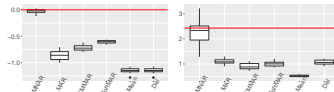


■ MNAR
■ EMMAR
■ SoftMAR
■ Mean

- Robustness to the percentage of missing values.
- Robustness to the noise.
- Model (PPCA) misspecification
- Rank misspecification.



■ MNAR
■ EMMAR
■ SoftMAR
■ Mean



Conclusion

Take-home messages

- MNAR is hard.
- Modeling the mechanism and using an EM algorithm is computationally costly.
- Our proposal: **new estimation and imputation method to perform PPCA with MNAR data,**
- **without any need of modeling the missing mechanism,**
- with strong theoretical guarantees as identifiability and consistency and efficient algorithm.

Perspectives

- Estimating the rank in the PPCA setting with MNAR data.
- Extension to the exponential family to process count data.

► Estimation and imputation in probabilistic principal component analysis with missing not at random data. [Sportisse et al., 2020]

References I



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



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



Little, R. J. and Rubin, D. B. (2014).
Statistical analysis with missing data, volume 333.
John Wiley & Sons.









Mazumder, R., Hastie, T., and Tibshirani, R. (2010).
Spectral regularization algorithms for learning large incomplete matrices.
Journal of machine learning research, 11(Aug):2287–2322.



Miao, W. and Tchetgen, E. T. (2018).
Identification and inference with nonignorable missing covariate data.
Statistica Sinica, 28(4):2049–2067.

References II

-  Mohan, K., Thoemmes, F., and Pearl, J. (2018). Estimation with incomplete data: The linear case. In *IJCAI*, pages 5082–5088.
-  Pearl, J. (2003). Causality: models, reasoning, and inference. *Econometric Theory*, 19(675-685):46.
-  Rubin, D. B. (1976). Inference and missing data. *Biometrika*, 63(3):581–592.
-  Sportisse, A., Boyer, C., and Josse, J. (2018). Imputation and low-rank estimation with missing non at random data. *arXiv preprint arXiv:1812.11409*.
-  Sportisse, A., Boyer, C., and Josse, J. (2020). Estimation and imputation in probabilistic principal component analysis with missing not at random data.
-  Tang, G., Little, R. J., and Raghunathan, T. E. (2003). Analysis of multivariate missing data with nonignorable nonresponse. *Biometrika*, 90(4):747–764.

Missing-data mechanism

ϕ : the unknown parameters of the missingness.

Y_1	Y_2
1	2
3	20
22	4

MCAR mechanism

The missingness does not depend on the variables.

$$p(\Omega|Y; \phi) = p(\Omega; \phi), \quad \forall Y, \phi.$$

MAR

1	2
NA	20
22	4

MAR mechanism

The missingness depends on the observed variables.

$$p(\Omega|Y; \phi) = p(\Omega|Y_{\text{obs}}; \phi), \quad \forall Y_{\text{mis}}, \phi.$$

MNAR

1	2
3	20
NA	4

MNAR mechanism

Other case, i.e.

$$p(\Omega|Y; \phi) = p(\Omega|Y_{\text{obs}}, Y_{\text{mis}}; \phi), \quad \forall \phi.$$

self-masked when the missingness of a variable depends on the variable itself.

...but hard to handle

- ▶ Likelihood approach: maximizing the joint log-likelihood,

$$\ell(\Theta, \phi; Y, M) = p(Y; \Theta)p(M|Y; \phi),$$

with Θ : unknown data distribution parameter.

- ▶ Missing data: maximizing the observed joint log-likelihood,

$$\ell(\Theta, \phi; Y_{\text{obs}}, M) = \int \ell(\Theta, \phi; Y, M) dY_{\text{mis}}.$$

- ✓ In the MAR setting, one can ignore the mechanism since

$$p(M|Y; \phi) = p(M|Y_{\text{obs}}; \phi).$$

$$\Rightarrow \ell(\Theta, \phi; Y_{\text{obs}}, M) \propto \ell(\Theta; Y_{\text{obs}}) = \int \ell(\Theta; Y) dY_{\text{mis}}.$$

- ✗ In the MNAR setting, one should consider the mechanism.

Existing works for handling MNAR data (1)

- ▶ Modeling the MNAR mechanism via a **Logistic Model**:

$\forall i \in [1, n]$, $\phi_j = (\phi_{1j}, \phi_{2j})$ denoting a parameter vector:

$$p(M_{ij}|y_{ij}; \phi) = [(1 + e^{-\phi_{1j}(y_{ij} - \phi_{2j})})^{-1}]^{(1 - M_{ij})} [1 - (1 + e^{-\phi_{1j}(y_{ij} - \phi_{2j})})^{-1}]^{M_{ij}}$$

(self-masked MNAR mechanism)

- ▶ Maximizing the joint log-likelihood with the **EM algorithm** in linear models [Ibrahim et al., 1999] or in low-rank models [S., Boyer, Josse 2018].

- ✓ Handling MNAR data.
- ✗ Often restricted to a limited number of MNAR variables.
- ✗ Parametric assumption for the mechanism distribution.
- ✗ Computationally costly.

General setting

- ▶ d MNAR variables indexed by $\mathcal{M} := \{m_1, \dots, m_d\} \subset \{1, \dots, p\}$ (with $d < p$).
- ▶ Other variables are observed or M(C)AR.
- ▶ The distribution of the mechanism may depend on all variables (missing or observed) except r pivot variables, indexed by \mathcal{J} .

$$\forall m \in \mathcal{M}, \quad \mathbb{P}(\Omega_{.m} = 1 | Y) = \mathbb{P}(\Omega_{.m} = 1 | (Y_{.k})_{k \in \bar{\mathcal{J}}}),$$
$$\bar{\mathcal{J}} = \{1, \dots, p\} \setminus \mathcal{J}.$$

Main assumptions

- A1.** $\forall m \in \mathcal{M}, \forall j \in \mathcal{J}, (B_{.m} \ (B_{.j'})_{j' \in \mathcal{J}_{-j}})$ is invertible.
(\rightsquigarrow implies that any variable is generated by all the latent variables)
- A2.** $\forall m \in \mathcal{M}, \forall j \in \mathcal{J}, Y_{.j} \perp\!\!\!\perp \Omega_{.m} | (Y_{.k})_{k \in \bar{\mathcal{J}}}$
(\rightsquigarrow follows from the missing-data mechanism above)

Imputation of the missing values in Y

► Impute the missing values Y_{i1} for $i \in \{1, \dots, n\}$ such that $M_{i1} = 0$ using the conditional expectation of (Y_{i1}) given Y_{i2} and Y_{i3} .

$$\mathbb{E}[Y_{i1}|Y_{i2}, Y_{i3}] = \alpha_1 + (\Gamma_{12} \quad \Gamma_{13}) \begin{pmatrix} \Gamma_{22} & \Gamma_{23} \\ \Gamma_{32} & \Gamma_{33} \end{pmatrix}^{-1} \left(\begin{pmatrix} Y_{i2} \\ Y_{i3} \end{pmatrix} - \begin{pmatrix} \alpha_2 \\ \alpha_3 \end{pmatrix} \right),$$

with $\Gamma = B^T B + \sigma^2 \text{Id}_{3 \times 3}$.

$$\hat{Y}_{i1} = \hat{\alpha}_1 + (\hat{\Gamma}_{12} \quad \hat{\Gamma}_{13}) \begin{pmatrix} \hat{\Gamma}_{22} & \hat{\Gamma}_{23} \\ \hat{\Gamma}_{32} & \hat{\Gamma}_{33} \end{pmatrix}^{-1} \left(\begin{pmatrix} Y_{i2} \\ Y_{i3} \end{pmatrix} - \begin{pmatrix} \hat{\alpha}_2 \\ \hat{\alpha}_3 \end{pmatrix} \right).$$

$$Y = \begin{pmatrix} Y_{.1} & Y_{.2} & Y_{.3} \\ 12 & 28 & 31 \\ \text{NA} & 23 & 89 \\ 32 & 6 & 24 \\ \vdots & \vdots & \vdots \\ \text{NA} & 3 & 7 \end{pmatrix} \rightarrow Y = \begin{pmatrix} Y_{.1} & Y_{.2} & Y_{.3} \\ 12 & 28 & 31 \\ 16 & 23 & 89 \\ 32 & 6 & 24 \\ \vdots & \vdots & \vdots \\ 21 & 3 & 7 \end{pmatrix}$$

Imputation of the missing values in Y

► Impute the missing values Y_{i1} for $i \in \{1, \dots, n\}$ such that $M_{i1} = 0$ **using the conditional expectation of (Y_{i1}) given Y_{i2} and Y_{i3} .**

$$\mathbb{E}[Y_{i1}|Y_{i2}, Y_{i3}] = \alpha_1 + \begin{pmatrix} \Gamma_{12} & \Gamma_{13} \end{pmatrix} \begin{pmatrix} \Gamma_{22} & \Gamma_{23} \\ \Gamma_{32} & \Gamma_{33} \end{pmatrix}^{-1} \left(\begin{pmatrix} Y_{i2} \\ Y_{i3} \end{pmatrix} - \begin{pmatrix} \alpha_2 \\ \alpha_3 \end{pmatrix} \right),$$

with $\Gamma = B^T B + \sigma^2 \text{Id}_{3 \times 3}$.

$$\hat{Y}_{i1} = \hat{\alpha}_1 + \begin{pmatrix} \hat{\Gamma}_{12} & \hat{\Gamma}_{13} \end{pmatrix} \begin{pmatrix} \hat{\Gamma}_{22} & \hat{\Gamma}_{23} \\ \hat{\Gamma}_{32} & \hat{\Gamma}_{33} \end{pmatrix}^{-1} \left(\begin{pmatrix} Y_{i2} \\ Y_{i3} \end{pmatrix} - \begin{pmatrix} \hat{\alpha}_2 \\ \hat{\alpha}_3 \end{pmatrix} \right).$$

The methodology is extended to the general case, for any data with p covariates, r latent variables and d missing variables.

Numerical experiments

Toy example setting (1)

★ $r = 2$, $p = 10$, $n = 1000$, $\sigma = 0.1$, 7 self-masked MNAR variables.

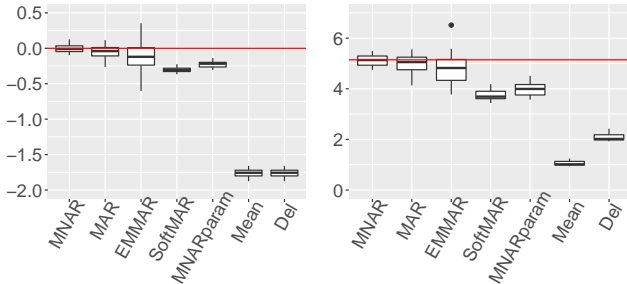


Figure: Mean and variance estimation for one MNAR variable.

Numerical experiments

Toy example setting (2)

★ $r = 2$, $p = 10$, $n = 1000$, $\sigma = 0.1$, 7 self-masked MNAR variables.

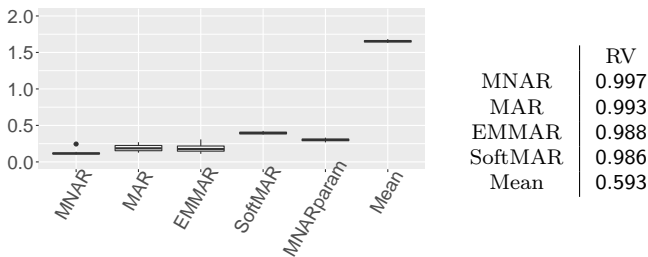


Figure: Prediction error (left) and median of the RV coefficients for the loading matrix (right).

Numerical experiments

Robustness to the noise (1)

Exogeneity does not hold but does it have an impact on the results when the noise increases?

★ $r = 2$, $p = 10$, $n = 1000$, 7 self-masked MNAR variables.

For different values of the noise level.

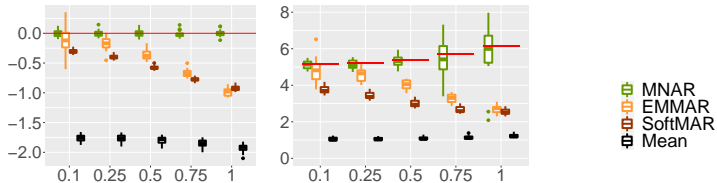


Figure: Mean and variance for one missing variable.

Numerical experiments

Robustness to the noise (2)

★ $r = 2$, $p = 10$, $n = 1000$, 7 self-masked MNAR variables.

For different values of the noise level.

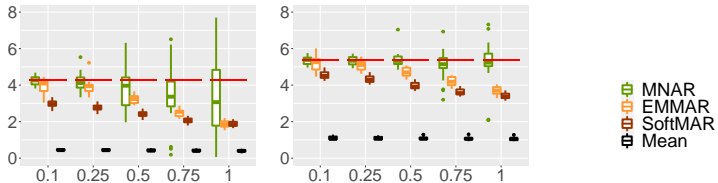


Figure: Covariances between a MNAR variable and a pivot one (left) and between two MNAR variables (right).

Numerical experiments

Robustness to the noise (3)

★ $r = 2$, $p = 10$, $n = 1000$, 7 self-masked MNAR variables.

For different values of the noise level.

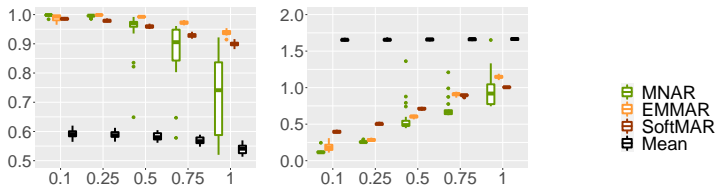


Figure: RV coefficient (left) and prediction error (right).

Numerical experiments

Robustness to the noise (3)

★ $r = 2$, $p = 10$, $n = 1000$, 7 self-masked MNAR variables.

For different values of the noise level.

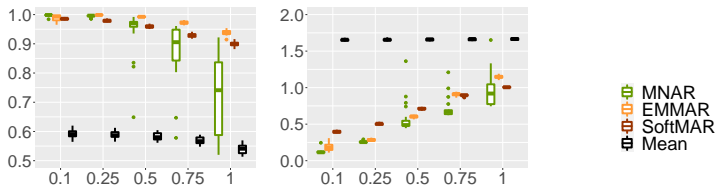


Figure: RV coefficient (left) and prediction error (right).

- ✗ Bias for the covariance MNAR/pivot variables.
- ✗ Loading matrix estimation deteriorates as the data gets noisier.
- ✓ Some estimations (mean, variance) are unbiased.
- ✓ In term of prediction error, it remains competitive.

Numerical experiments

Rank misspecification

★ $r = 3$, $p = 20$, $n = 1000$, $\sigma = 0.8$, 10 self-masked MNAR variables.

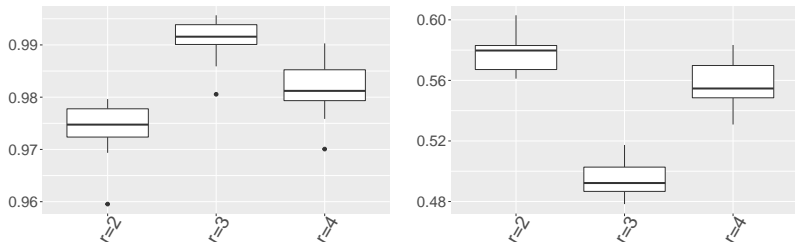


Figure: RV coefficients for the loading matrix (left graphic) and prediction error (right graphic) for different cases where the rank is either underestimated, well estimated or overestimated.

Numerical experiments

Rank misspecification

★ $r = 3$, $p = 20$, $n = 1000$, $\sigma = 0.8$, 10 self-masked MNAR variables.

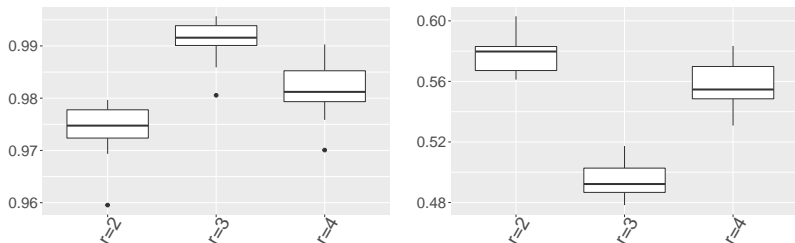


Figure: RV coefficients for the loading matrix (left graphic) and prediction error (right graphic) for different cases where the rank is either underestimated, well estimated or overestimated.

✓ Stability to rank misspecification.