

Gestion des données manquantes en/par ACM et analyse de données mixtes

François Husson

UP de mathématiques appliquées - l'institut Agro



Journées d'études en statistique – SFdS 2021

- 1 Introduction
- 2 ACM particulières
- 3 Imputation par ACM itérative
- 4 Imputation simple pour données mixtes
- 5 Données multi-niveaux

Les méthodes d'analyse factorielle

- Analyse exploratoire de tableaux de données
- Dépend de la structure et de la nature des variables :
 - ACP : variables quantitatives
 - ACM : variables qualitatives
 - AFDM : variables quantitatives et qualitatives
 - AFM : structure avec des groupes de variables
 - ...

Toutes les méthodes d'analyse factorielle peuvent être vues comme une ACP sur une matrice particulière avec des poids spécifiques pour les lignes et les colonnes

« Doing a data analysis, in good mathematics, is simply searching eigenvectors, all the science of it (the art) is just to find the right matrix to diagonalize » (Benzécri)

Rappels d'ACM

- Analyse exploratoire d'un tableau de variables qualitatives
- Analyse de questionnaires

$T =$

1	0	0	1	0	0	1	...	0	1	p		
1	0	0	1	0	1	0	...	NA	NA		p	
NA	NA	NA	0	1	0	0	...	0	1			p
1	0	0	1	0	0	1	...	0	1			
t_{ik}												
0	0	1	NA	NA	0	1	...	0	1	p		
1	0	0	1	0	0	1	...	0	1		p	
n_1	n_2	n_3						n_q	np		

$$P_{\Sigma} = \begin{array}{ccccccc} n_1 & & & & & & \\ & n_2 & & & & & \\ & & n_3 & & & & \\ & & & \dots & & & \\ & & & & \dots & & \\ & & & & & \dots & \\ & & & & & & \dots \\ & & & & & & & n_q \end{array}$$

L'ACM comme une ACP pondérée

ACM vue comme l'ACP du triplet

$$\left(n\mathbf{T}\mathbf{P}_{\Sigma}^{-1}, \frac{1}{np}\mathbf{P}_{\Sigma}, \frac{1}{n}\mathbf{I}_n \right)$$

$$\left(n\mathbf{T}\mathbf{P}_{\Sigma}^{-1}, \frac{1}{np}\mathbf{P}_{\Sigma}, \frac{1}{n}I_n \right)$$

ACP d'un triplet $(\mathbf{A}, \mathbf{M}, \mathbf{P})$

L'ACP d'un triplet $(\mathbf{A}, \mathbf{M}, \mathbf{P})$ est la SVD suivante :

$$A = UDV'$$

avec \mathbf{U} les vecteurs propres de $\mathbf{A}\mathbf{M}\mathbf{A}'\mathbf{P}$ et tels que $\mathbf{U}'\mathbf{P}\mathbf{U} = Id$

et \mathbf{V} les vecteurs propres de $\mathbf{A}'\mathbf{P}\mathbf{A}\mathbf{M}$ et tels que $\mathbf{V}'\mathbf{M}\mathbf{V} = Id$

$$\left(n\mathbf{T}\mathbf{P}_{\Sigma}^{-1}, \frac{1}{np}\mathbf{P}_{\Sigma}, \frac{1}{n}I_n \right)$$

ACP d'un triplet $(\mathbf{A}, \mathbf{M}, \mathbf{P})$

L'ACP d'un triplet $(\mathbf{A}, \mathbf{M}, \mathbf{P})$ est la SVD suivante :

$$\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}'$$

avec \mathbf{U} les vecteurs propres de $\mathbf{A}\mathbf{M}\mathbf{A}'\mathbf{P}$ et tels que $\mathbf{U}'\mathbf{P}\mathbf{U} = Id$

et \mathbf{V} les vecteurs propres de $\mathbf{A}'\mathbf{P}\mathbf{A}\mathbf{M}$ et tels que $\mathbf{V}'\mathbf{M}\mathbf{V} = Id$

\mathbf{U} ; \mathbf{D} ; \mathbf{V} minimisent le critère d'erreur de reconstitution :

$$\mathcal{C} = \|\mathbf{A} - \mathbf{U}\mathbf{D}\mathbf{V}'\|_{\mathbf{M},\mathbf{P}}^2$$

- ① Introduction
- ② ACM particulières
- ③ Imputation par ACM itérative
- ④ Imputation simple pour données mixtes
- ⑤ Données multi-niveaux

Traitement d'un questionnaire avec *missing single*

Les données

1232 répondants, 14 questions, 35 modalités, 9% de NA pour 42% des répondants

Traitement d'un questionnaire avec *missing single*

Les données

1232 répondants, 14 questions, 35 modalités, 9% de NA pour 42% des répondants

Création de nouvelles modalités

Création d'une modalité NA pour chaque variable ayant au moins une valeur manquante

	V1	V2	V3		V1_a	V1_b	V1_c	V1_NA	V2_e	V2_f	V2_NA	V3_g	V3_h
ind 1	a	NA	g		1	0	0	0	0	0	1	1	0
ind 2	NA	f	g		0	0	0	1	0	1	0	1	0
ind 3	a	e	h		1	0	0	0	1	0	0	0	1
ind 4	a	e	h		1	0	0	0	1	0	0	0	1
ind 5	b	f	h		0	1	0	0	0	1	0	0	1
ind 6	c	f	h		0	0	1	0	0	1	0	0	1
ind 7	c	f	h		0	0	1	0	0	1	0	0	1

Traitement d'un questionnaire avec *missing single*

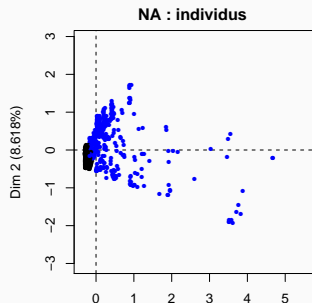
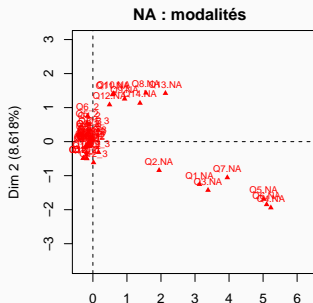
Les données

1232 répondants, 14 questions, 35 modalités, 9% de NA pour 42% des répondants

Création de nouvelles modalités

Création d'une modalité NA pour chaque variable ayant au moins une valeur manquante

	V1	V2	V3		V1_a	V1_b	V1_c	V1_NA	V2_e	V2_f	V2_NA	V3_g	V3_h
ind 1	a	NA	g	ind 1	1	0	0	0	0	0	1	1	0
ind 2	NA	f	g	ind 2	0	0	0	1	0	1	0	1	0
ind 3	a	e	h	ind 3	1	0	0	0	1	0	0	0	1
ind 4	a	e	h	ind 4	1	0	0	0	1	0	0	0	1
ind 5	b	f	h	ind 5	0	1	0	0	0	1	0	0	1
ind 6	c	f	h	ind 6	0	0	1	0	0	1	0	0	1
ind 7	c	f	h	ind 7	0	0	1	0	0	1	0	0	1



Traitement d'un questionnaire avec la méthode *missing passive*

Les données

1232 répondants, 14 questions, 35 modalités, 9% de NA pour 42% des répondants

Methode missing passive

Si un NA, considérer que l'individu n'a choisi aucune modalité

	V1	V2	V3		V1 a	V1 b	V1 c	V2 e	V2 f	V3 g	V3 h
ind 1	a	NA	g		1	0	0	0	0	1	0
ind 2	NA	f	g		0	0	0	0	1	1	0
ind 3	a	e	h		1	0	0	1	0	0	1
ind 4	a	e	h		1	0	0	1	0	0	1
ind 5	b	f	h		0	1	0	0	1	0	1
ind 6	c	f	h		0	0	1	0	1	0	1
ind 7	c	f	h		0	0	1	0	1	0	1

Traitement d'un questionnaire avec la méthode *missing passive*

Les données

1232 répondants, 14 questions, 35 modalités, 9% de NA pour 42% des répondants

Methode missing passive

Si un NA, considérer que l'individu n'a choisi aucune modalité

	V1	V2	V3		V1 a	V1 b	V1 c	V2 e	V2 f	V3 g	V3 h
ind 1	a	NA	g	ind 1	1	0	0	0	0	1	0
ind 2	NA	f	g	ind 2	0	0	0	0	1	1	0
ind 3	a	e	h	ind 3	1	0	0	1	0	0	1
ind 4	a	e	h	ind 4	1	0	0	1	0	0	1
ind 5	b	f	h	ind 5	0	1	0	0	1	0	1
ind 6	c	f	h	ind 6	0	0	1	0	1	0	1
ind 7	c	f	h	ind 7	0	0	1	0	1	0	1

Problème

Marges lignes inégales \implies perte de nombreuses propriétés de l'ACM
 \implies imposer les marges par méthode missing passive modified margin

Methode missing passive modified margin

Prendre la méthode missing passive et imposer l'égalité des marges

	V1	V2	V3		V1 a	V1 b	V1 c	V2 e	V2 f	V3 g	V3 h
ind 1	a	NA	g	ind 1	1	0	0	0	0	1	0
ind 2	NA	f	g	ind 2	0	0	0	0	1	1	0
ind 3	a	e	h	ind 3	1	0	0	1	0	0	1
ind 4	a	e	h	ind 4	1	0	0	1	0	0	1
ind 5	b	f	h	ind 5	0	1	0	0	1	0	1
ind 6	c	f	h	ind 6	0	0	1	0	1	0	1
ind 7	c	f	h	ind 7	0	0	1	0	1	0	1

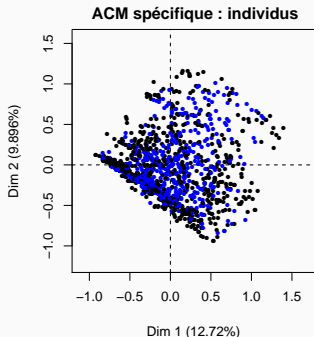
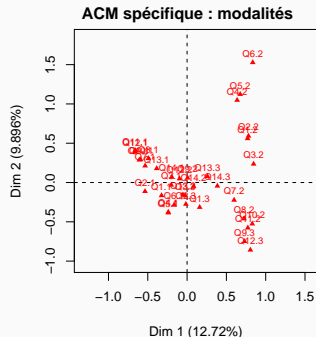
ACM spécifique – missing passive modified margin

Methode missing passive modified margin

Prendre la méthode missing passive et imposer l'égalité des marges

	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	h

	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	1

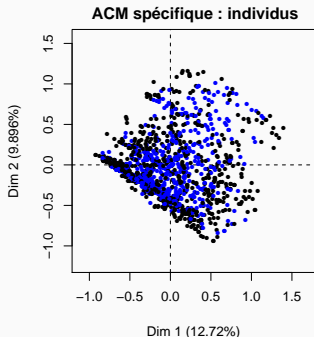
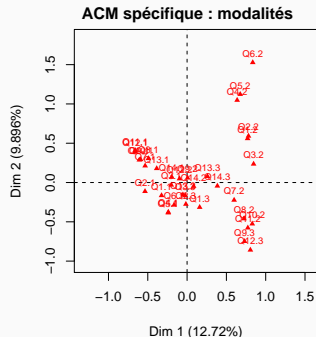


Methode missing passive modified margin

Prendre la méthode missing passive et imposer l'égalité des marges

	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	h

	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	1



Inconvénient

revient à considérer que les individus auraient pris des modalités différentes de celles proposées

Equivalence avec une ACM spécifique

ACM spécifique qui construit les axes en mettant en supplémentaire les modalités NA

- ① Introduction
- ② ACM particulières
- ③ Imputation par ACM itérative**
- ④ Imputation simple pour données mixtes
- ⑤ Données multi-niveaux

ACM itérative régularisée (Josse *et al.*, 2012)

- 1 Initialisation : imputation de la matrice indicatrice (proportion)
- 2 Itération jusqu'à convergence
 - (a) Estimation de $\mathbf{U}^\ell, \mathbf{D}^\ell, \mathbf{V}^\ell$: ACM sur le tableau complété, i.e. l'ACP du triplet

$$\left(n\mathbf{T}^{\ell-1}(\mathbf{P}_\Sigma^{\ell-1})^{-1}, \frac{1}{np}\mathbf{P}_\Sigma^{\ell-1}, \frac{1}{n}\mathbb{I}_n \right)$$

- (b) Utiliser la formule de reconstitution (prendre les valeurs singulières régularisées) :

$$(\hat{a}_{ik}^\ell - 1)\sqrt{\frac{n_k^{\ell-1}}{np}} = \left(\sum_{s=2}^S \hat{u}_{is} \left(\hat{d}_s - \frac{\hat{\sigma}^2}{\hat{d}_s} \right) \hat{v}_{ks} \right)$$


Calculer les valeurs reconstituées en utilisant les marges de l'étape $\ell - 1$: $\hat{\mathbf{T}}^\ell = \frac{1}{n}\hat{\mathbf{A}}^\ell\mathbf{P}_\Sigma^{\ell-1}$
et le nouveau tableau disjonctif complété est $\mathbf{T}^\ell = \mathbf{R} * \mathbf{T} + (1 - \mathbf{R}) * \hat{\mathbf{T}}^\ell$

- (c) **Mise à jour des marges** : les marges colonnes n_k^ℓ du nouveau tableau complété \mathbf{T}^ℓ sont calculées et enregistrées dans \mathbf{P}_Σ^ℓ ;
- 3 les étapes (2.a), (2.b) et (2.c) sont répétées jusqu'à convergence.

ACM itérative régularisée (Josse *et al.*, 2012)

	V1	V2	V3	...	V14
ind 1	a	NA	g	...	u
ind 2	NA	f	g		u
ind 3	a	e	h		v
ind 4	a	e	h		v
ind 5	b	f	h		u
ind 6	c	f	h		u
ind 7	c	f	NA		v
...
ind 1232	c	f	h		v

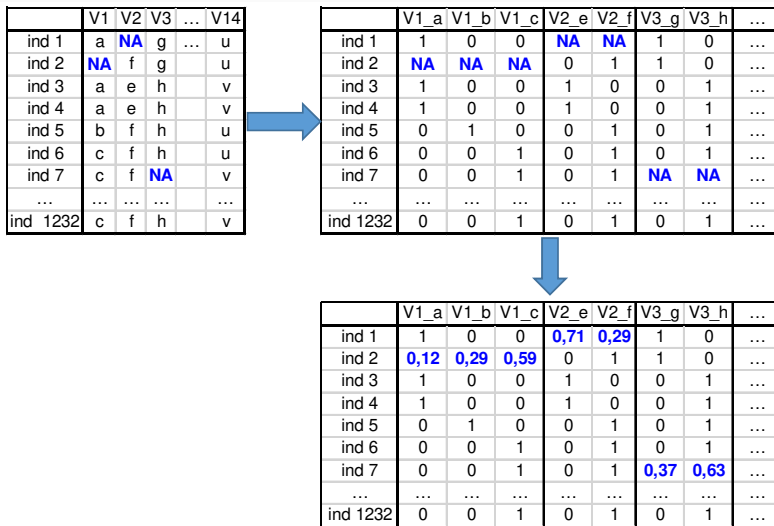
ACM itérative régularisée (Josse *et al.*, 2012)



	V1	V2	V3	...	V14
ind 1	a	NA	g	...	u
ind 2	NA	f	g		u
ind 3	a	e	h		v
ind 4	a	e	h		v
ind 5	b	f	h		u
ind 6	c	f	h		u
ind 7	c	f	NA		v
...
ind 1232	c	f	h		v

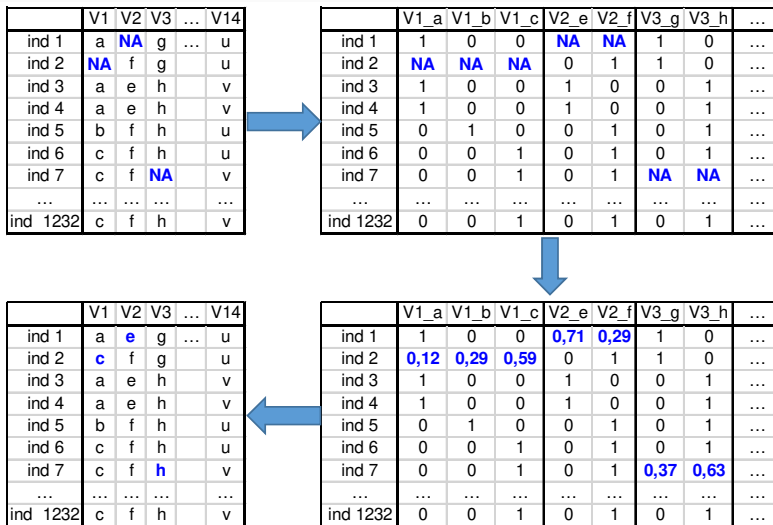
	V1_a	V1_b	V1_c	V2_e	V2_f	V3_g	V3_h	...
ind 1	1	0	0	NA	NA	1	0	...
ind 2	NA	NA	NA	0	1	1	0	...
ind 3	1	0	0	1	0	0	1	...
ind 4	1	0	0	1	0	0	1	...
ind 5	0	1	0	0	1	0	1	...
ind 6	0	0	1	0	1	0	1	...
ind 7	0	0	1	0	1	NA	NA	...
...
ind 1232	0	0	1	0	1	0	1	...

ACM itérative régularisée (Josse *et al.*, 2012)



Les valeurs imputées peuvent être vues comme des degrés d'appartenance

ACM itérative régularisée (Josse *et al.*, 2012)



Les valeurs imputées peuvent être vues comme des degrés d'appartenance

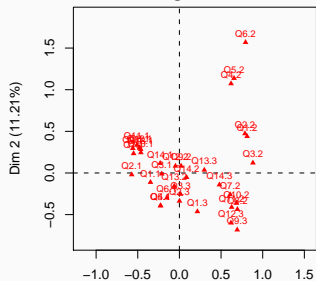
Imputation du tableau disjonctif

```
> library(missMDA)
> data(vnf)
> ncp <- estim_ncpMCA(vnf)
> res.impute <- imputeMCA(vnf, ncp=4)
```

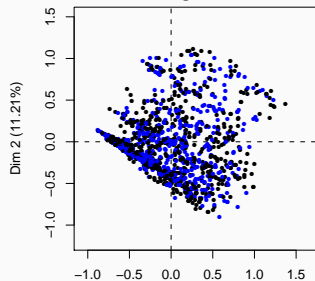
ACM sur le tableau complété (utilisation de l'argument tab.disj)

```
> res.mca <- MCA(vnf, tab.disj = res.impute$tab.disj)
```

ACM itérative régularisée : modalités



ACM itérative régularisée : individus



- ① Introduction
- ② ACM particulières
- ③ Imputation par ACM itérative
- ④ Imputation simple pour données mixtes**
- ⑤ Données multi-niveaux

Analyse Factorielle de Données Mixtes (cas complet)

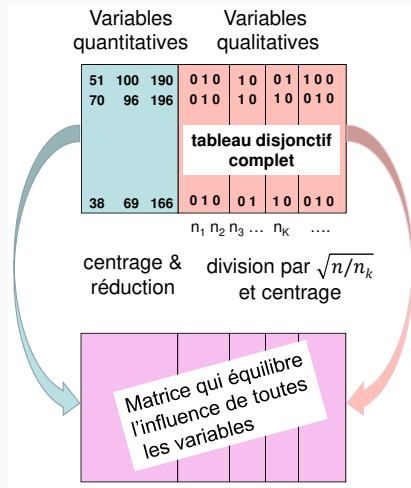
AFDM (Escofier, 1979), PCAMIX (Kiers, 1991)

- ACP sur une matrice pondérée
- La distance entre individus s'écrit :

$$d^2(i, l) = \sum_{j=1}^{p_1} (t_{ik} - t_{lk})^2 + \sum_{j=1}^{p_2} \sum_{k=1}^{K_j} \frac{1}{n_{k_j}} (t_{ij} - t_{lj})^2$$

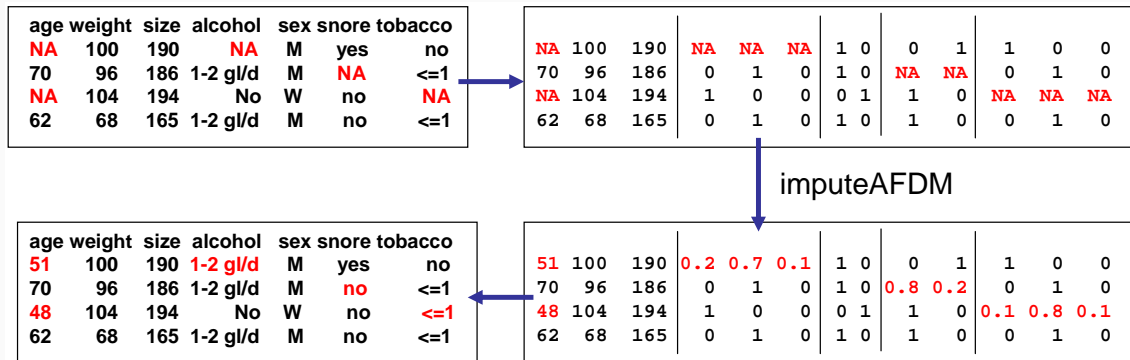
- Les composantes principales \mathbf{F}_s maximisent :

$$\sum_{j=1}^{p_1} r^2(\mathbf{F}_s, v_j) + \sum_{j=1}^{p_2} \eta^2(\mathbf{F}_s, v_j)$$



Algorithme d'AFDM itératif

- 1 Initialisation : imputation par la moyenne (quanti) et la proportion (quali)
- 2 Itérer jusqu'à convergence
 - (a) estimation : AFDM sur le jeu complété $\Rightarrow \mathbf{U}, \mathbf{D}, \mathbf{V}$
 - (b) imputation des valeurs manquantes avec le modèle de reconstitution
 - (c) moyennes, écarts-types et marges sont mis à jour



- Dispositif de simulations
 - 2 variables indépendantes provenant d'une distribution normale
 - 1 variable répétée 4 fois, l'autre 8 \Rightarrow 2 dimensions
 - Bruit ajouté
 - La moitié des variables sur chaque dimension sont découpées en 3 classes
 - 10%, 20% or 30% de données manquantes au hasard

\Rightarrow Données sont construites pour être en 4 dimensions

- Critère
 - pour données quantitatives :

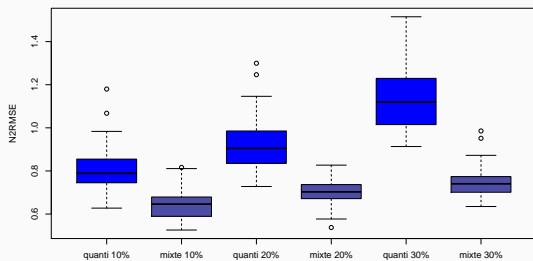
$$N2RMSE = \sqrt{\sum_{i \in \text{manquant}} \frac{\text{moyenne} \left((X_i^{\text{vrai}} - X_i^{\text{imp}})^2 \right)}{\text{var} (X_i^{\text{true}})}}$$

- pour données qualitatives : proportion de modalités mal prédites

Imputation avec var. quanti uniquement

Imputation avec variables quanti et quali

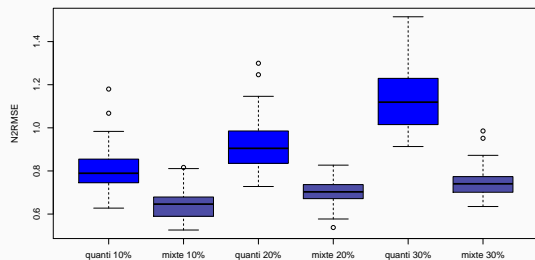
Error on continous data



Imputation avec var. quanti uniquement

Imputation avec variables quanti et quali

Error on continous data

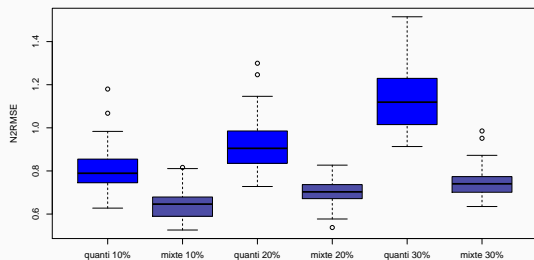


Variables quali améliorent
l'imputation sur variables quanti ...

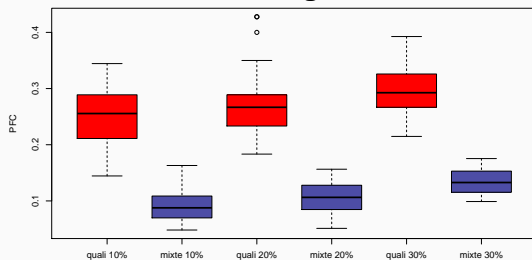
Simulations

Imputation avec var. quanti uniquement Imputation avec var. quali uniquement
Imputation avec variables quanti et quali

Error on continous data



Error on categorical data

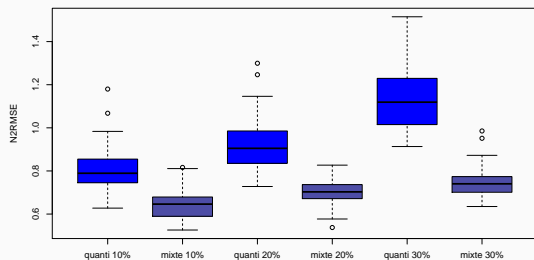


Variables quali améliorent
l'imputation sur variables quanti ...

Simulations

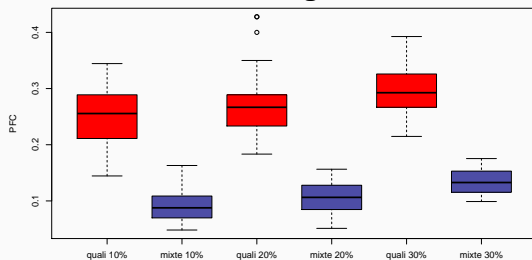
Imputation avec var. quanti uniquement Imputation avec var. quali uniquement
Imputation avec variables quanti et quali

Error on continous data



Variables quali améliorent
l'imputation sur variables quanti ...

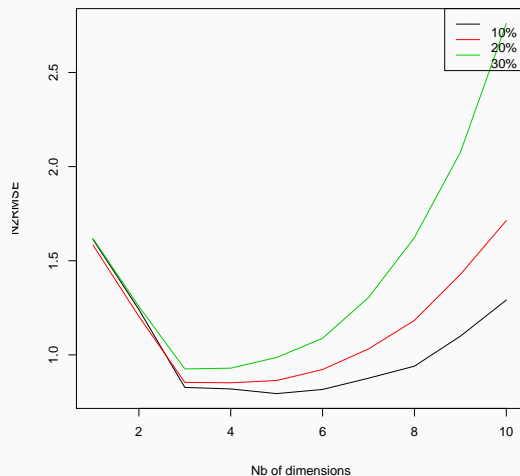
Error on categorical data



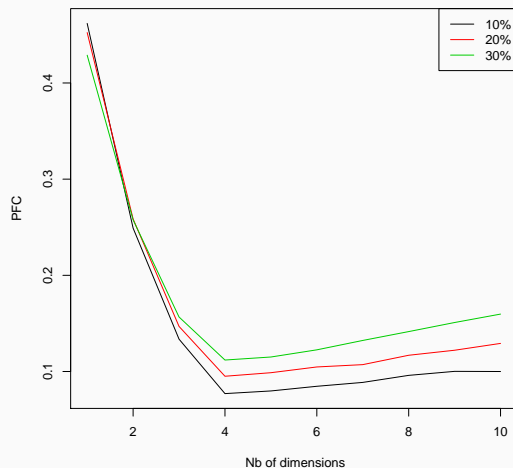
... et variables quanti améliorent l'imputation des variables quali

Simulations

Error on continuous variables



Error on the qualitative variables

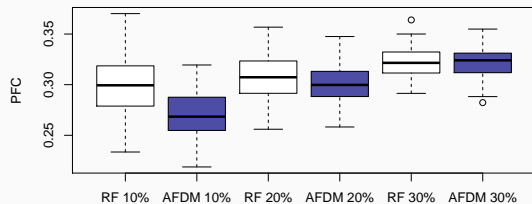
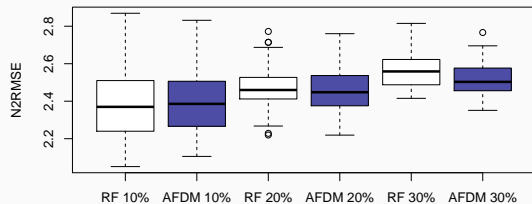


⇒ L'erreur sur le choix du nombre de dimensions a un impact faible sur l'erreur d'imputation
... si l'estimation n'est pas trop mauvaise

Comparaison avec forêts aléatoires

Imputations obtenues par forêts aléatoires & [ACP itérative](#)

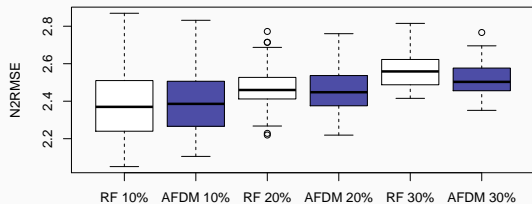
GBSG2



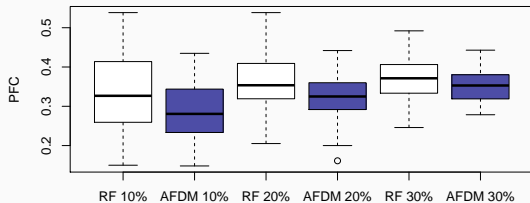
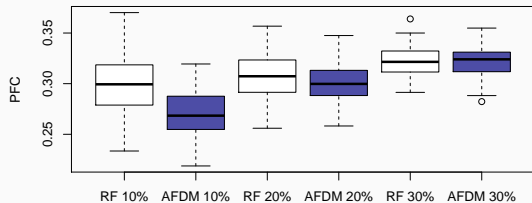
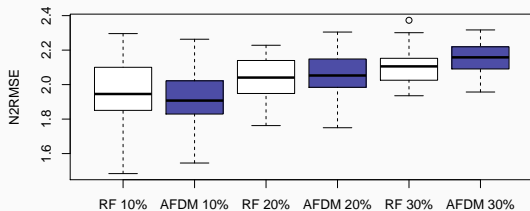
Comparaison avec forêts aléatoires

Imputations obtenues par forêts aléatoires & ACP itérative

GBSG2



Ozone



Imputation de données mixtes en pratique

```
> library(missMDA)
> nb <- estim_ncpFAMD(mydata)  ## tps de calcul long
> res.imp <- imputeFAMD(mydata, ncp = nb$ncp)
> res.famd <- FAMD(mydata, ,tab.disj = res.imp$tab.disj)

> library(missForest)
> missForest(mydata)

> library(mice)
> mice(mydata)
> mice(mydata, defaultMethod = "rf") ## mice avec forêts aléatoires
```

Analyse Factorielle Multiple

Même principe avec mise à jour des premières valeurs propres de chaque groupe en plus

Cas de groupes quantitatifs uniquement : le tableau est complété et l'AFM est lancée sur le tableau complété :

```
> data(orange)
> res.comp <- imputeMFA(orange, group=c(5,3), type=rep("s",2), ncp=2)
> res.mfa <- MFA(res.comp$completeObs, group=c(5,3), type=rep("s",2))
```

Cas où au moins un groupe qualitatif : le "tableau disjonctif" complété est fourni à l'AFM avec l'argument `tab.comp` :

```
> data(vnf)
> res.comp <- imputeMFA(vnf,group=c(6,5,3),type=c("n","n","n"),ncp=2)
> res.mfa <- MFA(vnf,group=c(6,5,3),type=c("n","n","n"), tab.comp=res.comp)
```

Bilan sur l'imputation simple

⇒ Données manquantes en analyse factorielle

- tableau simple : ACP, ACM, analyse fact. de données mixtes
- tableaux multiples (AFM)

⇒ Pré-traitement avant classification (avec données manquantes)

⇒ package R `missMDA` – `Factoshiny`

⇒ Imputation des données quantitatives, qualitatives, mixtes

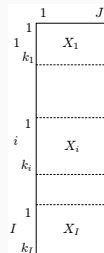
- basée sur la reconstitution de l'ACP (axes et composantes)
- prise en compte des liaisons entre var. quantitatives et qualitatives
- bonne alternative aux méthodes d'imputation (forêts aléatoires, etc.) si liaisons linéaires, pour les variables qualitatives (notamment les modalités rares)

- ① Introduction
- ② ACM particulières
- ③ Imputation par ACM itérative
- ④ Imputation simple pour données mixtes
- ⑤ **Données multi-niveaux**

Analyse en composantes multi-niveaux

Ex : patients hiérarchisés dans hôpitaux $X \in \mathbb{R}^{K \times J}$

- similarités entre hôpitaux ? niveau 1
- similarités entre patients dans un même hôpital ? niveau 2
- relations entre variables à chaque niveau



$$x_{ijk_i} = x_{.j.} + (x_{ij.} - x_{.j.}) + (x_{ijk_i} - x_{ij.})$$

Between + Within

Analysis de variance : décomposer la somme des carrés pour chaque variable j

$$\sum_{i=1}^I \sum_{k=1}^{k_i} (x_{ijk_i})^2 = \sum_{i=1}^I k_i (x_{.j.})^2 + \sum_{i=1}^I k_i (x_{ij.} - x_{.j.})^2 + \sum_{i=1}^I \sum_{k=1}^{k_i} (x_{ijk_i} - x_{ij.})^2$$

⇒ Modèle pour la partie between et within $i = 1, \dots, I$ groupes, J var

$$X_{i(k_i \times J)} = \mathbf{1}_{k_i} m' + \mathbf{1}_{k_i} U_i^b D^b V^{b'} + F_i^w D^w V^{w'} + E_i$$

- F_i^b ($Q_b \times 1$) between component scores of group i
- V^b ($J \times Q_b$) between loading matrix
- F_i^w ($k_i \times Q_w$) within component scores of group i
- V_w ($J \times Q_w$) within loading matrix. **Constant across groups**

Solution obtenue par moindres carrés (Timmerman, 2006)

Possibilité de faire des calculs distribués

Remarque

“The idea of imputation is both seductive and dangerous. It is seductive because it can lull the user into the pleasurable state of believing that the data are complete after all, and it is dangerous because it lumps together situations where the problem is sufficiently minor that it can be legitimately handled in this way and situations where standard estimators applied to the real and imputed data have substantial biases.” (Dempster & Rubin, 1983)

Remarque

“The idea of imputation is both seductive and dangerous. It is seductive because it can lull the user into the pleasurable state of believing that the data are complete after all, and it is dangerous because it lumps together situations where the problem is sufficiently minor that it can be legitimately handled in this way and situations where standard estimators applied to the real and imputed data have substantial biases.” (Dempster & Rubin, 1983)

Imputation simple *versus* imputation multiple

On ne peut accorder la même confiance à une valeur imputée et une valeur observée

L'imputation simple retourne 1 seule valeur pour chaque valeur manquante et 1 seule valeur ne permet pas de connaître l'incertitude sur la prédiction de cette valeur

⇒ Imputation multiple