

# Multiple imputation with principal component methods

Vincent Audigier<sup>1</sup>, François Husson<sup>2</sup>, Julie Josse<sup>3</sup>

1. CNAM, MSDMA team, Paris
2. Agrocampus Ouest, Rennes
3. Ecole Polytechnique, Paris

Telecom ParisTech, November 16, 2017

- ① Introduction
- ② Single imputation based on principal component methods
- ③ Multiple imputation for categorical data with MCA
- ④ Conclusion

# Missing values

NA					NA	NA	
			NA				
	NA						NA
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
		NA	NA	NA			

- Aim: inference on a quantity  $\theta$  from incomplete data  
→ point estimate  $\hat{\theta}$  and associated variability  $T$
- Complete-case analysis
- Likelihood approaches
- Multiple Imputation



# Principal component methods: aims

From multidimensional data, principal component methods:

- summarize
- describe
- visualize

Several objectives:

- identify the similarities between individuals
- highlight the relationships between variables
- describe some groups of individuals by a set of relevant variables

# How does it work?

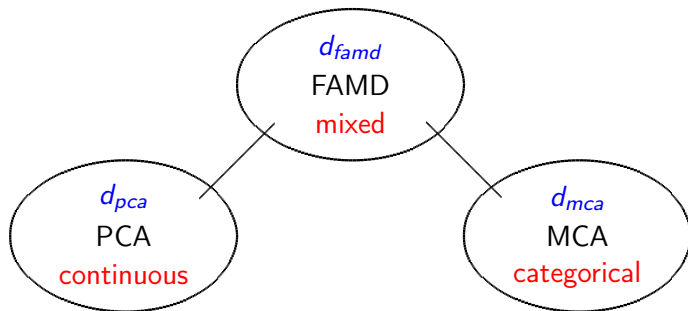
- The  $I$  individuals are seen as elements of  $\mathbb{R}^K$



- A distance  $d$  on  $\mathbb{R}^K$  to define proximities between individuals
- $\text{Vect}(v_1, \dots, v_S)$  maximising the projected inertia

⇒ Dimensionality reduction methods

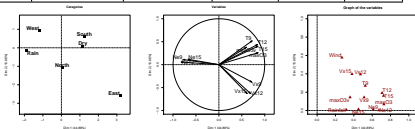
# A set of methods



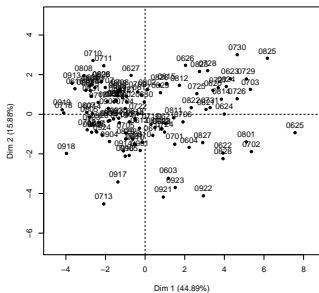
$$d_{famd}^2 = d_{pca}^2 + d_{mca}^2$$

## FAMD

	Wind	Rainfall	maxO3	T9	T12	...
0602	North	Dry	82	17.0	18.4	...
0603	East	Dry	92	15.3	17.6	...
0604	North	Dry	114	16.2	19.7	...
0605	West	Dry	94	17.4	20.5	...
0606	West	Rain	80	17.7	19.8	...
...	...	...	...	...	...	...



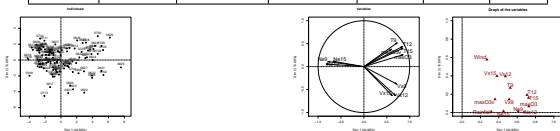
## Individuals



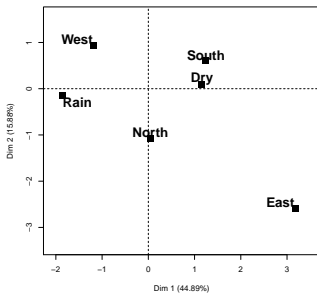


## FAMD

	Wind	Rainfall	maxO3	T9	T12	...
0602	North	Dry	82	17.0	18.4	...
0603	East	Dry	92	15.3	17.6	...
0604	North	Dry	114	16.2	19.7	...
0605	West	Dry	94	17.4	20.5	...
0606	West	Rain	80	17.7	19.8	...
...	...	...	...	...	...	...

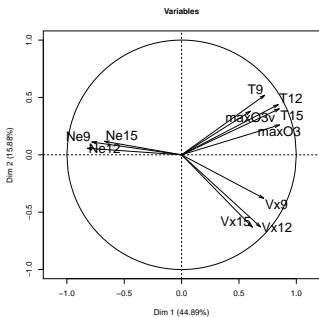
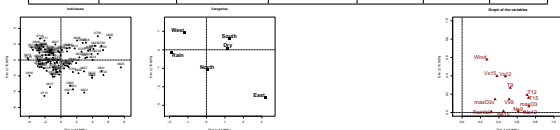


Categories



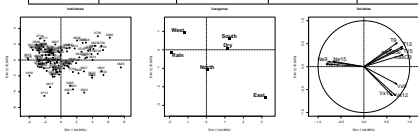
## FAMD

	Wind	Rainfall	maxO3	T9	T12	...
0602	North	Dry	82	17.0	18.4	...
0603	East	Dry	92	15.3	17.6	...
0604	North	Dry	114	16.2	19.7	...
0605	West	Dry	94	17.4	20.5	...
0606	West	Rain	80	17.7	19.8	...
...	...	...	...	...	...	...

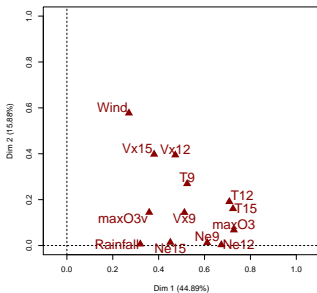


## FAMD

	Wind	Rainfall	maxO3	T9	T12	...
0602	North	Dry	82	17.0	18.4	...
0603	East	Dry	92	15.3	17.6	...
0604	North	Dry	114	16.2	19.7	...
0605	West	Dry	94	17.4	20.5	...
0606	West	Rain	80	17.7	19.8	...
...	...	...	...	...	...	...



Graph of the variables



# Why principal component methods for MI?

- ① Lack of imputation models for non-continuous data
  - mixed data: various relationships between variables
  - categorical data: combinatorial explosion

A large range of models

- ② Imputation models with too many parameters:
  - overfitting
  - storage issue

A small number of parameters

- ③ Models based on regressions:
  - number of individuals less than the number of variables
  - collinearity

No inversion issues

- ① Introduction
- ② Single imputation based on principal component methods
- ③ Multiple imputation for categorical data with MCA
- ④ Conclusion

# How to perform FAMD?

FAMD can be seen as the SVD of  $\mathbf{X}$  with weights for

- the continuous variables and categories:  $(\mathbf{D}_\Sigma)^{-1}$
- the individuals:  $\frac{1}{I} \mathbb{1}_I$

$$\rightarrow SVD \left( \mathbf{X}, (\mathbf{D}_\Sigma)^{-1}, \frac{1}{I} \mathbb{1}_I \right)$$

$X =$

11.04	...	2.07	1	0	...	1	0	0
10.76	...	1.86	1	0	...	1	0	0
11.02	...	2.04	1	0	...	1	0	0
11.02	...	1.92	0	1	...	0	1	0
11.06		2.01	0	1		0	0	1
10.95		1.67	0	1		0	1	0

$D_\Sigma =$

$\sigma_{x_1}$			
	$\ddots$		
		$\sigma_{x_k}$	
			0
0			$I_{k+1}$
			$\ddots$
			$I_K$

# How to perform FAMD?

$$SVD \left( \mathbf{X}, (\mathbf{D}_\Sigma)^{-1}, \frac{1}{I} \mathbb{1}_I \right) \longrightarrow \mathbf{X}_{I \times K} = \mathbf{U}_{I \times K} \mathbf{\Lambda}_{K \times K}^{1/2} \mathbf{V}_{K \times K}^\top$$

with  $\mathbf{U}^\top \left( \frac{1}{I} \mathbb{1}_I \right) \mathbf{U} = \mathbb{1}_K$   
 $\mathbf{V}^\top \mathbf{D}_\Sigma^{-1} \mathbf{V} = \mathbb{1}_K$

- principal components:  $\hat{\mathbf{F}}_{I \times S} = \hat{\mathbf{U}}_{I \times S} \hat{\mathbf{\Lambda}}_{S \times S}^{1/2}$
- loadings:  $\hat{\mathbf{V}}_{K \times S}^\top$
- fitted matrix:  $\hat{\mathbf{X}}_{I \times K} = \hat{\mathbf{U}}_{I \times S} \hat{\mathbf{\Lambda}}_{S \times S}^{1/2} \hat{\mathbf{V}}_{K \times S}^\top$

$$\| \hat{\mathbf{X}} - \mathbf{X} \|_{\mathbf{D}_\Sigma^{-1} \otimes \frac{1}{I} \mathbb{1}}^2 = \text{tr} \left( (\hat{\mathbf{X}} - \mathbf{X}) \mathbf{D}_\Sigma^{-1} (\hat{\mathbf{X}} - \mathbf{X})^\top \frac{1}{I} \mathbb{1}_I \right)$$

minimized under the constraint of rank  $S$

# FAMD with missing values

⇒ FAMD: minimize

$$\|\mathbf{X}_{I \times K} - \hat{\mathbf{X}}_{I \times K}\|_{\mathbf{D}_{\Sigma}^{-1} \otimes \frac{1}{J} \mathbf{1}}^2$$

⇒ FAMD with missing values:

$$\|\mathbf{W}_{I \times K} * (\mathbf{X}_{I \times K} - \hat{\mathbf{X}}_{I \times K})\|_{\mathbf{D}_{\Sigma}^{-1} \otimes \frac{1}{J} \mathbf{1}}^2$$

with  $w_{ij} = 0$  if  $x_{ij}$  is missing,  $w_{ij} = 1$  otherwise

Many algorithms developed for PCA such as NIPALS (Christoffersson, 1970) or iterative PCA (Kiers, 1997)



# FAMD with missing values

Iterative FAMD algorithm:

- 1 initialization: imputation by mean/proportion
- 2 iterate until convergence
  - (a) estimation of the parameters of FAMD  
 $\rightarrow$  SVD of  $(\mathbf{X}, (\mathbf{D}_\Sigma)^{-1}, \frac{1}{I} \mathbb{1}_I)$
  - (b) imputation of the missing values with  
 $\hat{\mathbf{X}}_{I \times K} = \hat{\mathbf{U}}_{I \times S} \hat{\Lambda}_{S \times S}^{1/2} \hat{\mathbf{V}}_{K \times S}^\top$
  - (c)  $\mathbf{D}_\Sigma$  is updated

NA	...	2.07	A	...	A
10.76	...	1.86	A	...	A
11.02	...	NA	A	...	NA
11.02	...	1.92	B	...	B
11.06		2.01	NA	...	C
NA		1.67	B	...	B

$\rightarrow$

NA	...	2.07	1	0	...	1	0	0
10.76	...	1.86	1	0	...	1	0	0
11.02	...	NA	1	0	...	NA	NA	NA
11.02	...	1.92	0	1	...	0	1	0
11.06		2.01	NA	NA		0	0	1
NA		1.67	0	1		0	1	0

# FAMD with missing values

Iterative FAMD algorithm:

- 1 initialization: imputation by mean/proportion
- 2 iterate until convergence
  - (a) estimation of the parameters of FAMD  
 $\rightarrow$  SVD of  $(\mathbf{X}, (\mathbf{D}_\Sigma)^{-1}, \frac{1}{I} \mathbb{1}_I)$
  - (b) imputation of the missing values with  
 $\hat{\mathbf{X}}_{I \times K} = \hat{\mathbf{U}}_{I \times S} \hat{\Lambda}_{S \times S}^{1/2} \hat{\mathbf{V}}_{K \times S}^\top$
  - (c)  $\mathbf{D}_\Sigma$  is updated

NA	...	2.07	A	...	A
10.76	...	1.86	A	...	A
11.02	...	NA	A	...	NA
11.02	...	1.92	B	...	B
11.06		2.01	NA	...	C
NA		1.67	B	...	B

$\rightarrow$

11.01	...	2.07	1	0	...	1	0	0
10.76	...	1.86	1	0	...	1	0	0
11.02	...	1.89	1	0	...	0.61	0.19	0.20
11.02	...	1.92	0	1	...	0	1	0
11.06		2.01	0.32	0.68		0	0	1
11.01		1.67	0	1		0	1	0

# FAMD with missing values

Iterative FAMD algorithm:

- 1 initialization: imputation by mean/proportion
- 2 iterate until convergence
  - (a) estimation of the parameters of FAMD  
 $\rightarrow$  SVD of  $(\mathbf{X}, (\mathbf{D}_\Sigma)^{-1}, \frac{1}{I} \mathbb{1}_I)$
  - (b) imputation of the missing values with  
 $\hat{\mathbf{X}}_{I \times K} = \hat{\mathbf{U}}_{I \times S} \hat{\Lambda}_{S \times S}^{1/2} \hat{\mathbf{V}}_{K \times S}^\top$
  - (c)  $\mathbf{D}_\Sigma$  is updated

NA	...	2.07	A	...	A
10.76	...	1.86	A	...	A
11.02	...	NA	A	...	NA
11.02	...	1.92	B	...	B
11.06	...	2.01	NA	...	C
NA	...	1.67	B	...	B

$\rightarrow$

11.04	...	2.07	1	0	...	1	0	0
10.76	...	1.86	1	0	...	1	0	0
11.02	...	2.04	1	0	...	0.81	0.05	0.14
11.02	...	1.92	0	1	...	0	1	0
11.06	...	2.01	0.25	0.75	...	0	0	1
10.95	...	1.67	0	1	...	0	1	0

# Single imputation with FAMD (Audigier et al., 2016)

Iterative FAMD algorithm:

- 1 initialization: imputation by mean/proportion
- 2 iterate until convergence
  - (a) estimation of the parameters of FAMD  
 $\rightarrow$  SVD of  $(\mathbf{X}, (\mathbf{D}_\Sigma)^{-1}, \frac{1}{I} \mathbb{1}_I)$
  - (b) imputation of the missing values with  
 $\hat{\mathbf{X}}_{I \times K} = \hat{\mathbf{U}}_{I \times S} \hat{\Lambda}_{S \times S}^{1/2} \hat{\mathbf{V}}_{K \times S}^\top$
  - (c)  $\mathbf{D}_\Sigma$  is updated

11.04	...	2.07	A	...	A
10.76	...	1.86	A	...	A
11.02	...	2.04	A	...	A
11.02	...	1.92	B	...	B
...	...	...	...	...	...
11.06	...	2.01	B	...	C
10.95	...	1.67	B	...	B



11.04	...	2.07	1	0	...	1	0	0
10.76	...	1.86	1	0	...	1	0	0
11.02	...	2.04	1	0	...	0.81	0.05	0.14
11.02	...	1.92	0	1	...	0	1	0
...	...	...	...	...	...	...	...	...
11.06	...	2.01	0.25	0.75	...	0	0	1
10.95	...	1.67	0	1	...	0	1	0

$\Rightarrow$  the imputed values can be seen as degree of membership

# Single imputation with FAMD (Audigier et al., 2016)

Iterative FAMD algorithm:

- 1 initialization: imputation by mean/proportion
- 2 iterate until convergence
  - (a) estimation of the parameters of FAMD  
 $\rightarrow$  SVD of  $(\mathbf{X}, (\mathbf{D}_\Sigma)^{-1}, \frac{1}{J} \mathbb{1}_I)$
  - (b) imputation of the missing values with  
 $\hat{\mathbf{X}}_{I \times K} = \hat{\mathbf{U}}_{I \times S} f(\hat{\lambda}_{S \times S}^{1/2}) \hat{\mathbf{V}}_{K \times S}^\top$
  - (c)  $\mathbf{D}_\Sigma$  is updated

$$f(\hat{\lambda}_s^{1/2}) = \hat{\lambda}_s^{1/2} - \frac{\hat{\sigma}^2}{\hat{\lambda}_s^{1/2}}$$

11.04	...	2.07	A	...	A
10.76	...	1.86	A	...	A
11.02	...	2.04	A	...	A
11.02	...	1.92	B	...	B
...	...	...	...	...	...
11.06	...	2.01	B	...	C
10.95	...	1.67	B	...	B

←

11.04	...	2.07	1	0	...	1	0	0
10.76	...	1.86	1	0	...	1	0	0
11.02	...	2.04	1	0	...	0.81	0.05	0.14
11.02	...	1.92	0	1	...	0	1	0
...	...	...	...	...	...	...	...	...
11.06	...	2.01	0.25	0.75	...	0	0	1
10.95	...	1.67	0	1	...	0	1	0

$\Rightarrow$  the imputed values can be seen as degree of membership

# Imputation with Random Forests (Bühlmann & Stekhoven, 2011)

A random forest:

- Each tree is built from a subset of observations
- Fitted values are pooled

Imputation with RF:

- 1 Initial imputation: mean imputation - frequent category
- 2 Fit a RF  $X_j^{obs}$  to the other variables  $X_{-j}^{obs}$   
Predict  $X_j^{miss}$  using the trained R on  $X_{-j}^{miss}$
- 3 Cycle through variables
- 4 Repeat step 2 and 3 until convergence

Implemented: R package `missForest` (Stekhoven)

# Simulation study

Several data sets

- relationships between variables
- number of categories
- percentage of missing values (10%,20%,30%)

Criteria:

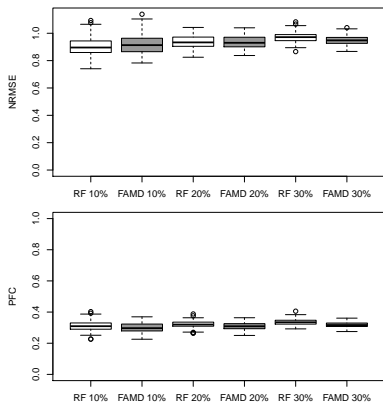
- for continuous data:

$$NRMSE = \sqrt{\frac{\sum_k \sum_{i=1}^I (1 - w_{ik}) \left( \frac{x_{ik} - \hat{x}_{ik}}{\hat{\sigma}_{x_k}} \right)^2}{\sum_k \sum_{i=1}^I (1 - w_{ik})}}$$

- for categorical data: proportion of falsely classified entries

# Comparisons from real data sets

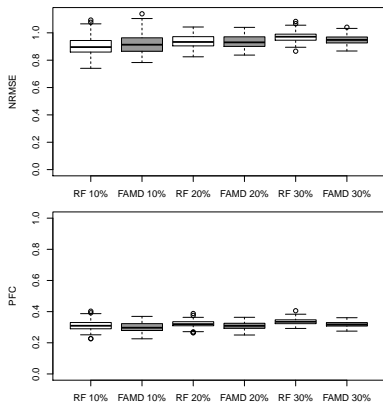
## GBSG2



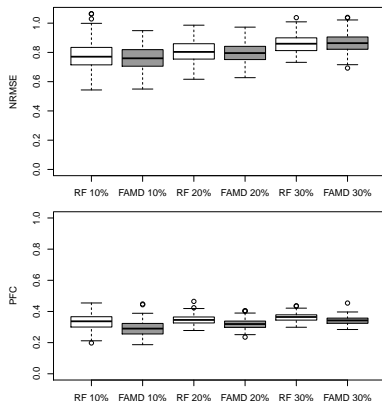


# Comparisons from real data sets

## GBSG2



## Tips



# Simulation results

Single imputation with FAMD shows a high quality of prediction compared to random forests (Stekhoven and Bühlmann, 2012)

- on real data
- when the relationships between continuous variables are linear
- for rare categories
- with MAR/MCAR mechanism

Can impute mixed, continuous or categorical data

# Simulation results

Single imputation with FAMD shows a high quality of prediction compared to random forests (Stekhoven and Bühlmann, 2012)

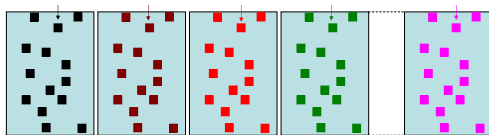
- on real data
- when the relationships between continuous variables are linear
- for rare categories
- with MAR/MCAR mechanism

Can impute mixed, continuous or categorical data

But a single imputation method only

# From single imputation to multiple imputation

$$P(X^{miss}|X^{obs}, \psi_1) \dots P(X^{miss}|X^{obs}, \psi_M)$$



- 1 Reflect the variability on the parameters of the imputation model

$$\rightarrow \left( \left( \hat{\mathbf{U}}_{I \times S}, \hat{\Lambda}_{S \times S}^{1/2}, \hat{\mathbf{V}}_{K \times S}^T \right)_1, \dots, \left( \hat{\mathbf{U}}_{I \times S}, \hat{\Lambda}_{S \times S}^{1/2}, \hat{\mathbf{V}}_{K \times S}^T \right)_M \right)$$

Bayesian or Bootstrap

- 2 Add a disturbance on the prediction by  $\hat{X}_m = \hat{\mathbf{U}}_m \hat{\Lambda}_m^{1/2} \hat{\mathbf{V}}_m^T$   
 → need to distinguish continuous and categorical data

- ① Introduction
- ② Single imputation based on principal component methods
- ③ Multiple imputation for categorical data with MCA
- ④ Conclusion

# Stochastic single imputation

$\pi_b$	0.4
$\pi_a$	0.6
$\pi_{b A}$	0.2
$\pi_{a A}$	0.8
$\pi_{a B}$	0.4
$\pi_{b B}$	0.6

→

$V_1$	$V_2$
A	a
B	b
B	a
B	b
$\vdots$	$\vdots$

→

$V_1$	$V_2$
A	a
B	NA
B	a
B	NA
$\vdots$	$\vdots$

**MCA majority**

$\pi_{b A}$	0.14
$\pi_{a A}$	0.86
$\pi_{a B}$	0.27
$\pi_{b B}$	0.73

**MCA draw**

$\pi_{b A}$	0.18
$\pi_{a A}$	0.82
$\pi_{a B}$	0.41
$\pi_{b B}$	0.59

$$\text{cov}_{95\%}(\pi_b) = 51.5 \quad \text{cov}_{95\%}(\pi_b) = 89.9$$

⇒ Standard errors of the parameters ( $\hat{\sigma}_{\hat{\pi}_b}$ ) calculated from the imputed data set are underestimated



# Properties

MI for categorical data is **challenging** for a moderate number of variables

- estimation issues
- storage issues

MCA address the categorical data challenge by

- requiring a small number of parameters
- preserving the essential data structure
- using a regularisation strategy

MIMCA can be applied on various data sets

- small or large number of variables/categories
- small or large number of individuals



# MI using the loglinear model (Schafer, 1997)

- Hypothesis on  $X = (x_{ijk})_{i,j,k}$ :  $X|\psi \sim \mathcal{M}(n, \psi)$

$$\log(\psi_{ijk}) = \lambda_0 + \lambda_i^A + \lambda_j^B + \lambda_k^C + \lambda_{ij}^{AB} + \lambda_{ik}^{AC} + \lambda_{jk}^{BC} + \lambda_{ijk}^{ABC}$$

- 1 Variability of the parameter  $\psi$ : Bayesian formulation
  - 2 Imputation using the set of  $M$  parameters
- Implemented: R package `cat` (J.L. Schafer)

Properties:

- Captures all the data relationships
- A number of parameters very large  $\rightarrow$  fails on large data sets

# MI using a latent class model (Si and Reiter, 2013)

- Hypothesis:  $\mathbb{P}(X = (x_1, \dots, x_K); \psi) = \sum_{\ell=1}^L \left( \psi_{\ell} \prod_{k=1}^K \psi_{x_k}^{(\ell)} \right)$
- ① Variability of the parameters  $\psi_L$  and  $\psi_X$  : Bayesian formulation
- ② Imputation using the set of  $M$  parameters
- Implemented: R package `mi` (Gelman *et al.*)

## Properties:

- Local independence assumption
- Captures complex relationships
- A small number of parameters

## Conditional modelling (van Buuren, 2006)

- A recent one: one **random forest**/variable (Doove *et al.*, 2014)

Properties:

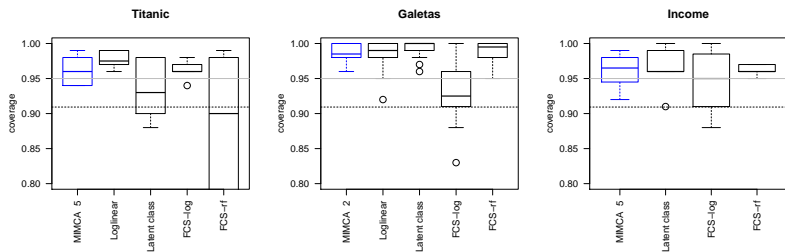
- non-parametric modelling
  - captures complex relationships between variables
- 
- A standard one: one **logistic regression** model/variable without interaction

Properties: captures relationships between pairs of variables

# Simulations from real data sets

- Quantities of interest:  $\theta$  = parameters of a logistic model
- Simulation design (repeated 200 times)
  - the real data set is considered as a population
  - drawn one sample from the data set
  - generate 20% of missing values
  - multiple imputation using  $M = 5$  imputed arrays
- Criteria
  - bias
  - CI width, coverage
- Comparison with :
  - JM: log-linear model, latent class model
  - FCS: logistic regression, random forests

# Results - Inference



	Titanic	Galetas	Income
Number of variables	4	4	14
Number of categories	$\leq 4$	$\leq 11$	$\leq 9$

## Results - Time

	Titanic	Galetas	Income
MIMCA	2.750	8.972	<b>58.729</b>
Loglinear	0.740	4.597	NA
Latent class model	10.854	17.414	143.652
FCS logistic	4.781	38.016	881.188
FCS forests	265.771	112.987	6329.514

Table: Time consumed in second

	Titanic	Galetas	Income
Number of individuals	2201	1192	6876
Number of variables	4	4	14

# Conclusion

MI methods using dimensionality reduction method

- captures the relationships between variables
- captures the similarities between individuals
- requires a small number of parameters

Address some imputation issues:

- can be applied on various data sets
- provide correct inferences for analysis model based on relationships between pairs of variables

Available in the R package missMDA, with a user guide on <http://vincentaudigier.weebly.com/links.html>

# Perspectives

- MI for mixed data using FAMD
- uncertainty on the number of dimensions  $S$
- promising methods for very large datasets
  - dimensionality reduction methods
  - SVD can be quickly performed on large matrices
  - should we preferably use single imputation?



# References I

- V. Audigier, F. Husson, and J. Josse. Mimca: multiple imputation for categorical variables with multiple correspondence analysis. *Statistics and Computing*, 27(2):501–518, 2017. ISSN 1573-1375.
- V. Audigier, F. Husson, and J. Josse. Multiple imputation for continuous variables using a bayesian principal component analysis. *Journal of Statistical Computation and Simulation*, 86(11):2140–2156, 2016a.
- V. Audigier, F. Husson, and J. Josse. A principal component method to impute missing values for mixed data. *Advances in Data Analysis and Classification*, 10(1):5–26, 2016b.
- D. B. Rubin. *Multiple Imputation for Non-Response in Survey*. Wiley, New-York, 1987.
- J. L. Schafer. *Analysis of Incomplete Multivariate Data*. Chapman & Hall/CRC, London, 1997.

## How to choose the number of dimensions?

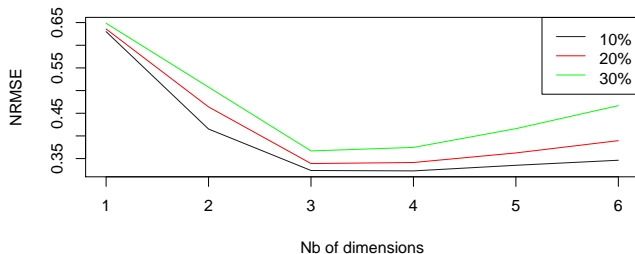
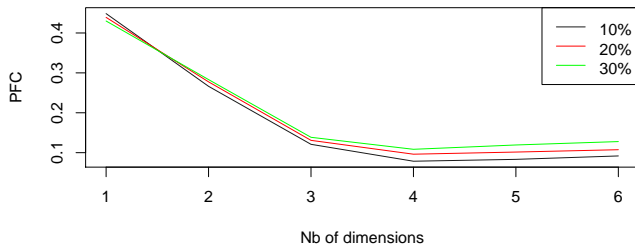
By cross-validation procedures:

- adding missing values on the incomplete data set
- predicting each of them using FAMD for several number of dimensions
- calculating the prediction error

Several ways:

- Leave-one-out (Bro et al., 2008)
- Repeated cross-validation

# Misspecification of the number of dimensions



# Single imputation MAR

- A mixed data set is simulated by splitting normal data
- Missing values are added on one variable  $Y$  according to a MAR mechanism:  $\mathbb{P}(Y = NA) = \frac{\exp(\beta_0 + \beta_1 X_1)}{1 + \exp(\beta_0 + \beta_1 X_1)}$
- Data are imputed using FAMD and RF

