

The Inria logo is written in a red, cursive script font.The Sartorius logo is written in a bold, black, sans-serif font inside a light gray rectangular box.

HOW TO DEAL WITH MISSING VALUES IN THE HIGH DIMENSIONAL SUPERVISED CONTEXT?

Hadrien Lorenzo & Olivier Cloarec & Jérôme Saracco

hadrien.lorenzo@u-bordeaux.fr [orator]

Inria, IMB, Sartorius

CMSTATISTICS 2021, KING'S COLLEGE LONDON, UK

December 19, 2021

Why talking about supervised imputation ?

Vaccine research context

- Important Ebola vaccine data-set (see Rechten et al. (2017)).
- $n = 20$, 4 blocks of high dimensions ($p_k = 20.000$, RNA-Seq).
- \mathbf{y} , 5-dimensional to be predicted.
- Many missing values: 30%.

Why talking about supervised imputation ?

Vaccine research context

- Important Ebola vaccine data-set (see Rechten et al. (2017)).
- $n = 20$, 4 blocks of high dimensions ($p_k = 20.000$, RNA-Seq).
- \mathbf{y} , 5-dimensional to be predicted.
- Many missing values: 30%.

Classical imputation

- Hypothesis: multivariate Gaussian.
- A few variables are imputed, let say $\mathcal{V}_{\text{imputation}}$.

Why talking about supervised imputation ?

Vaccine research context

- Important Ebola vaccine data-set (see Rechten et al. (2017)).
- $n = 20$, 4 blocks of high dimensions ($p_k = 20.000$, RNA-Seq).
- \mathbf{y} , 5-dimensional to be predicted.
- Many missing values: 30%.

Classical imputation

- Hypothesis: multivariate Gaussian.
- A few variables are imputed, let say $\mathcal{V}_{\text{imputation}}$.

Then supervised analysis (sparse PLS)

A few variables are selected for prediction, let say $\mathcal{V}_{\text{prediction}}$.

Why talking about supervised imputation ? (II)

Observation

$$\mathcal{V}_{\text{imputation}} \cap \mathcal{V}_{\text{prediction}} = \emptyset.$$

Why talking about supervised imputation ? (II)

Observation

$$\mathcal{V}_{\text{imputation}} \cap \mathcal{V}_{\text{prediction}} = \emptyset.$$

Interpretation

No imputed variables are of interest to answer the question.

Why talking about supervised imputation ? (II)

Observation

$$\mathcal{V}_{\text{imputation}} \cap \mathcal{V}_{\text{prediction}} = \emptyset.$$

Interpretation

No imputed variables are of interest to answer the question.

Causation

Imputed variables are associated to each other but not to \mathbf{y} .

Why talking about supervised imputation ? (II)

Observation

$$\mathcal{V}_{\text{imputation}} \cap \mathcal{V}_{\text{prediction}} = \emptyset.$$

Interpretation

No imputed variables are of interest to answer the question.

Causation

Imputed variables are associated to each other but not to \mathbf{y} .

Solution

Hide useless variables: regularization through supervised analysis.

Imputation and data augmentation

If \mathbf{x} is a random vector of parameter θ , $\mathbf{x}^{(obs)}$ is observed and $\mathbf{x}^{(miss)}$ is missing, a prior on θ is chosen by the user.

Imputation and data augmentation

If \mathbf{x} is a random vector of parameter θ , $\mathbf{x}^{(obs)}$ is observed and $\mathbf{x}^{(miss)}$ is missing, a prior on θ is chosen by the user.

The two posterior distributions $p(\mathbf{x}^{(miss)}|\mathbf{x}^{(obs)})$ and $p(\theta|\mathbf{x}^{(obs)})$ verify

$$p(\mathbf{x}^{(miss)}|\mathbf{x}^{(obs)}) = \int_{\theta \in \Theta} p(\mathbf{x}^{(miss)}|\theta, \mathbf{x}^{(obs)})p(\theta|\mathbf{x}^{(obs)})d\theta,$$

$$p(\theta|\mathbf{x}^{(obs)}) = \int_{\mathbf{x}^{(miss)} \in \mathbf{x}^{(miss)}} p(\theta|\mathbf{x}^{(miss)}, \mathbf{x}^{(obs)})p(\mathbf{x}^{(miss)}|\mathbf{x}^{(obs)})d\theta,$$

Imputation and data augmentation

If \mathbf{x} is a random vector of parameter θ , $\mathbf{x}^{(obs)}$ is observed and $\mathbf{x}^{(miss)}$ is missing, a prior on θ is chosen by the user.

The two posterior distributions $p(\mathbf{x}^{(miss)}|\mathbf{x}^{(obs)})$ and $p(\theta|\mathbf{x}^{(obs)})$ verify

$$p(\mathbf{x}^{(miss)}|\mathbf{x}^{(obs)}) = \int_{\theta \in \Theta} p(\mathbf{x}^{(miss)}|\theta, \mathbf{x}^{(obs)})p(\theta|\mathbf{x}^{(obs)})d\theta,$$

$$p(\theta|\mathbf{x}^{(obs)}) = \int_{\mathbf{x}^{(miss)} \in \mathcal{X}^{(miss)}} p(\theta|\mathbf{x}^{(miss)}, \mathbf{x}^{(obs)})p(\mathbf{x}^{(miss)}|\mathbf{x}^{(obs)})d\theta,$$

leading to stationary equations, if $g(\theta) = p(\theta|\mathbf{x}^{(obs)})$:

$$g(\theta) = \int_{\phi \in \Theta} K(\theta, \phi)g(\phi)d\phi = \mathcal{T}g(\theta)$$

where $K(\theta, \phi) = \int_{z \in \mathcal{X}^{(miss)}} p(\theta|z, \mathbf{x}^{(obs)})p(z|\mathbf{x}^{(obs)}, \phi)dz$.

See for example Tanner and Wong (1987) on data augmentation and Rubin (1996) for global review.

Algorithms

Use Markov Chain Monte Carlo methods to converge to the posteriors such as at each iteration

$$\theta^* \sim p(\theta | \mathbf{x}^{(obs)}, \mathbf{x}^{(miss),*}),$$

$$\mathbf{x}^{(miss),*} \sim p(\mathbf{x}^{(miss)} | \theta^*, \mathbf{x}^{(obs)}),$$

and at ∞ , draws follow $p(\mathbf{x}^{(miss)}, \theta | \mathbf{x}^{(obs)})$.

Weakness?

Necessity to define a joint model. Hard for mixed data.

Solution? FAMD (PCA for mixed data, see Audigier, Husson, and Josse (2016)) treats categorical variables through dummy variables.

Other solution? Define conditional models such as $p(\mathbf{x}_1 | \mathbf{x}_2, \dots, \mathbf{x}_p, \theta_1)$ where variables can be qualitative/quantitative. Iterate per coordinate in cycles until convergence.

Chained Equations/Fully Conditional Specification (FCS)

And for supervised problems ?

Idea: Partially (?) Conditional Specification: cycles of

$$\begin{aligned}
 \boldsymbol{\theta}_x^{*(t)} &\sim p(\boldsymbol{\theta}_x | \mathbf{x}^{\text{obs}}, \mathbf{y}^{\text{obs}}, \mathbf{y}^{\text{miss}*(t-1)}), \\
 \mathbf{x}^{\text{miss}*(t)} &\sim p(\mathbf{x}^{(\text{miss})} | \mathbf{x}^{\text{obs}}, \mathbf{y}^{\text{obs}}, \mathbf{y}^{\text{miss}*(t-1)}, \boldsymbol{\theta}_x^{*(t)}), \\
 \boldsymbol{\theta}_y^{*(t)} &\sim p(\boldsymbol{\theta}_y | \mathbf{x}^{\text{obs}}, \mathbf{x}^{\text{miss}*(t)}, \mathbf{y}^{\text{obs}}), \\
 \mathbf{y}^{\text{miss}*(t)} &\sim p(\mathbf{y}^{(\text{miss})} | \mathbf{x}^{\text{obs}}, \mathbf{x}^{\text{miss}*(t)}, \mathbf{y}^{\text{obs}}, \boldsymbol{\theta}_y^{*(t)}),
 \end{aligned} \tag{1}$$

And for supervised problems ?

Idea: Partially (?) Conditional Specification: cycles of

$$\begin{aligned}
 \theta_x^{*(t)} &\sim p(\theta_x | \mathbf{x}^{\text{obs}}, \mathbf{y}^{\text{obs}}, \mathbf{y}^{\text{miss}*(t-1)}), \\
 \mathbf{x}^{\text{miss}*(t)} &\sim p(\mathbf{x}^{(\text{miss})} | \mathbf{x}^{\text{obs}}, \mathbf{y}^{\text{obs}}, \mathbf{y}^{\text{miss}*(t-1)}, \theta_x^{*(t)}), \\
 \theta_y^{*(t)} &\sim p(\theta_y | \mathbf{x}^{\text{obs}}, \mathbf{x}^{\text{miss}*(t)}, \mathbf{y}^{\text{obs}}), \\
 \mathbf{y}^{\text{miss}*(t)} &\sim p(\mathbf{y}^{(\text{miss})} | \mathbf{x}^{\text{obs}}, \mathbf{x}^{\text{miss}*(t)}, \mathbf{y}^{\text{obs}}, \theta_y^{*(t)}),
 \end{aligned} \tag{1}$$

Statistical model, close to Trygg and Wold (2003)

$$\begin{aligned}
 \mathbf{x} &= f_{\theta_x}(\mathbf{x}_y, \mathbf{x}_x), \\
 \mathbf{y} &= g_{\theta_y}(\mathbf{y}_x, \mathbf{y}_y), \\
 \mathbf{x} &\perp\!\!\!\perp \mathbf{y} | \mathbf{x}_y, \theta_x, \\
 \mathbf{y} &\perp\!\!\!\perp \mathbf{x} | \mathbf{y}_x, \theta_y,
 \end{aligned} \tag{2}$$

where $\mathbf{x}_x \perp\!\!\!\perp \mathbf{y}$, $\mathbf{x}_x \perp\!\!\!\perp \mathbf{x}_y$, $\mathbf{y}_y \perp\!\!\!\perp \mathbf{x}$ and $\mathbf{y}_y \perp\!\!\!\perp \mathbf{y}_x$.

And for supervised problems ? (II)

Observation: Impute \mathbf{x} (resp. \mathbf{y}) uses information from:

- \mathbf{x} (resp. \mathbf{y}) through \mathbf{x}_x (resp. \mathbf{y}_y) \rightarrow JM part,
- \mathbf{y} (resp. \mathbf{x}) through \mathbf{x}_y (resp. \mathbf{y}_x) \rightarrow FCS part.

Remark: If \mathbf{y} is univariate, \mathbf{y}_y does not exist.

And for supervised problems ? (II)

Observation: Impute \mathbf{x} (resp. \mathbf{y}) uses information from:

- \mathbf{x} (resp. \mathbf{y}) through \mathbf{x}_x (resp. \mathbf{y}_y) \rightarrow JM part,
- \mathbf{y} (resp. \mathbf{x}) through \mathbf{x}_y (resp. \mathbf{y}_x) \rightarrow FCS part.

Remark: If \mathbf{y} is univariate, \mathbf{y}_y does not exist.

Koh-Lanta

Hypothesis for this presentation: No missing values in \mathbf{y} .

Parameter variability: Bootstrap of the input dataset.

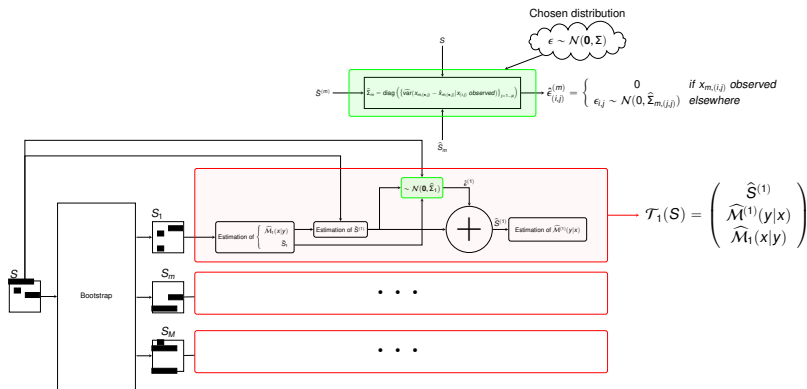
Distribution hypothesis: Multivariate Gaussian.

Simplifying assumption: \mathbf{x}_x and \mathbf{y}_y with respective diagonal variance matrices: only noise, not really realistic.

$$\begin{aligned}\mathbf{x} &= \mathbf{M}\mathbf{y} + \boldsymbol{\epsilon}, \\ \mathbf{y} &= \mathbf{N}\mathbf{x} + \boldsymbol{\xi}.\end{aligned}\tag{3}$$

$$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma} = \text{diag}(\sigma_j^2)_{j=1..p})$$

The Koh-Lanta algorithm for imputation



Which statistical model ?

Our approach uses only FCS part equivalents: needs of prediction models only.

We have chosen **ddsPLS** (see Lorenzo et al. (2021)):

- PLS based approach,
- Single penalization for sparsity per component (sparse both in \mathbf{x} and \mathbf{y}),
- Automatic fix the number of components (R) through bootstrap operations and minimizing $\bar{R}_B^2 - \bar{Q}_B^2$ while $\bar{Q}_B^2 > 0$,
- “**Koh-Lanta (in ddsPLS)**” : adapt to important noise context (see Cai and Liu (2011)),
- **LD** mode for **Low Dimensional** contexts: “**Koh-Lanta (in ddsPLS LD)**”

Competitive approaches

Two types of imputation:

- Single imputation: “**MEAN-PLS**” , “**NIPALS**” .
- Multiple imputation (MI): “**MI-NIPALS**” , “**missMDA-PLS**” .

Competitive approaches

Two types of imputation:

- Single imputation: “**MEAN-PLS**” , “**NIPALS**” .
- Multiple imputation (MI): “**MI-NIPALS**” , “**missMDA-PLS**” .
- “**NIPALS**” performs PLS analysis estimating covariance matrices only on the observed values (no real imputation, generalizable ?).

Competitive approaches

Two types of imputation:

- Single imputation: “**MEAN-PLS**” , “**NIPALS**” .
- Multiple imputation (MI): “**MI-NIPALS**” , “**missMDA-PLS**” .
- “**NIPALS**” performs PLS analysis estimating covariance matrices only on the observed values (no real imputation, generalizable ?).
- “**missMDA-PLS**” (see Josse, Pagès, and Husson (2011)) performs regularized PCA, function MIPCA. R is chosen through the function `estim_ncp` (cross-validation and minimizing MSE).

Model

Let us suppose the following latent variable model

$$x_j = \begin{cases} \sqrt{1 - \sigma^2} \phi_1 + \sigma \varepsilon_j & \text{for } j = 1 \dots p_1 \\ \sqrt{1 - \sigma^2} \phi_2 + \sigma \varepsilon_j & \text{for } j = p_1 + 1 \dots 2p_1 \\ \sqrt{1 - \sigma^2} \phi_3 + \sigma \varepsilon_j & \text{for } j = 2p_1 + 1 \dots 2p_1 + p_2 \\ \varepsilon_j & \text{for } j = 2p_1 + p_2 + 1 \dots 2p_1 + p_2 + p_3 \end{cases} \quad (4)$$

$$\mathbf{y} : \begin{cases} y_1 = \sqrt{1 - \sigma^2} \phi_1 + \sigma \xi_1 = y_1^* + \sigma \xi_1 \\ y_2 = \sqrt{1 - \sigma^2} (\phi_1 + 2\phi_2) / \sqrt{5} + \sigma \xi_2 = y_2^* + \sigma \xi_2 \\ y_3 = \xi_3 \end{cases} \quad (5)$$

where $\sigma = \sqrt{0.1} \approx 0.316$ and

$$(\phi_1, \phi_2, \phi_3, \varepsilon'_{1\dots p}, \xi'_{1\dots 3})' \sim \mathcal{N}(\mathbf{0}_{3+p+3}, \mathbb{I}_{3+p+3}).$$

Note that $\sqrt{\mathbb{E}(y_j - y_j^*)^2} = \begin{cases} \sigma \approx 0.316 & \text{if } j = 1, 2 \\ 1 & \text{if } j = 3 \end{cases}$

Simulation parameters

Idea: test both the impacts of high noise (structured and unstructured) and of sample size.

Simulation parameters

Idea: test both the impacts of high noise (structured and unstructured) and of sample size.

- $p_1 = 10$: information variance,
- $(p_2, p_3) \in \{(1, 1), (100, 500)\}$: low noise and high noise variances (structured and unstructured),
- $n \in \{100, 50, 20\}$: hard to very hard problems.

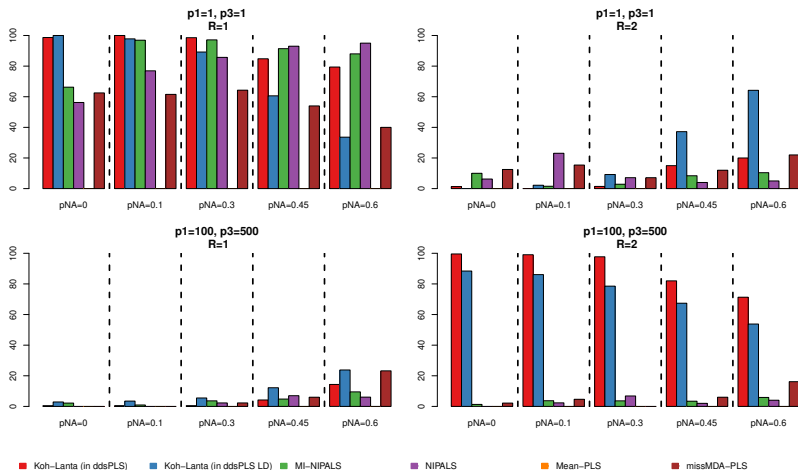
Simulation parameters

Idea: test both the impacts of high noise (structured and unstructured) and of sample size.

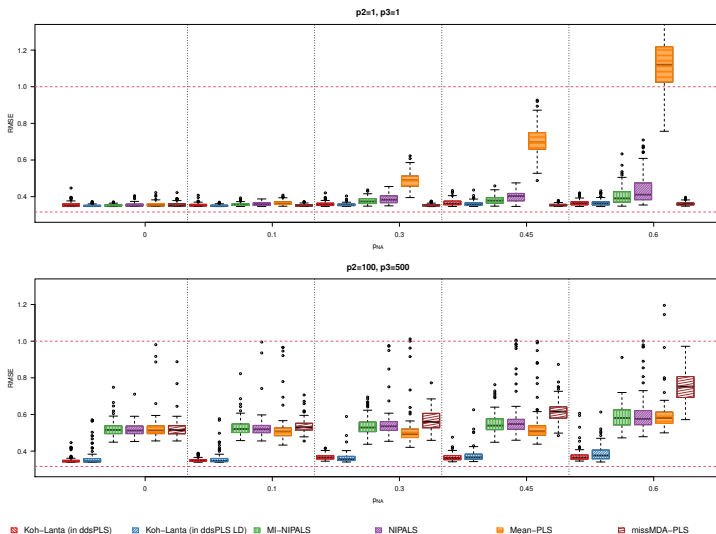
- $p_1 = 10$: information variance,
- $(p_2, p_3) \in \{(1, 1), (100, 500)\}$: low noise and high noise variances (structured and unstructured),
- $n \in \{100, 50, 20\}$: hard to very hard problems.

Remark: Use of factorial linear methods, the number of components R should be equal to $R = 2$.

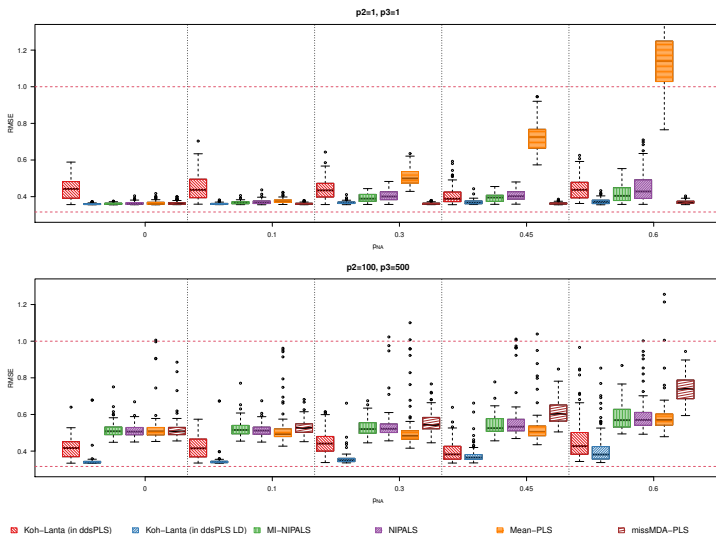
$n = 100, R$



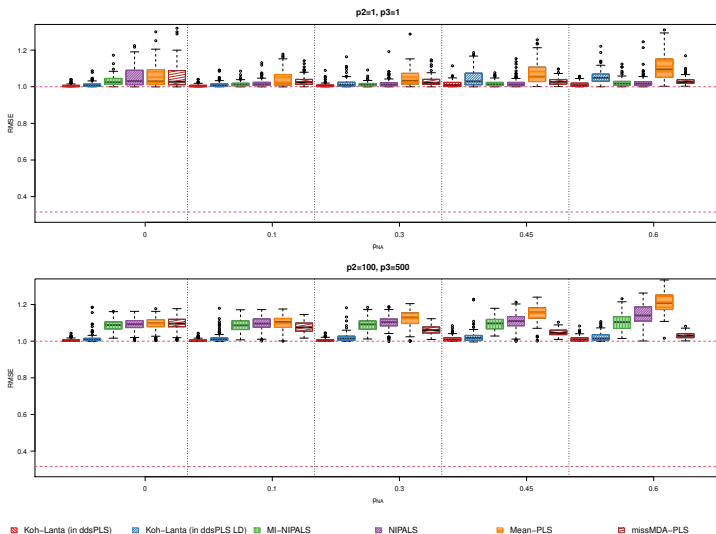
$n = 100$, RMSE for y_1



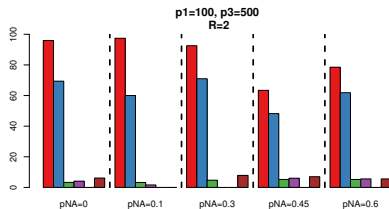
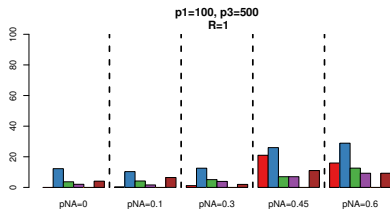
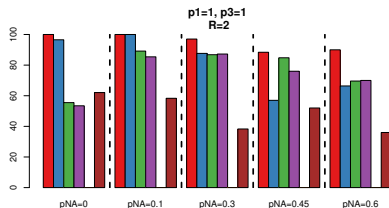
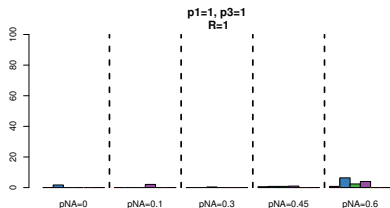
$n = 100$, RMSE for y_2



$n = 100$, RMSE for y_3

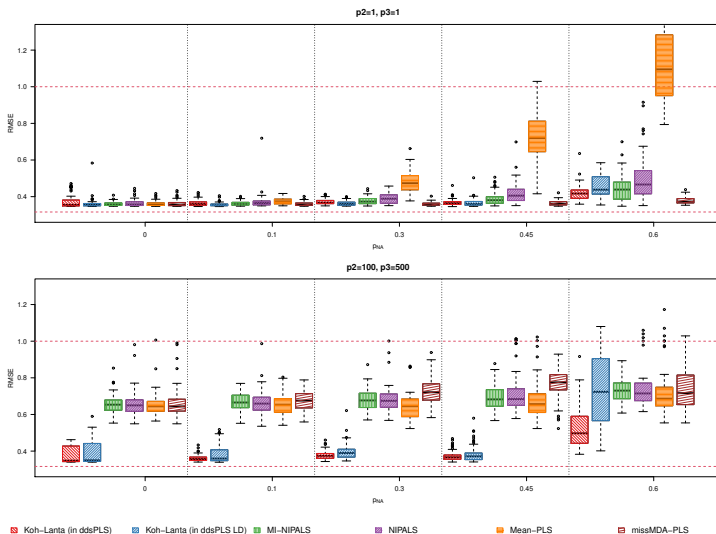


$n = 50, R$

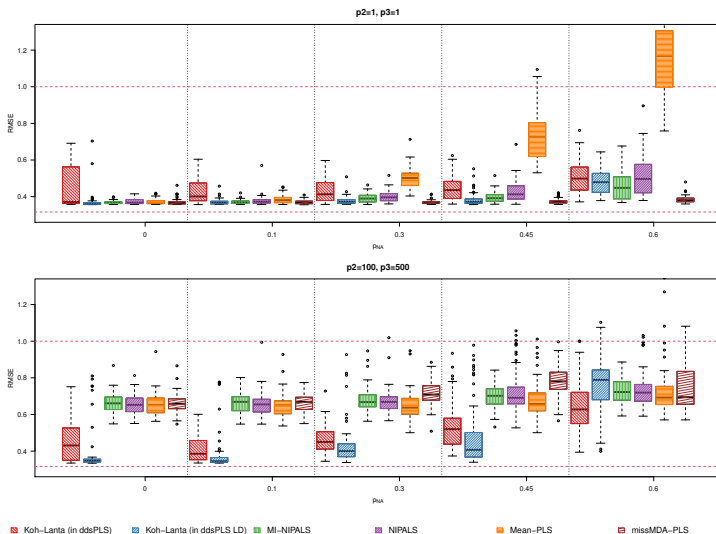


■ Koh-Lanta (in ddsPLS)
 ■ Koh-Lanta (in ddsPLS LD)
 ■ MI-NIPALS
 ■ NIPALS
 ■ Mean-PLS
 ■ missMDA-PLS

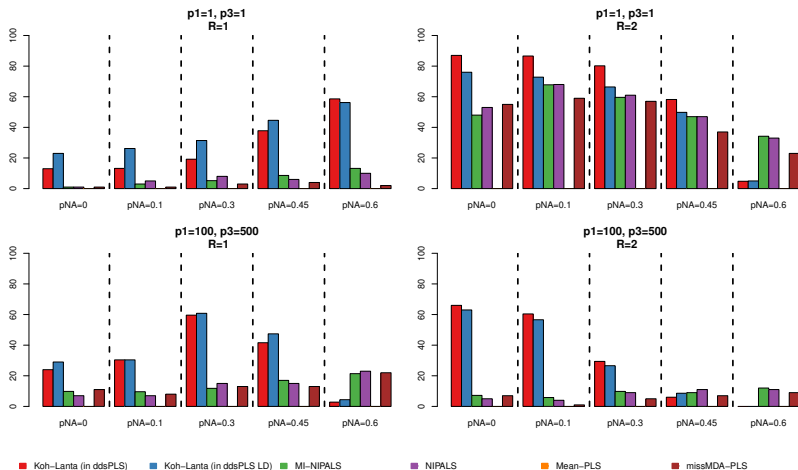
$n = 50$, RMSE for y_1

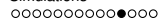


$n = 50$, RMSE for y_2

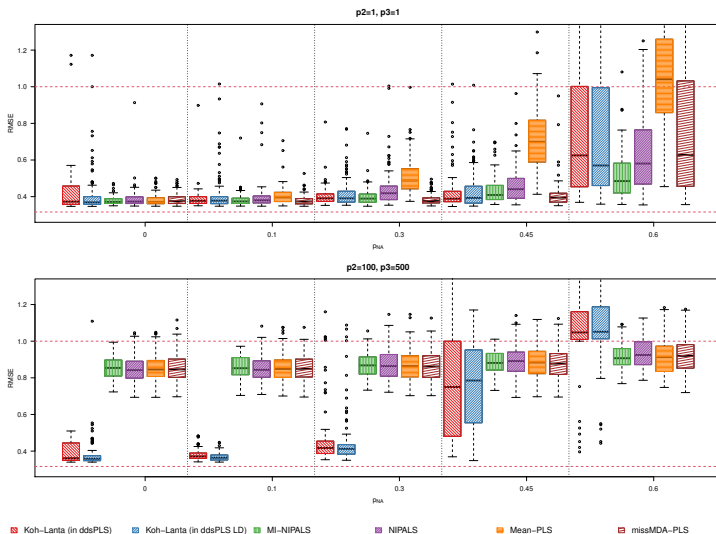


$n = 20, R$

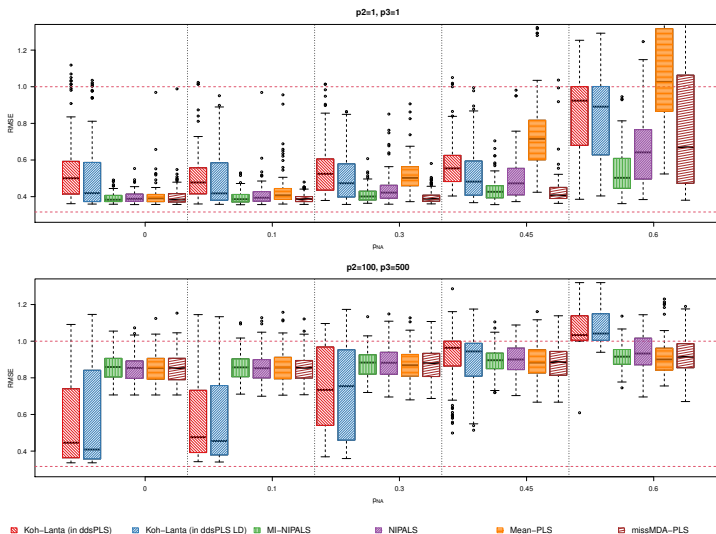




$n = 20$, RMSE for y_1



$n = 20$, RMSE for y_2



Summary

- “**MEAN-PLS**” suffers for large proportion of **NA**.

Summary

- “**MEAN-PLS**” suffers for large proportion of **NA**.
- “**missMDA-PLS**” is efficient for low noise but not efficient in high dimension.

Summary

- “**MEAN-PLS**” suffers for large proportion of **NA**.
- “**missMDA-PLS**” is efficient for low noise but not efficient in high dimension.
- “**Koh-Lanta (in ddsPLS)**” suffers for complex structures but “**Koh-Lanta (in ddsPLS LD)**” shows higher FPR in variable selection.

Summary

- “**MEAN-PLS**” suffers for large proportion of **NA**.
- “**missMDA-PLS**” is efficient for low noise but not efficient in high dimension.
- “**Koh-Lanta (in ddsPLS)**” suffers for complex structures but “**Koh-Lanta (in ddsPLS LD)**” shows higher FPR in variable selection.
- Competitors suffer more for high noise: no better than “**MEAN-PLS**” .

Summary

- “**MEAN-PLS**” suffers for large proportion of **NA**.
- “**missMDA-PLS**” is efficient for low noise but not efficient in high dimension.
- “**Koh-Lanta (in ddsPLS)**” suffers for complex structures but “**Koh-Lanta (in ddsPLS LD)**” shows higher FPR in variable selection.
- Competitors suffer more for high noise: no better than “**MEAN-PLS**” .
- $n = 20$ very hard but new methodologies suffer when they do not build models ($R = 0$).

Conclusion & future works

- Even in imputation, regularization helps for hard settings (low n and/or large p).
- Dig in "Partially Conditional Specification".
- Koh-Lanta is presented though **ddsPLS** but can be generalized to any supervised context.

References I

- [1] Vincent Audigier, François Husson, and Julie Josse. “A principal component method to impute missing values for mixed data”. In: *Advances in Data Analysis and Classification* 10.1 (2016), pp. 5–26.
- [2] Tony Cai and Weidong Liu. “Adaptive Thresholding for Sparse Covariance Matrix Estimation”. In: *Journal of the American Statistical Association* 106.494 (2011), pp. 672–684. doi: 10.1198/jasa.2011.tm10560. eprint: <https://doi.org/10.1198/jasa.2011.tm10560>. URL: <https://doi.org/10.1198/jasa.2011.tm10560>.
- [3] Julie Josse, Jérôme Pagès, and François Husson. “Multiple imputation in principal component analysis”. In: *Advances in data analysis and classification* 5.3 (2011), pp. 231–246.
- [4] Hadrien Lorenzo et al. “Data-driven sparse partial least squares”. In: *Statistical Analysis and Data Mining: The ASA Data Science Journal* (2021).
- [5] Anne Rechten et al. “Systems Vaccinology Identifies an Early Innate Immune Signature as a Correlate of Antibody Responses to the Ebola Vaccine rVSV-ZEBOV”. In: *Cell Reports* 20.9 (Sept. 2017), pp. 2251–2261. ISSN: 2211-1247. doi: 10.1016/j.celrep.2017.08.023.

References II

- [6] Donald B. Rubin. “Multiple Imputation After 18+ Years”. In: *Journal of the American Statistical Association* 91.434 (1996), pp. 473–489. ISSN: 01621459. URL: <http://www.jstor.org/stable/2291635>.
- [7] Martin A. Tanner and Wing Hung Wong. “The Calculation of Posterior Distributions by Data Augmentation”. In: *Journal of the American Statistical Association* 82.398 (1987), pp. 528–540. ISSN: 01621459. URL: <http://www.jstor.org/stable/2289457>.
- [8] Johan Trygg and Svante Wold. “O2-PLS, a two-block (X–Y) latent variable regression (LVR) method with an integral OSC filter”. In: *Journal of chemometrics* 17.1 (2003), pp. 53–64.