1 We would like to thank the reviewers for their thorough reading of the article and their many pertinent remarks, which
2 help to improve the clarity. In the final version, we will address all comments on form. To address the reviewers' main
3 concerns and better show the extent and feasibility of our methodology, we respond by adding an application on a
4 recommendation system data. We consider the Jester dataset [2] of 5000 users who rated jokes, with 27% of missing
5 values. The low-rank assumption for the loading matrix (allowed by Assump. 1) makes sense: any variable (i.e. user
6 preferences) can be expressed as a linear combination of $r$ latent variables[1] (hence, a "fully connected PPCA"). The
7 first latent variable opposes individuals who like jokes about physics but dislike jokes about sexuality, and conversely.

8 **MNAR mechanism. (R1, R3)**   Considering MNAR and self-masking values is plausible because users only rate
9 jokes they like or dislike strongly or might be ashamed to assume their taste for sexual jokes. Note that the **self-masked**
10 **assumption** is required only for the identifiability but the estimation strategy is also derived for **general MNAR**
11 **mechanisms** (allowed by Assump. A2) where the missingness may depend on other missing variables. Assump. A2
12 means that a user's non-response for the sexual joke given all jokes may depend on the scores of the sexual and physical
13 jokes but not on the musical and computer jokes.

14 **Selecting the number $r$ of latent variables and estimating the noise variance. (R1)**   To select $r$, one could use
15 complete observations only but this is not possible when the number of features is large. As an alternative, we used
16 both a cross-validation strategy assuming M(C)AR mechanism as detailed in [3] and also a beta implementation (that
17 we coded) of a CV assuming MNAR mechanism. The second one is dependent on the chosen mechanism. As noted by
18 the reviewers, Algorithm 1 is robust to a misspecification of the rank and thus a reasonable heuristic may already be
19 enough. Both approaches estimate $r = 5$. CV was also used for Traumabase where oracle values were only used for
20 synthetic data. With $r$ at hand, the noise variance is obtained directly using weighted residual sum of squares as in [3].

21 **Selecting the $r$ pivot variables. (R1,R3)**   The next step consists in selecting $r$ (M(C)AR) pivot variables (observed
22 or M(C)AR variables imply Assump. A4) on which regressions[2] are performed [3]. Here, because we do not have further
23 information on the missing mechanisms, we select the variables with the lowest missing rate. In Traumabase, the
24 selection was discussed with experts (doctors) who identified M(C)AR variables. To reduce the error committed by a
25 wrong selection of pivot variables, we suggest selecting a bigger set ($> 5$) and computing the final estimator with the
26 median of the estimators over all possible combinations. In Fig. 2, by discarding outliers, this **aggregation approach**
27 is more robust than selecting only $r$ pivot variables.

28 **Additional experiments. (R1,R4)**   Then, we test our method by introducing additional MNAR values on one variable
29 (containing 33% NA) using a self-masked mechanism leading to 65% NA. In Fig. 1, **our method (`MNAR`) outperforms**
30 **all the others on rating data** including `Deep` [1] which imputes MNAR values using deep generative models (R4)[4].
31 The parametric method `MNARparam` is not displayed as it does not scale on such large data. The code for the whole
32 methodology was already available, but now recast as a **beta version of package** (R1) and submitted soon on CRAN.
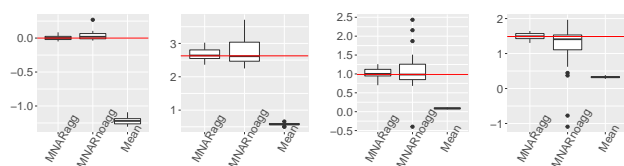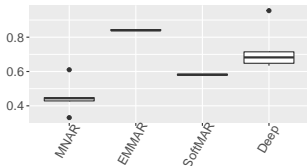
33


Fig. 1: Prediction error (difference between true values
and predicted ones) for the Jester dataset, the mean imputa-
tion corresponding to an error of 1. The process of drawing
additional MNAR values and predicting them is repeated
10 times which gives the stochasticity.



Fig. 2: Synthetic data from Section 4.1, with Algorithm 1 performed with aggregation (`MNARagg`) or not
(`MNARnoagg`). True values in red, estimated values (means, variance, cov) in boxplot. For a given set of
PPCA parameter, the stochasticity comes from the process of drawing 20 times the latent variables, the
additive noise and the missing-data pattern (R3).

34 **Comparison with Miao et al. [4] (R2,R3)**   For one variable $Y \sim \mathcal{N}(\mu, \sigma^2)$, Miao et al. prove identifiability of the
35 variance and the absolute value of the mean, assuming a self-masked mechanism with a known strictly monotone form
36 (including classical Probit and Logit). They cannot get identifiability for the mean (not the absolute value) with Logit.
37 We have used their result to prove variance identifiability in PPCA and provide a genuine proof for the mean without
38 discarding Logit. Secondly, for a specific setting of an heteroscedastic regression model with missing values only in $Y$,
39 where the variance of $Y$ given the observed covariates is injective, they provide identifiability results for the conditional
40 distribution with general MNAR. This setting and the proof are too restricted to be considered in PPCA.

41 **Supervised learning task on Traumabase. (R3)**   To predict the administration or not of
42 the tranexomic acid (binary variable), we impute explanatory variables before proceeding
43 to the classification task. In Tab. 1, our method gives the smallest prediction error.

| | |
|---|---|
| MNAR | 5.06% |
| EMMAR | 5.82% |
| SoftMAR | 5.45% |
| MNARparam | 5.39% |
| Mean | 5.27% |

Tab. 1: Mean of prediction error
over 10 repetitions.

44 [1] L. Gondara and K. Wang. Mida: Multiple imputation using denoising autoencoders. In *PAKDD*, 2018.
45 [2] M. Hahsler. recommenderlab: A framework for developing and testing recommendation algorithms. Technical report, 2015.
46 [3] J. Josse and F. Husson. Selecting the number of components in principal component analysis using cross-validation approximations. *CS&DA*, 2012.
47 [4] Wang Miao, Peng Ding, and Zhi Geng. Identifiability of normal and normal mixture models with nonignorable missing data. *JASA*, 2016.

[1] It does not require that the linear combination coefficients are non zero.   [2] assumed to be consistent by Assump. A3, which holds
as the noise tends to 0. (R1)   [3] Note that our method is not based on the complete-case of the dataset but on the complete-case of
the $r$ ($\ll p$) pivot variables (R3).   [4] Note that this method requires to be trained on a complete dataset.