

# Handling missing values with regularized iterative multiple correspondence analysis

J. Josse<sup>1</sup>, M. Chavent<sup>2</sup>, B. Liquet<sup>3</sup> & F. Husson<sup>1</sup>

<sup>1</sup> Agrocampus Rennes

<sup>2</sup> University of Bordeaux 2

<sup>3</sup> ISPED Bordeaux

ICC conference, St Andrew, 11 July 2011

## Handling missing values in PCA

⇒ Minimization of:

$$\mathcal{C} = \|\mathbf{X}_{I \times J} - \mathbf{F}_{I \times S} \mathbf{U}_{S \times J}^t\|^2$$

⇒ With missing values:

$$\mathcal{C} = \|\mathbf{W} * (\mathbf{X} - \mathbf{F}\mathbf{U}^t)\|^2,$$

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

⇒ Iterative PCA (Kiers, 1997)

## Iterative PCA

- 1 initialization  $\ell = 0$ :  $\mathbf{X}^0$  (mean imputation)
- 2 step  $\ell$ :
  - (a) PCA is performed on the completed data set  $\rightarrow (\hat{\mathbf{F}}^\ell, \hat{\mathbf{U}}^\ell)$ ;  $S$  dimensions are kept
  - (b) missing values are imputed with the model matrix  $\hat{\mathbf{X}}^\ell = \hat{\mathbf{F}}^\ell \hat{\mathbf{U}}^{\ell'}$ ; the new imputed dataset is  $\mathbf{X}^\ell = \mathbf{W} * \mathbf{X} + (1 - \mathbf{W}) * \hat{\mathbf{X}}^\ell$
- 3 steps (a) and (b) are repeated until convergence.

$\Rightarrow$  The number of dimensions  $S$  has to be chosen *a priori*

$\Rightarrow$  EM algorithm of  $x_{ij} = \sum_{s=1}^S f_{is} u_{js} + \varepsilon_{ij}, \varepsilon_{ij} \sim \mathcal{N}(0, \sigma^2)$

$\Rightarrow$  Nora-Chouteau in CA (1974)



## Iterative MCA

- 1 initialization  $\ell = 0$ :  $\mathbf{X}^0$  missing values are imputed with the proportion of the category  $\frac{I_k}{I}$  (the sum must equal one);
- 2 step  $\ell$ :
  - (a) MCA on  $\mathbf{X}^{\ell-1}$ : PCA on  $(I\mathbf{X}^{\ell-1}(\mathbf{D}_{\Sigma}^{\ell-1})^{-1}, \frac{1}{IJ}\mathbf{D}_{\Sigma}^{\ell-1}, \frac{1}{I}\mathbb{I}_I) \rightarrow (\hat{\mathbf{F}}^{\ell}, \hat{\mathbf{U}}^{\ell})$ ;  $S$  dimensions are kept
  - (b) impute the indicator matrix using the model matrix:

$$\hat{x}_{ik}^{\ell} = \frac{1}{I} \left( 1 + \sum_{s=1}^S \hat{f}_{is}^{\ell} \hat{u}_{ks}^{\ell} \right) \mathbf{D}_{\Sigma}^{\ell-1}$$

The new imputed dataset is  $\mathbf{X}^{\ell} = \mathbf{W} * \mathbf{X} + (1 - \mathbf{W}) * \hat{\mathbf{X}}^{\ell}$

- (c)  $\mathbf{D}_{\Sigma}^{\ell}$  is updated with the new column margins  $I_k^{\ell}$  of  $\mathbf{X}^{\ell}$ ;
- 3 steps (a), (b) and (c) are repeated until convergence

# Iterative MCA

- 1 initialization  $\ell = 0$ :  $\mathbf{X}^0$  missing values are imputed with the proportion of the category  $\frac{I_k}{I}$  (the sum must equal one);
- 2 step  $\ell$ :
  - (a) MCA on  $\mathbf{X}^{\ell-1}$ : PCA on  $(I\mathbf{X}^{\ell-1}(\mathbf{D}_{\Sigma}^{\ell-1})^{-1}, \frac{1}{IJ}\mathbf{D}_{\Sigma}^{\ell-1}, \frac{1}{I}\mathbb{I}_I) \rightarrow (\hat{\mathbf{F}}^{\ell}, \hat{\mathbf{U}}^{\ell})$ ;  $S$  dimensions are kept
  - (b) impute the indicator matrix using the model matrix:

$$\hat{x}_{ik}^{\ell} = \frac{1}{I} \left( 1 + \sum_{s=1}^S \hat{f}_{is}^{\ell} \hat{u}_{ks}^{\ell} \right) \mathbf{D}_{\Sigma}^{\ell-1}$$

The new imputed dataset is  $\mathbf{X}^{\ell} = \mathbf{W} * \mathbf{X} + (1 - \mathbf{W}) * \hat{\mathbf{X}}^{\ell}$  ;

- (c)  $\mathbf{D}_{\Sigma}^{\ell}$  is updated with the new column margins  $I_k^{\ell}$  of  $\mathbf{X}^{\ell}$
- 3 steps (a), (b) and (c) are repeated until convergence

## Iterative MCA

- Step 0: missing fuzzy average = reconstruction of order 0

$$\hat{x}_{ik}^0 = \frac{1}{I} \left( 1 + \sum_{s=1}^S \hat{f}_{is}^l \hat{u}_{ks}^l \right) \mathbf{D}_{\Sigma}^0 = \mathbf{D}_{\Sigma}^0 / I$$

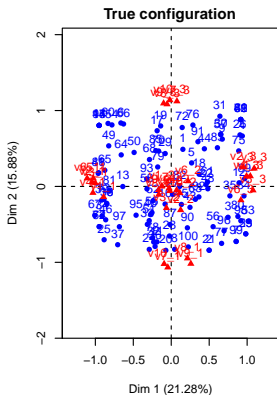
- The algorithm returns a completed indicator matrix

	V1	V2	V3
ind 1	a	NA	g
ind 2	NA	f	g
ind 3	a	e	h
ind 4	a	e	h
ind 5	b	f	h
ind 6	c	f	h
ind 7	c	f	NA

	V1_a	V1_b	V1_c	V2_e	V2_f	V3_g	V3_h
ind 1	1	0	0	<b>0.71</b>	<b>0.29</b>	1	0
ind 2	<b>0.13</b>	<b>0.29</b>	<b>0.59</b>	0	1	1	0
ind 3	1	0	0	1	0	0	1
ind 4	1	0	0	1	0	0	1
ind 5	0	1	0	0	1	0	1
ind 6	0	0	1	0	1	0	1
ind 7	0	0	1	0	1	<b>0.37</b>	<b>0.63</b>

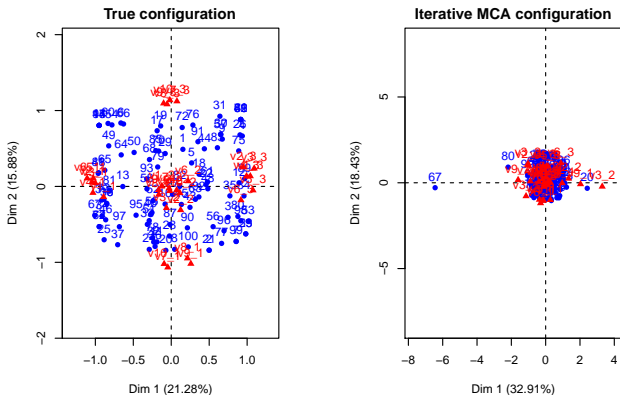
⇒ imputed values can be seen as degree of membership

# Overfitting

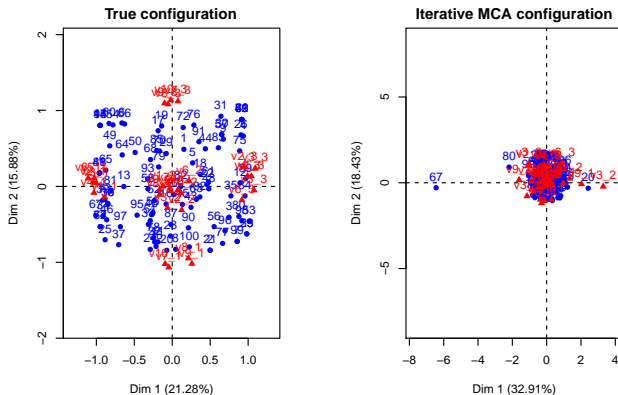




## Overfitting



# Overfitting



$$\text{mean}_{i,k \in \text{obs}}(x_{ik} - \hat{x}_{ik})^2 = 0.026 \text{ whereas } \text{mean}_{i,k \notin \text{obs}}(x_{ik} - \hat{x}_{ik})^2 = 0.12$$

⇒ observed values are well-fitted but the prediction of the missing values and the estimation of the axes and components are very poor  
 ⇒ convergence problems of EM → overfitting problems

# Overfitting

⇒ Fitting is good, prediction is bad

- Many parameters are estimated with respect to the number of observed values: the number of dimensions  $S$  and the number of missing values are important
- The relationship between variables are not strong

① Reduce  $S$

② Early stopping

③ Shrinkage methods

## Regularized Iterative MCA

⇒ Initialization, estimation step - imputation step.

The imputation step:

$$\hat{x}_{ik}^{\ell} = \frac{1}{I} \left( 1 + \sum_{s=1}^S \hat{f}_{is}^{\ell} \hat{u}_{ks}^{\ell} \right) \mathbf{D}_{\Sigma}^{\ell-1} = \frac{1}{I} \left( 1 + \sum_{s=1}^S \frac{\hat{f}_{is}^{\ell}}{\|\hat{\mathbf{f}}_s^{\ell}\|} (\sqrt{\lambda_s}) \hat{u}_{ks}^{\ell} \right) \mathbf{D}_{\Sigma}^{\ell-1},$$

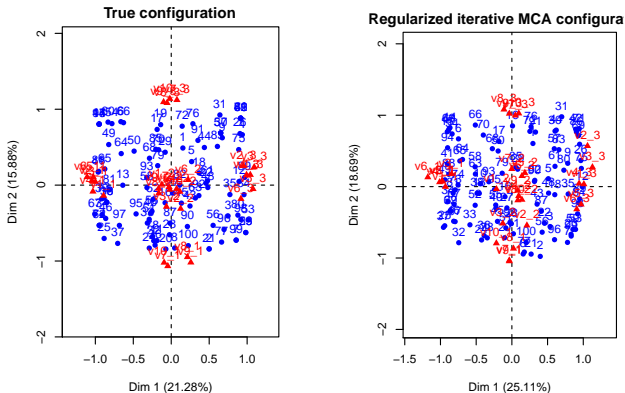
is replaced by a "shrunk" imputation step:

$$\hat{x}_{ik}^{\ell} = \frac{1}{I} \left( 1 + \sum_{s=1}^S \frac{\hat{f}_{is}^{\ell}}{\|\hat{\mathbf{f}}_s^{\ell}\|} \left( \sqrt{\lambda_s} - \frac{\hat{\sigma}^2}{\sqrt{\lambda_s}} \right) \hat{u}_{ks}^{\ell} \right) \mathbf{D}_{\Sigma}^{\ell-1}$$

with  $\hat{\sigma}^2 = \frac{1}{K-J-S} \sum_{s=S+1}^{K-J} \lambda_s$ .

⇒ Remove the noise to avoid instability on the predictions

## Regularized iterative MCA



$$\text{mean}_{i,k \notin \text{obs}} (x_{ik} - \hat{x}_{ik})^2 = 0.056$$

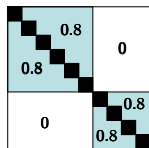
# Simulations

Many scenarios are considered:

- percentage of missing values: small, medium
- missing values mechanism: MCAR, MAR
- pattern of missing values: random or not random
- relationship between variables: low or strong
- 1000 simulations

The simulated data:

- 100 individuals
- 10 variables from a normal distribution
- each variable is cut in 3 equal-count categories  
⇒ By construction, 4 underlying dimensions



## Missing single

	V1	V2	V3
ind 1	a	NA	g
ind 2	NA	f	g
ind 3	a	e	h
ind 4	a	e	h
ind 5	b	f	h
ind 6	c	f	h
ind 7	c	f	NA

	V1_a	V1_b	V1_c	V1_NA	V2_e	V2_f	V2_NA	V3_g	V3_h	V3_NA
ind 1	1	0	0	0	0	0	1	1	0	0
ind 2	0	0	0	1	0	1	0	1	0	0
ind 3	1	0	0	0	1	0	0	0	1	0
ind 4	1	0	0	0	1	0	0	0	1	0
ind 5	0	1	0	0	0	1	0	0	1	0
ind 6	0	0	1	0	0	1	0	0	1	0
ind 7	0	0	1	0	0	1	0	0	0	1

- A new category is added for missing values  
 ⇒ well-adapted for “not really missing” or MNAR

## Missing passive modified margin

- Missing passive (Meulman, 1982)

	V1	V2	V3
ind 1	a	NA	g
ind 2	NA	f	g
ind 3	a	e	h
ind 4	a	e	h
ind 5	b	f	h
ind 6	c	f	h
ind 7	c	f	NA

	V1_a	V1_b	V1_c	V2_e	V2_f	V3_g	V3_h
ind 1	1	0	0	0	0	1	0
ind 2	0	0	0	0	1	1	0
ind 3	1	0	0	1	0	0	1
ind 4	1	0	0	1	0	0	1
ind 5	0	1	0	0	1	0	1
ind 6	0	0	1	0	1	0	1
ind 7	0	0	1	0	1	0	0

Row margins are not equal  $\Rightarrow$  many properties of MCA are lost

- Missing passive modified margin (Escofier, 1987)
  - $\Rightarrow$  row margins are fixed to  $J$
  - $\Rightarrow$  Equivalence with subset MCA (Greenacre & Pardo, 2006)



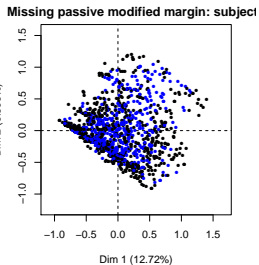
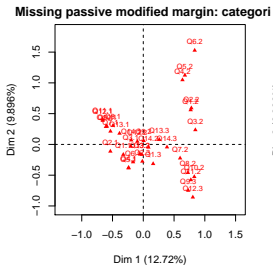
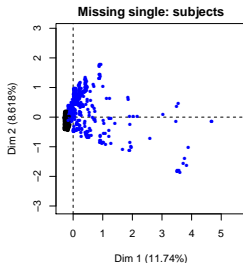
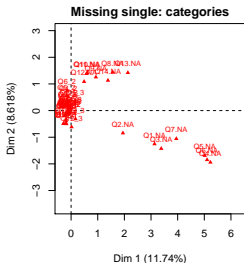
## Simulations

The criterion used is the RV coefficient between the configuration without missing values and the one obtained from the algorithm

Missing	Link	Missing Passive Modified Margin R - NR	Missing Fuzzy Average R - NR	Missing single R - NR	RiMCA R - NR
10% MCAR	low	<b>0.94</b> - 0.91	<b>0.94</b> - 0.92	0.87 - 0.47	<b>0.94</b> - <b>0.93</b>
10% MCAR	strong	0.97 - 0.94	0.97 - 0.95	0.96 - 0.68	<b>0.98</b> - <b>0.97</b>
30% MCAR	low	<b>0.77</b> - 0.44	0.77 - 0.77	0.67 - 0.32	0.76 - <b>0.78</b>
30% MCAR	strong	0.88 - 0.71	0.88 - <b>0.91</b>	0.86 - 0.46	<b>0.91</b> - 0.90
8% MAR	low	0.94 - 0.91	0.94 - 0.91	0.72 - 0.28	<b>0.95</b> - <b>0.92</b>
8% MAR	strong	0.96 - 0.91	0.96 - 0.90	0.96 - 0.54	<b>0.98</b> - <b>0.96</b>
16% MAR	low	0.86 - 0.80	0.83 - 0.79	0.50 - 0.29	<b>0.88</b> - <b>0.83</b>
16% MAR	strong	0.89 - 0.80	0.84 - 0.78	0.88 - 0.55	<b>0.95</b> - <b>0.90</b>

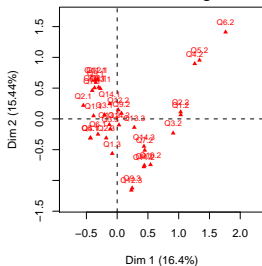
## A real example

- 1232 respondents, 14 questions, 35 categories, 9% of missing values concerning 42% of respondents

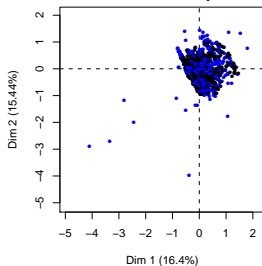


# A real example

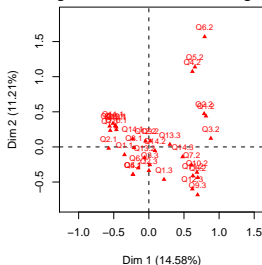
iterative MCA: categories



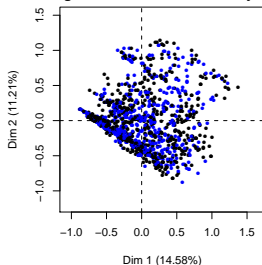
iterative MCA: subjects



Regularized iterative MCA: categories



Regularized iterative MCA: subjects



# Conclusion

## Regularized iterative MCA

- gives “good” results
- is efficient when strong relationships between variables (you learn from the other variables) ...
- ... but needs tuning parameters
- can be used as an imputation method?
- can be used to perform a clustering on categorical variables with missing values
- is available in the `missMDA` package that imputes the indicator matrix and the `FactoMineR` package that performs the MCA from an indicator matrix