

# Gestion des données manquantes en/par analyse factorielle

F. Husson

UP de mathématiques appliquées - l'institut Agro

husson@agrocampus-ouest.fr

*"Je ne suis pas en train de dire que pour être heureux il faut faire tout le temps des mathématiques. Mais toute personne qui s'y exerce sérieusement fait forcément l'expérience du bonheur."*

Alain Badiou – philosophe

# Plan

- 1 Introduction
- 2 Imputation simple pour variables quantitatives
- 3 Imputation simple pour variables qualitatives
- 4 Imputation simple pour données mixtes
- 5 Imputation multiple
- 6 Conclusion

## Les données manquantes



Gertrude Mary Cox

*“The best thing to do about missing values is not to have any”*

Les données manquantes sont très présentes en pratique :  
non-réponse à un questionnaire, données perdues ou détruites, appareils qui tombent en panne, plantes détruisent (maladie, ravageurs, etc.) ...

## Les données manquantes

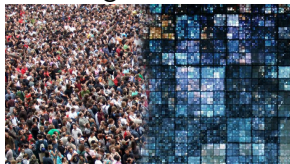


Gertrude Mary Cox

*“The best thing to do about missing values is not to have any”*

Les données manquantes sont très présentes en pratique :  
non-réponse à un questionnaire, données perdues ou détruites, appareils qui tombent en panne, plantes détruisent (maladie, ravageurs, etc.) ...

Et en big data ?



“One of the ironies of Big Data is that missing data play an ever more significant role” (R. Sameworth, 2019)

Une matrice  $n \times p$ , avec chaque cellule ayant une proba 0.01 d'être manquante

$p = 5 \Rightarrow \approx 95\%$  de lignes conservées

$p = 300 \Rightarrow \approx 5\%$  de lignes conservées

# Les données manquantes

A diagram illustrating a data matrix with missing values. The matrix is represented as a light blue rectangle. Above the rectangle, the word "Variables" is centered. Below it, the column indices are labeled: "1", "j", and "p". To the left of the rectangle, the word "Individus" is centered. To its left, the row indices are labeled: "1", "i", and "n". Inside the rectangle, question marks ("?", " ") are placed at various positions to indicate missing data. For example, in the first row, there is a "?" at column 1, and in the first column, there is a "?" at row i. The matrix is not fully populated with question marks, but several are scattered throughout to represent missing data points.

- Etude et mise en œuvre des méthodes factorielles en présence de données manquantes : ACP (variables quantitatives), ACM (variables qualitatives), AFDM (données mixtes), AFM (tableaux multiples)
- Imputation de données

## Exemple sur des données ozone

Code disponible : <http://factominer.free.fr/missMDA/ozone.R>

```
> don <- read.table("http://factominer.free.fr/missMDA/ozoneNA.csv",
  header=TRUE, sep=" ", row.names=1)
```

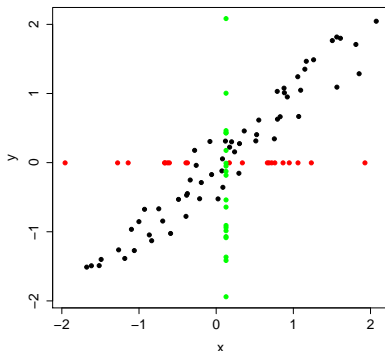
	O3	T9	T12	T15	Ne9	Ne12	Ne15	Vx9	Vx12	Vx15	O3v
0601	82	15.6	18.5	NA	4	4	8	NA	-1.7101	-0.6946	84
0602	82	NA	NA	NA	5	5	7	NA	NA	NA	87
0603	92	NA	17.6	19.5	2	5	4	2.9544	1.8794	0.5209	82
0604	114	16.2	NA	NA	1	1	0	NA	NA	NA	92
0605	94	17.4	20.5	NA	8	8	7	-0.5	NA	-4.3301	114
0606	80	17.7	NA	18.3	NA	NA	NA	-5.6382	-5	-6	94
0607	NA	16.8	15.6	14.9	7	8	8	-4.3301	-1.8794	-3.7588	80
0610	79	14.9	17.5	18.9	5	5	4	0	-1.0419	-1.3892	NA
0611	101	NA	19.6	21.4	2	4	4	-0.766	NA	-2.2981	79
0612	NA	18.3	21.9	22.9	5	6	8	1.2856	-2.2981	-3.9392	101
0613	101	17.3	19.3	20.2	NA	NA	NA	-1.5	-1.5	-0.8682	NA
.	.	.	.	.	.	.	.	.	.	.	.
.	.	.	.	.	.	.	.	.	.	.	.
0919	NA	14.8	16.3	15.9	7	7	7	-4.3301	-6.0622	-5.1962	42
0920	71	15.5	18	17.4	7	7	6	-3.9392	-3.0642	0	NA
0921	96	NA	NA	NA	3	3	3	NA	NA	NA	71
0922	98	NA	NA	NA	2	2	2	4	5	4.3301	96
0923	92	14.7	17.6	18.2	1	4	6	5.1962	5.1423	3.5	98
0924	NA	13.3	17.7	17.7	NA	NA	NA	-0.9397	-0.766	-0.5	92
0925	84	13.3	17.7	17.8	3	5	6	0	-1	-1.2856	NA
0927	NA	16.2	20.8	22.1	6	5	5	-0.6946	-2	-1.3681	71
0928	99	16.9	23	22.6	NA	4	7	1.5	0.8682	0.8682	NA
0929	NA	16.9	19.8	22.1	6	5	3	-4	-3.7588	-4	99
0930	70	15.7	18.6	20.7	NA	NA	NA	0	-1.0419	-4	NA

## De (mauvaises) solutions faciles à mettre en œuvre

- Suppression des données manquantes : rarement intéressant ... mais souvent utilisée (fonction `lm` de R)

## De (mauvaises) solutions faciles à mettre en œuvre

- Suppression des données manquantes : rarement intéressant ... mais souvent utilisée (fonction `lm` de R)
- Imputation par la moyenne (par défaut dans certains logiciels dont FactoMineR)



Distorsion très importante des liaisons entre variables



# Etude du dispositif de données manquantes

Traitement des données manquantes dépend du :

- dispositif de données manquantes : structuré/non structuré
- mécanisme conduisant à l'apparition de données manquantes (Rubin, 1976)

# Etude du dispositif de données manquantes

Traitement des données manquantes dépend du :

- dispositif de données manquantes : structuré/non structuré
- mécanisme conduisant à l'apparition de données manquantes (Rubin, 1976)
  - MCAR : probabilité ne dépend pas de cette valeur ni des autres
  - MAR : probabilité peut dépendre des valeurs d'autres variables
  - MNAR : probabilité dépend de la valeur elle-même

(Ex : Revenu - âge)

# Etude du dispositif de données manquantes

Traitement des données manquantes dépend du :

- dispositif de données manquantes : structuré/non structuré
- mécanisme conduisant à l'apparition de données manquantes (Rubin, 1976)
  - MCAR : probabilité ne dépend pas de cette valeur ni des autres
  - MAR : probabilité peut dépendre des valeurs d'autres variables
  - MNAR : probabilité dépend de la valeur elle-même

(Ex : Revenu - âge)

⇒ Visualisation des données manquantes

## Décompte des valeurs manquantes

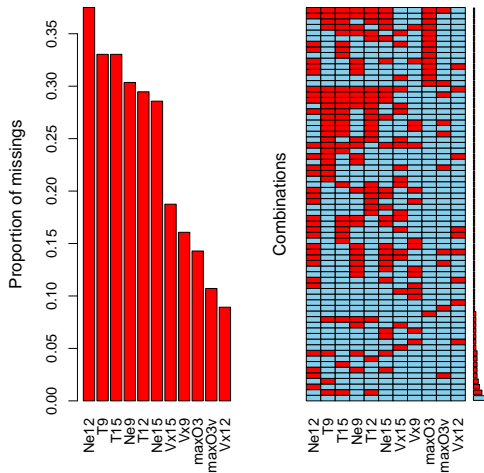
```
> don <- read.table("http://factominer.free.fr/missMDA/ozoneNA.csv",
  header=TRUE, sep=" ", row.names=1)
> library(VIM)
> res <- summary(aggr(don, prop=TRUE, combined=TRUE))$combinations
> res[rev(order(res[,2])),]
```

Variables sorted by

number of missings:

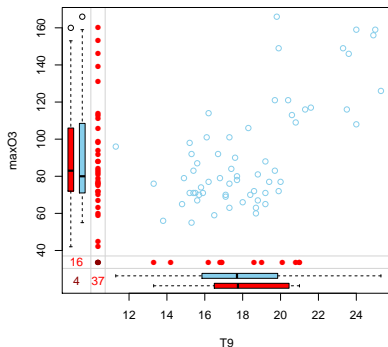
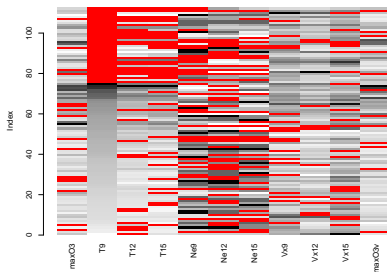
Variable	Count	Combinations	Count	Percent
		0:0:0:0:0:0:0:0:0:0:0	13	11.6071429
Ne12	0.37500000	0:1:1:1:0:0:0:0:0:0:0	7	6.2500000
T9	0.33035714	0:0:0:0:0:1:0:0:0:0:0	5	4.4642857
T15	0.33035714	0:1:0:0:0:0:0:0:0:0:0	4	3.5714286
Ne9	0.30357143	0:1:0:0:1:1:1:0:0:0:0	3	2.6785714
T12	0.29464286	0:0:1:0:0:0:0:0:0:0:0	3	2.6785714
Ne15	0.28571429	0:0:0:1:0:0:0:0:0:0:0	3	2.6785714
Vx15	0.18750000	0:0:0:0:1:1:1:0:0:0:0	3	2.6785714
Vx9	0.16071429	0:0:0:0:0:1:0:0:0:0:1	3	2.6785714
max03	0.14285714	0:1:1:1:1:0:0:0:0:0:0	2	1.7857143
max03v	0.10714286	0:0:0:0:1:0:0:0:0:1:0	2	1.7857143
Vx12	0.08928571	0:0:0:0:0:0:1:1:0:0:0	2	1.7857143
		0:0:0:0:0:0:1:0:0:0:0	2	1.7857143
		.....	.	...

# Visualisation du dispositif de données manquantes



```
> library(VIM)
> aggr(don, only.miss=TRUE, sortVar=TRUE)
```

# Visualisation du dispositif de données manquantes



```
> library(VIM)
> matrixplot(don, sortby=2)
> marginplot(don[, c("T9", "maxO3")])
```

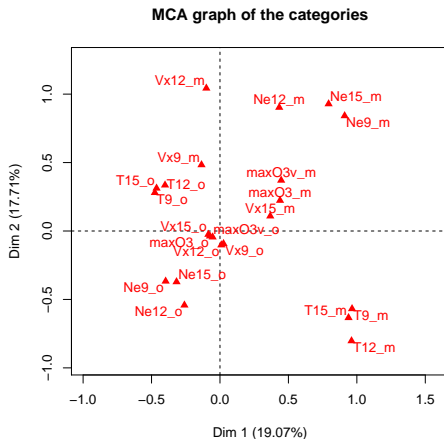
## Visualisation par l'ACM

⇒ Créer une matrice de présence-absence

```
> mis.ind <- matrix("o",nrow=nrow(don),ncol=ncol(don))
> mis.ind[is.na(don)]="m"
> dimnames(mis.ind)=dimnames(don)
> mis.ind
```

	max03	T9	T12	T15	Ne9	Ne12	Ne15	Vx9	Vx12	Vx15	max03v
20010601	"o"	"o"	"o"	"m"	"o"	"o"	"o"	"o"	"o"	"o"	"o"
20010602	"o"	"m"	"m"	"m"	"o"	"o"	"o"	"o"	"o"	"o"	"o"
20010603	"o"	"o"	"o"	"o"	"o"	"m"	"m"	"o"	"m"	"o"	"o"
20010604	"o"	"o"	"o"	"m"	"o"	"o"	"o"	"m"	"o"	"o"	"o"
20010605	"o"	"m"	"o"	"o"	"m"	"m"	"m"	"o"	"o"	"o"	"o"
20010606	"o"	"o"	"o"	"o"	"o"	"m"	"o"	"o"	"o"	"o"	"o"
20010607	"o"	"o"	"o"	"o"	"o"	"o"	"m"	"o"	"o"	"o"	"o"
20010610	"o"	"o"	"o"	"o"	"o"	"o"	"m"	"o"	"o"	"o"	"o"

# Visualisation par l'ACM



```
> library(FactoMineR)
> resMCA <- MCA(mis.ind)
> plot(resMCA,invis="ind",title="MCA graph of the categories")
```



# Approches recommandées pour gérer les valeurs manquantes

⇒ Modifier la méthode, le processus d'estimation pour gérer les données manquantes

⇒ Imputation (multiple) pour obtenir un jeu de données complété à partir duquel toute analyse statistique peut être effectuée

## Algorithme EM (Dempster, Laird et Rubin, 1977)

Principe de l'algorithme d'espérance-maximisation :

- Etape E (Estimation) : remplacer les valeurs manquantes par des valeurs vraisemblables grâce aux données observées et aux paramètres  $\hat{\theta}$  (estimés à l'étape M)
- Etape M (Maximisation de la vraisemblance) : estimation des paramètres  $\theta$  par MV en considérant les données complétées à l'étape E comme de vraies valeurs

Itérer jusqu'à convergence.

La difficulté est de modifier le processus d'estimation permettant de remplacer les valeurs manquantes

## Approche du Maximum de vraisemblance

Hypothèse : en régression  $\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

⇒ Estimation des paramètres du modèle avec EM :

```
> library(norm)
> pre <- prelim.norm(as.matrix(don)) # manipulations préliminaires
> thetahat <- em.norm(pre)           # estimation par MV
> getparam.norm(pre, thetahat)      # résultats
```

## Approche du Maximum de vraisemblance

Hypothèse : en régression  $\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

⇒ Estimation des paramètres du modèle avec EM :

```
> library(norm)
> pre <- prelim.norm(as.matrix(don)) # manipulations préliminaires
> thetahat <- em.norm(pre)           # estimation par MV
> getparam.norm(pre, thetahat)      # résultats
```

⇒ Variances :

- Supplemented EM (Meng, 1991)
- Approche Bootstrap :
  - Bootstrap les lignes :  $\mathbf{X}^1, \dots, \mathbf{X}^B$
  - Algorithme EM :  $(\hat{\boldsymbol{\mu}}^1, \hat{\boldsymbol{\Sigma}}^1), \dots, (\hat{\boldsymbol{\mu}}^B, \hat{\boldsymbol{\Sigma}}^B)$

## Approche du Maximum de vraisemblance

Hypothèse : en régression  $\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

⇒ Estimation des paramètres du modèle avec EM :

```
> library(norm)
> pre <- prelim.norm(as.matrix(don)) # manipulations préliminaires
> thetahat <- em.norm(pre)           # estimation par MV
> getparam.norm(pre, thetahat)      # résultats
```

⇒ Variances :

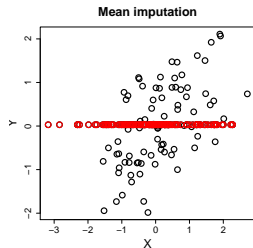
- Supplemented EM (Meng, 1991)
- Approche Bootstrap :
  - Bootstrap les lignes :  $\mathbf{X}^1, \dots, \mathbf{X}^B$
  - Algorithme EM :  $(\hat{\boldsymbol{\mu}}^1, \hat{\boldsymbol{\Sigma}}^1), \dots, (\hat{\boldsymbol{\mu}}^B, \hat{\boldsymbol{\Sigma}}^B)$

**Difficulté** : développer une méthode spécifique pour chaque méthode statistique

# Plan

- 1 Introduction
- 2 Imputation simple pour variables quantitatives**
- 3 Imputation simple pour variables qualitatives
- 4 Imputation simple pour données mixtes
- 5 Imputation multiple
- 6 Conclusion

# Méthodes d'imputation simple



$$\mu_y = 0$$

$$\sigma_y = 1$$

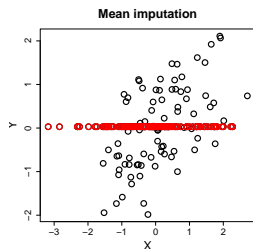
$$\rho = 0.6$$

0.01

0.5

0.30

## Méthodes d'imputation simple

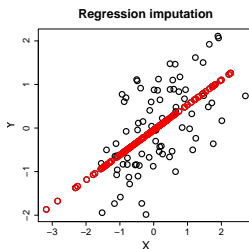


$$\mu_y = 0$$

$$\sigma_y = 1$$

$$\rho = 0.6$$

0.01
0.5
0.30



$$0.01$$

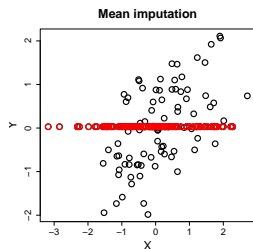
$$0.72$$

$$0.78$$

Imputer par régression  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \implies$  variance sous-estimée, corrélation sur-estimée

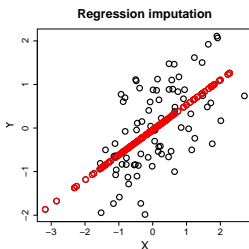


## Méthodes d'imputation simple

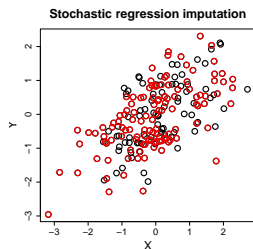


$$\begin{aligned}\mu_y &= 0 \\ \sigma_y &= 1 \\ \rho &= 0.6\end{aligned}$$

0.01
0.5
0.30



0.01
0.72
0.78



0.01
0.99
0.59

Imputer par régression  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \implies$  variance sous-estimée, corrélation sur-estimée

Imputer par **régression aléatoire**  $y_i \sim \mathcal{N}(\hat{\beta}_0 + \hat{\beta}_1 x_i, \hat{\sigma}^2) \implies$

Préserve la distribution

## Modèle joint

⇒ Hypothèse  $\mathbf{x}_i. \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

Généralise la régression aléatoire au cas multivarié :

- Utiliser un algorithme EM pour estimer  $\boldsymbol{\mu}$  et  $\boldsymbol{\Sigma}$  à partir d'un jeu incomplet
- Tirer à partir de  $\mathcal{N}(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}})$

```
> library(norm)
> pre <- prelim.norm(as.matrix(don))
> thetahat <- em.norm(pre)
> rngseed(123)
> imp <- imp.norm(pre, thetahat, don)
```

## Modèle conditionnel

⇒ Un modèle par variable

Exemple avec régression :

- 1 Initialisation de l'imputation : imputation par la moyenne
- 2 Ajuster une régression aléatoire de  $\mathbf{X}_j^{obs}$  en fonction des autres variables  $\mathbf{X}_{-j}^{obs}$   
Prédire  $\mathbf{X}_j^{miss}$  à partir du modèle ajusté
- 3 Boucler sur les variables

```
> library(mice)
> res.cm <- mice(don, m=1)
```

## Modèle conditionnel

⇒ Un modèle par variable

Exemple avec régression :

- 1 Initialisation de l'imputation : imputation par la moyenne
- 2 Ajuster une régression aléatoire de  $\mathbf{X}_j^{obs}$  en fonction des autres variables  $\mathbf{X}_{-j}^{obs}$   
Prédire  $\mathbf{X}_j^{miss}$  à partir du modèle ajusté
- 3 Boucler sur les variables

```
> library(mice)
> res.cm <- mice(don, m=1)
```

⇒ Flexibilité : différents modèles pour chaque variable

## Autres méthodes d'imputation simple

- k-plus proches voisins (`class`, `FNN`)
- forêts aléatoires (`missForest`, Stekhoven & Bühlmann, 2011)
- ...

⇒ R CRAN task View: Missing Data

⇒ R-miss-tastic

⇒ Imputation par ACP

## Ajustement du nuage en ACP

Trouver le sous-espace qui fournit la meilleure représentation des données

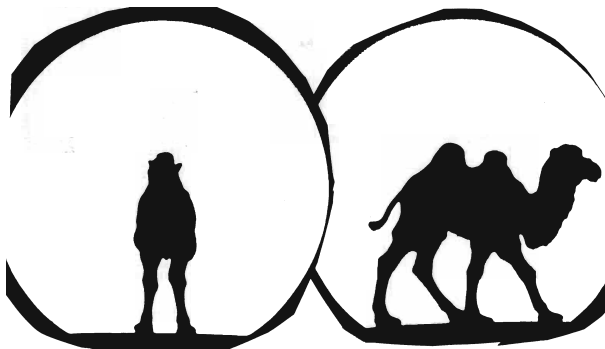


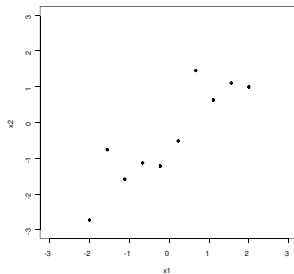
FIGURE – Chameau ou dromadaire ? source J.P. Fenelon

- ⇒ Meilleure approximation par projection
- ⇒ Meilleure représentation de la diversité, de la variabilité

## Ajustement du nuage en ACP

**X**

-2.00	-2.36
-1.56	-0.39
-1.11	-1.21
-0.67	-0.75
-0.22	-0.84
0.22	-0.14
0.67	1.84
1.11	1.01
1.56	1.48
2.00	1.38



⇒ Minimisation de la distance entre les individus et leur projection

⇒ Minimise  $||X - \hat{X}||$

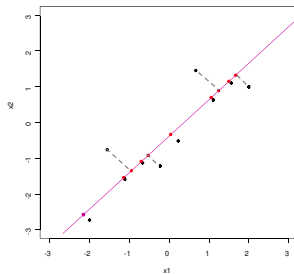
# Ajustement du nuage en ACP

**X**

-2.00	-2.36
-1.56	-0.39
-1.11	-1.21
-0.67	-0.75
-0.22	-0.84
0.22	-0.14
0.67	1.84
1.11	1.01
1.56	1.48
2.00	1.38

-2.16	-2.21
-0.96	-0.98
-1.15	-1.17
-0.70	-0.72
-0.53	-0.54
0.04	0.04
1.25	1.27
1.05	1.07
1.50	1.54
1.67	1.70

**$\hat{X}$**

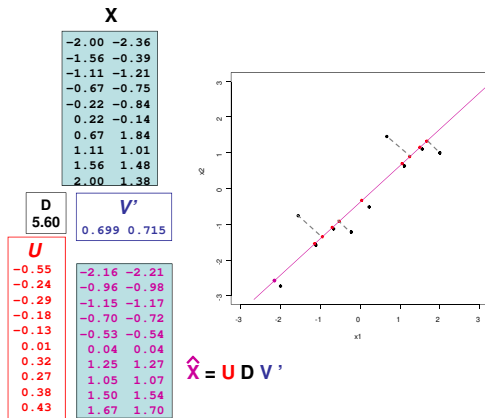


⇒ Minimisation de la distance entre les individus et leur projection

⇒ Minimise  $||X - \hat{X}||$



# Reconstitution en ACP



$\Rightarrow \hat{\mathbf{X}} = \mathbf{U} \mathbf{D} \mathbf{V}'$  (produit matriciel utilisant coordonnées des individus et des variables obtenues par ACP ;  $\mathbf{X}$  est supposé centré)

## ACP : cas complet

⇒ Point de vue géométrique : minimiser l'erreur de reconstitution

⇒ Approximation de  $\mathbf{X}$  par une matrice de rang  $S < p$  :

$$\|\mathbf{X}_{n \times p} - \hat{\mathbf{X}}_{n \times p}\|^2 \quad \text{SVD : } \hat{\mathbf{X}}^{\text{ACP}} = \mathbf{U}_{n \times S} \mathbf{D}_{S \times S} \mathbf{V}'_{p \times S}$$

$\mathbf{F} = \mathbf{U}\mathbf{D}$  composantes principales (scores)

$\mathbf{V}$  axes principaux (loadings)

## ACP : cas complet

⇒ Point de vue géométrique : minimiser l'erreur de reconstitution

⇒ Approximation de  $\mathbf{X}$  par une matrice de rang  $S < p$  :

$$\|\mathbf{X}_{n \times p} - \hat{\mathbf{X}}_{n \times p}\|^2 \quad \text{SVD : } \hat{\mathbf{X}}^{\text{ACP}} = \mathbf{U}_{n \times S} \mathbf{D}_{S \times S} \mathbf{V}'_{p \times S}$$

$\mathbf{F} = \mathbf{U}\mathbf{D}$  composantes principales (scores)

$\mathbf{V}$  axes principaux (loadings)

⇒ Point de vue modèle à effets fixes (Causinus, 1986)

$$\mathbf{X}_{n \times p} = \tilde{\mathbf{X}}_{n \times p} + \varepsilon_{n \times p}$$

$$x_{ij} = \sum_{s=1}^S d_s u_{is} v_{js} + \varepsilon_{ij} \quad \varepsilon_{ij} \sim \mathcal{N}(0, \sigma^2)$$

Estimateurs de maximum de vraisemblance = estimateurs des moindres carrés

## Imputation par ACP

⇒ ACP : moindres carrés

$$\|\mathbf{X}_{n \times p} - \mathbf{U}_{n \times S} \mathbf{D}_{S \times S} \mathbf{V}'_{p \times S}\|^2$$

## Imputation par ACP

⇒ ACP : moindres carrés

$$\|\mathbf{X}_{n \times p} - \mathbf{U}_{n \times S} \mathbf{D}_{S \times S} \mathbf{V}'_{p \times S}\|^2$$

⇒ ACP avec données manquantes : moindres carrés pondérés

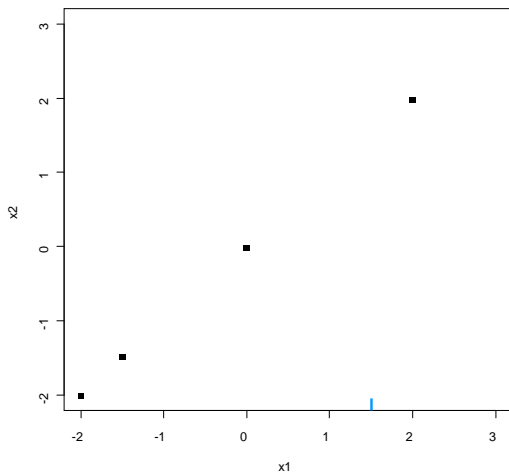
$$\|\mathbf{R}_{n \times p} * (\mathbf{X}_{n \times p} - \mathbf{U}_{n \times S} \mathbf{D}_{S \times S} \mathbf{V}'_{p \times S})\|^2$$

with  $r_{ij} = 0$  si  $x_{ij}$  manquant,  $r_{ij} = 1$  sinon

Beaucoup d'algorithmes : moindres carrés pondérés alterné (Gabriel & Zamir, 1979); ACP iterative (Kiers, 1997)

## ACP itérative

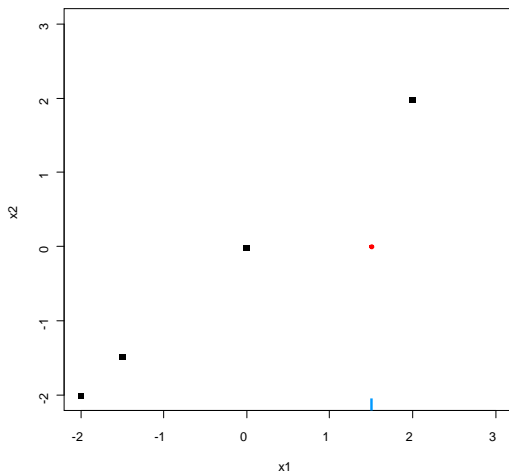
x1	x2
-2.0	-2.01
-1.5	-1.48
0.0	-0.01
1.5	NA
2.0	1.98



## ACP itérative

x1	x2
-2.0	-2.01
-1.5	-1.48
0.0	-0.01
1.5	NA
2.0	1.98

x1	x2
-2.0	-2.01
-1.5	-1.48
0.0	-0.01
1.5	0.00
2.0	1.98



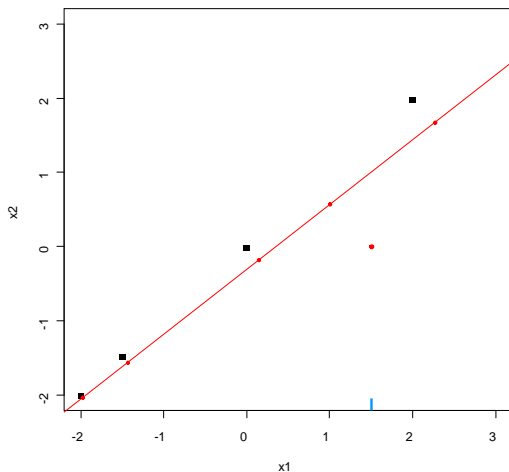
Initialisation  $\ell = 0$  :  $\mathbf{X}^0$  (imputation par la moyenne)

## ACP itérative

x1	x2
-2.0	-2.01
-1.5	-1.48
0.0	-0.01
1.5	NA
2.0	1.98

x1	x2
-2.0	-2.01
-1.5	-1.48
0.0	-0.01
1.5	0.00
2.0	1.98

$\hat{x}_1$	$\hat{x}_2$
-1.98	-2.04
-1.44	-1.56
0.15	-0.18
1.00	0.57
2.27	1.67



ACP sur le jeu de données complété  $\rightarrow (\mathbf{U}^\ell, \mathbf{D}^\ell, \mathbf{V}^\ell)$ ;

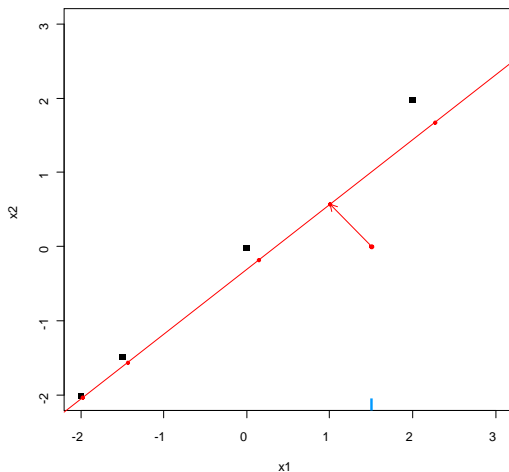


## ACP itérative

x1	x2
-2.0	-2.01
-1.5	-1.48
0.0	-0.01
1.5	NA
2.0	1.98

x1	x2
-2.0	-2.01
-1.5	-1.48
0.0	-0.01
1.5	0.00
2.0	1.98

$\hat{x}_1$	$\hat{x}_2$
-1.98	-2.04
-1.44	-1.56
0.15	-0.18
1.00	0.57
2.27	1.67



Valeurs manquantes imputées par le modèle  $\hat{\mathbf{X}}^\ell = \mathbf{U}^\ell \mathbf{D}^\ell \mathbf{V}^{\ell'}$

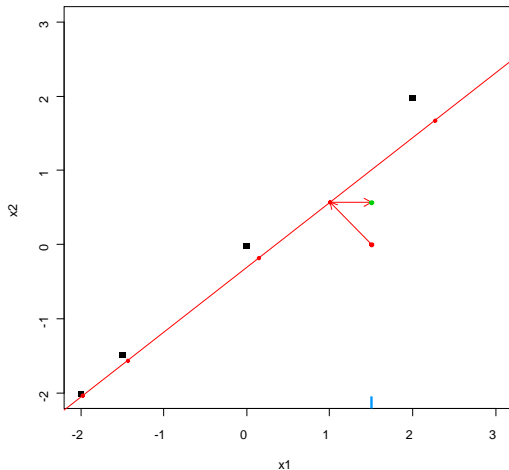
## ACP itérative

x1	x2
-2.0	-2.01
-1.5	-1.48
0.0	-0.01
1.5	NA
2.0	1.98

x1	x2
-2.0	-2.01
-1.5	-1.48
0.0	-0.01
1.5	0.00
2.0	1.98

$\hat{x}_1$	$\hat{x}_2$
-1.98	-2.04
-1.44	-1.56
0.15	-0.18
1.00	0.57
2.27	1.67

x1	x2
-2.0	-2.01
-1.5	-1.48
0.0	-0.01
1.5	0.57
2.0	1.98



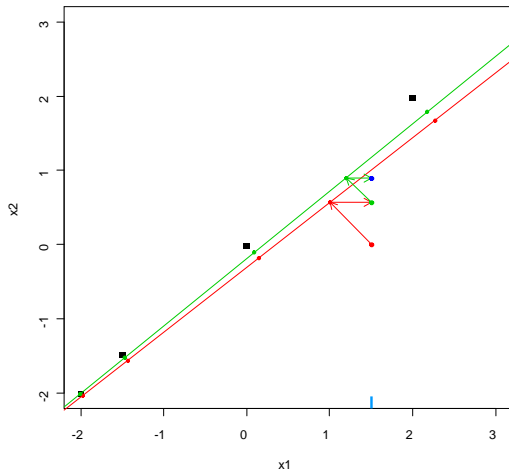
Nouveau jeu de données imputé  $\mathbf{X}^\ell = \mathbf{R} * \mathbf{X} + (1 - \mathbf{R}) * \hat{\mathbf{X}}^\ell$

# ACP itérative

x1	x2
-2.0	-2.01
-1.5	-1.48
0.0	-0.01
1.5	NA
2.0	1.98

x1	x2
-2.0	-2.01
-1.5	-1.48
0.0	-0.01
1.5	0.57
2.0	1.98

x1	x2
-2.0	-2.01
-1.5	-1.48
0.0	-0.01
1.5	0.57
2.0	1.98



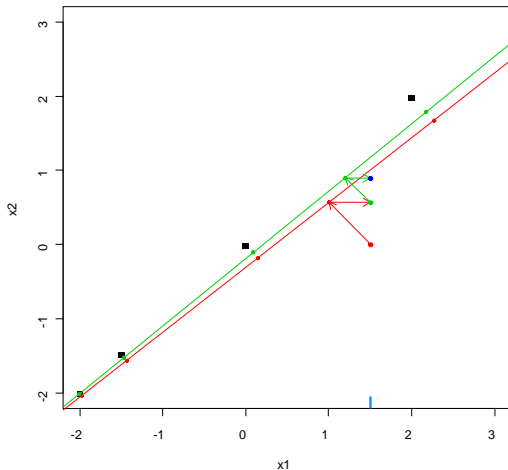
## ACP itérative

x1	x2
-2.0	-2.01
-1.5	-1.48
0.0	-0.01
1.5	NA
2.0	1.98

x1	x2
-2.0	-2.01
-1.5	-1.48
0.0	-0.01
1.5	0.57
2.0	1.98

$\hat{x}_1$	$\hat{x}_2$
-2.00	-2.01
-1.47	-1.52
0.09	-0.11
1.20	0.90
2.18	1.78

x1	x2
-2.0	-2.01
-1.5	-1.48
0.0	-0.01
1.5	0.90
2.0	1.98



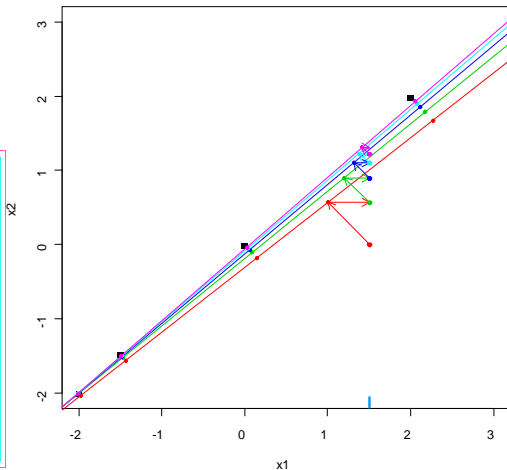
## ACP itérative

x1	x2
-2.0	-2.01
-1.5	-1.48
0.0	-0.01
1.5	NA
2.0	1.98

x1	x2
-2.0	-2.01
-1.5	-1.48
0.0	-0.01
1.5	0.00
2.0	1.98

$\hat{x}_1$	$\hat{x}_2$
-1.98	-2.04
-1.44	-1.56
0.15	-0.18
1.00	0.57
2.27	1.67

x1	x2
-2.0	-2.01
-1.5	-1.48
0.0	-0.01
1.5	0.57
2.0	1.98



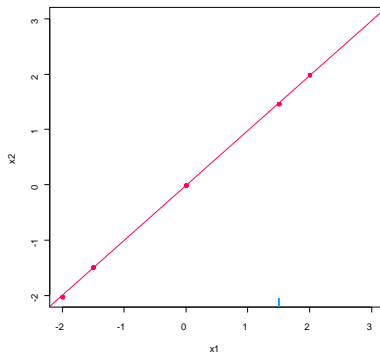
Les étapes sont répétées jusqu'à convergence

## ACP itérative

```

x1    x2
-2.0 -2.01
-1.5 -1.48
0.0  -0.01
1.5   NA
2.0  1.98

```



```

x1    x2
-2.0 -2.01
-1.5 -1.48
0.0  -0.01
1.5  1.46
2.0  1.98

```

ACP sur le jeu de données complété  $\rightarrow (\mathbf{U}^\ell, \mathbf{D}^\ell, \mathbf{V}^\ell)$   
 Valeurs manquantes imputées par le modèle  $\hat{\mathbf{X}}^\ell = \mathbf{U}^\ell \mathbf{D}^\ell \mathbf{V}^{\ell'}$

## ACP itérative

- ① initialisation  $\ell = 0$  :  $\mathbf{X}^0$  (imputation par la moyenne)
- ② step  $\ell$  :
  - (a) ACP sur le tableau complété  $\rightarrow (\mathbf{U}^\ell, \mathbf{D}^\ell, \mathbf{V}^\ell)$ ;  
 $S$  dimensions conservées
  - (b) valeurs manquantes imputées par  $\hat{\mathbf{X}}^\ell = \mathbf{U}^\ell \mathbf{D}^\ell \mathbf{V}^{\ell \prime}$ ;  
nouveau tableau imputé  $\mathbf{X}^\ell = \mathbf{R} * \mathbf{X} + (1 - \mathbf{R}) * \hat{\mathbf{X}}^\ell$
- ③ étapes répétées jusqu'à convergence

## ACP itérative

- ① initialisation  $\ell = 0$  :  $\mathbf{X}^0$  (imputation par la moyenne)
- ② step  $\ell$  :
  - (a) ACP sur le tableau complété  $\rightarrow (\mathbf{U}^\ell, \mathbf{D}^\ell, \mathbf{V}^\ell)$ ;  
 $S$  dimensions conservées
  - (b) valeurs manquantes imputées par  $\hat{\mathbf{X}}^\ell = \mathbf{U}^\ell \mathbf{D}^\ell \mathbf{V}^{\ell \prime}$ ;  
nouveau tableau imputé  $\mathbf{X}^\ell = \mathbf{R} * \mathbf{X} + (1 - \mathbf{R}) * \hat{\mathbf{X}}^\ell$
  - (c) moyennes (et écarts-types) sont mis à jour
- ③ étapes répétées jusqu'à convergence



## ACP itérative

- ① initialisation  $\ell = 0$  :  $\mathbf{X}^0$  (imputation par la moyenne)
- ② step  $\ell$  :
  - (a) ACP sur le tableau complété  $\rightarrow (\mathbf{U}^\ell, \mathbf{D}^\ell, \mathbf{V}^\ell)$ ;  
*S dimensions conservées*
  - (b) valeurs manquantes imputées par  $\hat{\mathbf{X}}^\ell = \mathbf{U}^\ell \mathbf{D}^\ell \mathbf{V}^{\ell'}$ ;  
nouveau tableau imputé  $\mathbf{X}^\ell = \mathbf{R} * \mathbf{X} + (1 - \mathbf{R}) * \hat{\mathbf{X}}^\ell$
  - (c) moyennes (et écarts-types) sont mis à jour
- ③ étapes répétées jusqu'à convergence

## ACP itérative

- ① initialisation  $\ell = 0$  :  $\mathbf{X}^0$  (imputation par la moyenne)
- ② step  $\ell$  :
  - (a) ACP sur le tableau complété  $\rightarrow (\mathbf{U}^\ell, \mathbf{D}^\ell, \mathbf{V}^\ell)$ ;  
**S dimensions conservées**
  - (b) valeurs manquantes imputées par  $\hat{\mathbf{X}}^\ell = \mathbf{U}^\ell \mathbf{D}^\ell \mathbf{V}^{\ell'}$ ;  
nouveau tableau imputé  $\mathbf{X}^\ell = \mathbf{R} * \mathbf{X} + (1 - \mathbf{R}) * \hat{\mathbf{X}}^\ell$
  - (c) **moyennes (et écarts-types) sont mis à jour**
- ③ étapes répétées jusqu'à convergence

$\Rightarrow$  algorithme EM pour le modèle à effets fixes

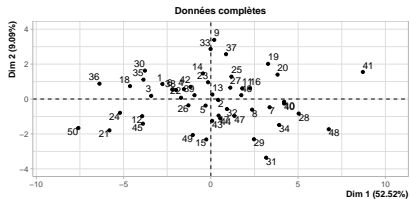
$\Rightarrow$  Imputation (complétion de matrice, Netflix)

$\Rightarrow$  Réduction de la variabilité (imputation par  $\mathbf{UDV}'$ )

$\Rightarrow$  Problème de surajustement

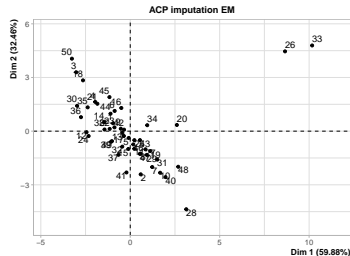
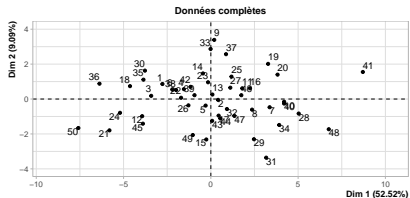
## Surajustement

$$X_{50 \times 10} = \mathbf{U}_{50 \times 2} \mathbf{D}\mathbf{V}'_{10 \times 2} + \mathcal{N}(0, 0.5);$$



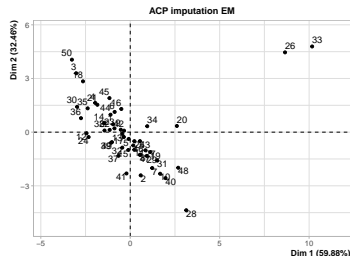
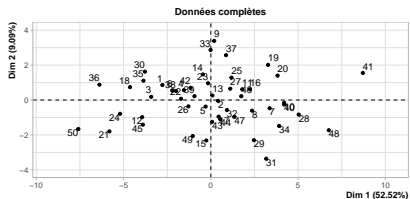
# Surajustement

$$X_{50 \times 10} = \mathbf{U}_{50 \times 2} \mathbf{D} \mathbf{V}'_{10 \times 2} + \mathcal{N}(0, 0.5); \text{ 50\% of NA}$$



## Surajustement

$$\mathbf{X}_{50 \times 10} = \mathbf{U}_{50 \times 2} \mathbf{D} \mathbf{V}'_{10 \times 2} + \mathcal{N}(0, 0.5); \text{ 50\% of NA}$$



⇒ erreur d'ajustement faible :  $\|\mathbf{R} * (\mathbf{X} - \hat{\mathbf{X}})\|^2 = 0.50$

⇒ erreur de prédiction élevée :  $\|(1 - \mathbf{R}) * (\mathbf{X} - \hat{\mathbf{X}})\|^2 = 16.98$

## Surajustement

⇒ Bon ajustement et mauvaise prédiction

- Trop de paramètres sont estimés par rapport au nombre de données observées : le nombre de dimension  $S$  et le nombre de données manquantes sont grands
- Faibles liaisons entre variables

① Diminuer le nombre  $S$

② Early stopping

③ Régularisation ⇒ ACP itérative régularisée

## ACP itérative régularisée (Josse *et al.*, 2009)

⇒ Initialisation - étape d'estimation - étape d'imputation

L'étape d'imputation :

$$\hat{x}_{ij}^{\text{ACP}} = \sum_{s=1}^S d_s u_{is} v_{js}$$

est remplacée par une étape d'imputation régularisée :

$$\hat{x}_{ij}^{\text{rACP}} = \sum_{s=1}^S \left( \frac{d_s^2 - \hat{\sigma}^2}{d_s^2} \right) d_s u_{is} v_{js} = \sum_{s=1}^S \left( d_s - \frac{\hat{\sigma}^2}{d_s} \right) u_{is} v_{js}$$

## ACP itérative régularisée (Josse *et al.*, 2009)

⇒ Initialisation - étape d'estimation - étape d'imputation

L'étape d'imputation :

$$\hat{x}_{ij}^{\text{ACP}} = \sum_{s=1}^S d_s u_{is} v_{js}$$

est remplacée par une étape d'imputation régularisée :

$$\hat{x}_{ij}^{\text{rACP}} = \sum_{s=1}^S \left( \frac{d_s^2 - \hat{\sigma}^2}{d_s^2} \right) d_s u_{is} v_{js} = \sum_{s=1}^S \left( d_s - \frac{\hat{\sigma}^2}{d_s} \right) u_{is} v_{js}$$

$$\hat{\sigma}^2 = \frac{RSS}{ddl} = \frac{n \sum_{s=S+1}^q d_s^2}{np - p - nS - pS + S^2 + S} \quad (\mathbf{X}_{n \times p}; \mathbf{U}_{n \times S}; \mathbf{V}_{p \times S})$$



## ACP itérative régularisée (Josse *et al.*, 2009)

⇒ Initialisation - étape d'estimation - étape d'imputation

L'étape d'imputation :

$$\hat{x}_{ij}^{\text{ACP}} = \sum_{s=1}^S d_s u_{is} v_{js}$$

est remplacée par une étape d'imputation régularisée :

$$\hat{x}_{ij}^{\text{rACP}} = \sum_{s=1}^S \left( \frac{d_s^2 - \hat{\sigma}^2}{d_s^2} \right) d_s u_{is} v_{js} = \sum_{s=1}^S \left( d_s - \frac{\hat{\sigma}^2}{d_s} \right) u_{is} v_{js}$$

$$\hat{\sigma}^2 = \frac{RSS}{ddl} = \frac{n \sum_{s=S+1}^q d_s^2}{np - p - nS - pS + S^2 + S} \quad (\mathbf{X}_{n \times p}; \mathbf{U}_{n \times S}; \mathbf{V}_{p \times S})$$

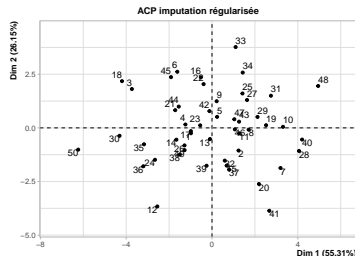
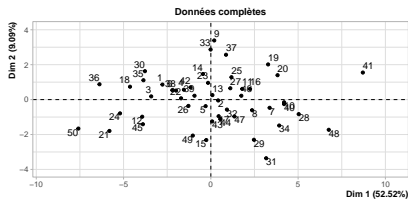
Compromis seuillage doux/dur (Mazumder, Hastie & Tibshirani, 2010)

$\sigma^2$  petit → ACP régularisée  $\approx$  ACP

$\sigma^2$  grand → imputation par la moyenne

## Surajustement

$$\mathbf{X}_{50 \times 10} = \mathbf{U}_{50 \times 2} \mathbf{D}\mathbf{V}'_{10 \times 2} + \mathcal{N}(0, 0.5); \text{ 50\% of NA}$$



$\Rightarrow$  erreur d'ajustement :  $\|\mathbf{R} * (\mathbf{X} - \hat{\mathbf{X}})\|^2 = 0.56$  (EM= 0.50)

$\Rightarrow$  erreur de prédiction :  $\|(1 - \mathbf{R}) * (\mathbf{X} - \hat{\mathbf{X}})\|^2 = 2.28$  (EM= 16.98)

## Propriétés de l'imputation

- Résultats de l'ACP obtenus à partir des données observées uniquement : graphe des individus et graphe des variables  
⇒ On "saute" les données manquantes, l'ACP itérative minimise  $\|\mathbf{R} * (\mathbf{X} - \mathbf{UDV}')\|^2$
- Bonne qualité d'imputation quand la structure dans le jeu de données est forte (imputation utilisant les ressemblances entre individus et les liaisons entre variables)
- le tableau imputé peut être utilisé (avec précaution) pour réaliser d'autres analyses
- Bien meilleur que l'algorithme Nipals (encore trop utilisé)
- Compétitif par rapport aux forêts aléatoires

# Imputation par ACP en pratique

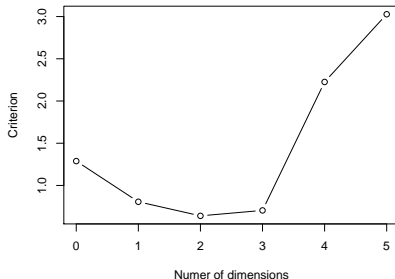
## Tutoriel sur l'ACP avec données manquantes

(données ozone, lignes de code)

⇒ Etape 1 : Estimation du nombre de dimensions

(Validation croisée, Bro, 2008 ; GCV, Josse & Husson, 2011)

```
> library(missMDA)
> nb <- estim_ncpPCA(don, method.cv="Kfold")
> nb$ncp      #2
> plot(0:5, nb$criterion, xlab="nb dim", ylab="MSEP")
```



# Imputation par ACP en pratique

⇒ Etape 2 : Imputation des données manquantes

```
> res.comp <- imputePCA(don, ncp = 2)
```

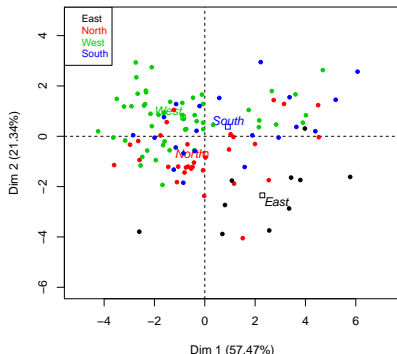
```
> res.comp$completeObs[1:3,]
```

	max03	T9	T12	T15	Ne9	Ne12	Ne15	Vx9	Vx12	Vx15	max03v
0601	87	15.60	18.50	20.47	4	4.00	8.00	0.69	-1.71	-0.69	84
0602	82	18.51	20.88	21.81	5	5.00	7.00	-4.33	-4.00	-3.00	87
0603	92	15.30	17.60	19.50	2	3.98	3.81	2.95	1.97	0.52	82

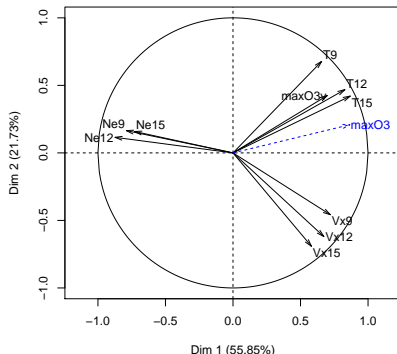
## ACP sur le tableau complété

⇒ Etape 3 : ACP sur le tableau complété

Individuals factor map (PCA)



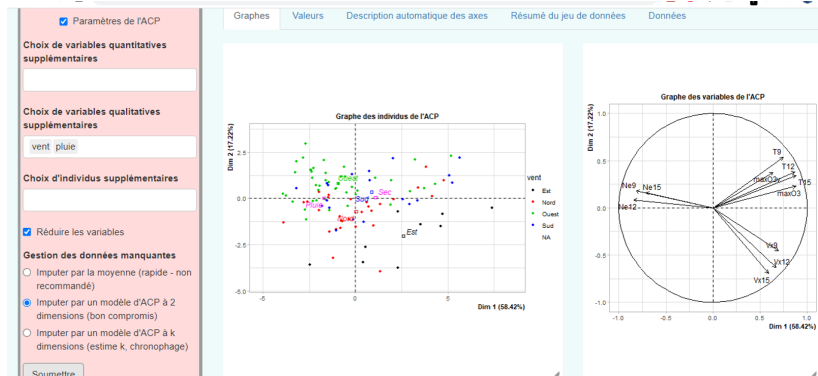
Variables factor map (PCA)



```
> imp <- cbind.data.frame(res.comp$completeObs, ozone[,12])
> res.pca <- PCA(imp, quanti.sup=1, quali.sup=12)
> plot(res.pca, hab=12, lab="quali")
> plot(res.pca, choix="var")
> res.pca$ind$coord #scores (principal components)
```

## 3 en 1 avec le package Factoshiny

- > library(Factoshiny)
- > Factoshiny(ozone)



## Jeu de données en écologie

Données Glopnet : 2494 espèces décrites par 6 variables quantitatives ([données](#), [lignes de code](#))

- LMA (leaf mass per area)
- LL (leaf lifespan)
- Amass (photosynthetic assimilation)
- Nmass (leaf nitrogen)
- Pmass (leaf phosphorus)
- Rmass (dark respiration rate)

et 1 variable qualitative : le biome (macro-écosystème)

Wright IJ, et al. (2004). The worldwide leaf economics spectrum.  
*Nature*, 428 :821.

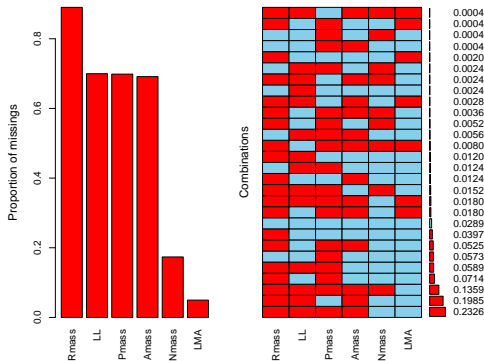
[www.nature.com/nature/journal/v428/n6985/extref/nature02403-s2.xls](http://www.nature.com/nature/journal/v428/n6985/extref/nature02403-s2.xls)



## Jeu de données en écologie

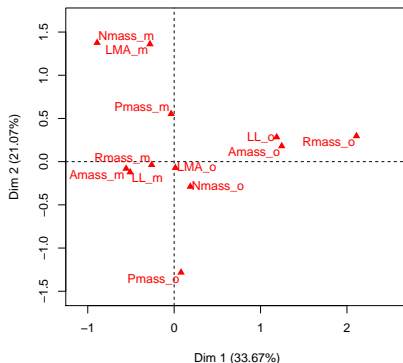
```
> sum(is.na(don))/(nrow(don)*ncol(don)) # 53% de données manquantes
[1] 0.5338145
> dim(na.omit(don)) ## suppression des espèces avec données manquantes
[1] 72 6
      ## reste seulement 72 espèces!

> library(VIM)
> aggr(don,numbers=TRUE,sortVar=TRUE)
```



## Jeu de données en écologie

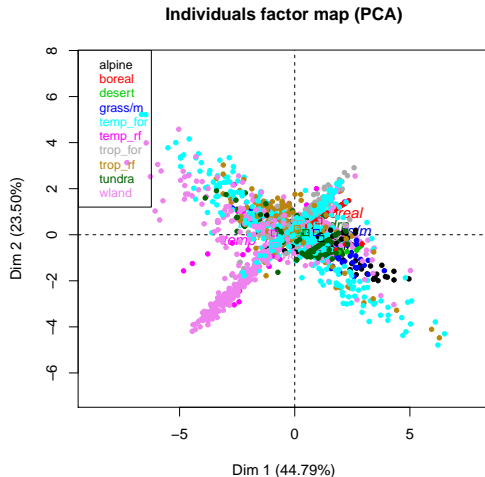
MCA graph of the categories



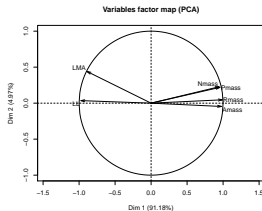
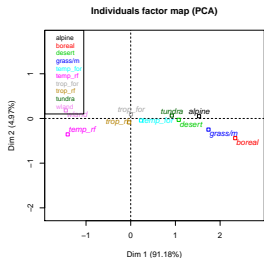
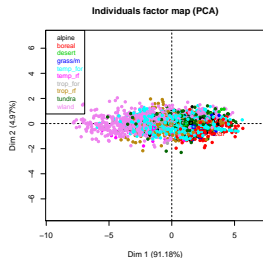
```
> mis.ind <- matrix("o",nrow=nrow(don),ncol=ncol(don))  
> mis.ind[is.na(don)] <- "m"  
> dimnames(mis.ind) <- dimnames(don)  
> library(FactoMineR)  
> resMCA <- MCA(mis.ind)  
> plot(resMCA,invis="ind",title="MCA graph of the categories")
```

# Jeu de données en écologie

Quid de l'imputation par la moyenne ?



# Jeu de données en écologie



```
> library(missMDA)
> nb <- estim_ncpPCA(don,method.cv="Kfold",nbsim=100)
> res.comp <- imputePCA(don,ncp=2)
> imp <- cbind.data.frame(res.comp$completeObs,tab.init[,1:4])
> res.pca <- PCA(imp,quanti.sup=1,quali.sup=12)
> plot(res.pca, hab=12, lab="quali")
> plot(res.pca, choix="var")
> res.pca$ind$coord #scores (principal components)
```

# Plan

- 1 Introduction
- 2 Imputation simple pour variables quantitatives
- 3 Imputation simple pour variables qualitatives**
- 4 Imputation simple pour données mixtes
- 5 Imputation multiple
- 6 Conclusion

## Imputation simple basée sur l'ACM

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

L'ACM peut être vue comme une ACP de la matrice indicatrice  $\mathbf{X}$  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)

$$X = \begin{array}{c|c|c|c|c} & I_1 & I_k & \dots & I_K \\ \hline I_1 & 1 & 0 & 0 & 0 & 1 \\ I_2 & 1 & 0 & 0 & 1 & 0 & \dots & NA & NA \\ I_3 & NA & NA & NA & 0 & 1 & 0 & 0 & \dots & 0 & 1 \\ I_4 & 1 & 0 & 0 & 1 & 0 & 0 & 1 & \dots & 0 & 1 \end{array} \begin{array}{l} J \\ J \\ J \\ J \\ J \\ J \\ J \\ J \\ J \\ J \\ J \end{array}$$

$x_{ik}$

$I_1 \quad I_k \quad I_K \quad IJ$

$$D_{\Sigma} = \begin{array}{c} I_1 \\ \dots \\ I_k \\ \dots \\ I_K \end{array} \begin{array}{c} 0 \\ \dots \\ 0 \end{array}$$

## 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
  - (b) Imputation des données manquantes par les données reconstituées
  - (c) Mise à jour des marges

	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	0.71	0.29	1	0	...
ind 2	0.12	0.29	0.59	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	0.37	0.63	...
...	...	...	...	...	...	...	...	...
ind 1232	0	0	1	0	1	0	1	...

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

## Imputation de la matrice indicatrice

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

	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	0.71	0.29	1	0	...
ind 2	0.12	0.29	0.59	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	0.37	0.63	...
...	...	...	...	...	...	...	...	...
ind 1232	0	0	1	0	1	0	1	...

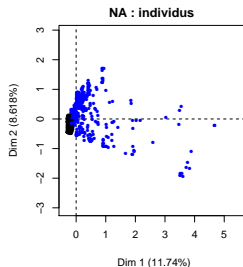
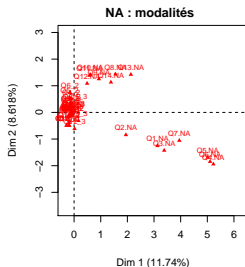
ACM sur le tableau complété (utilisation de l'argument `tab.disj`)

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



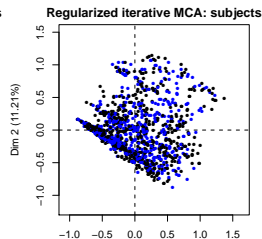
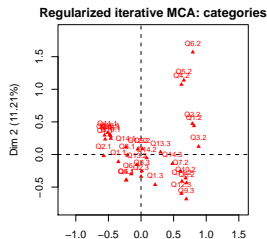
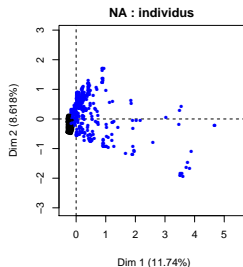
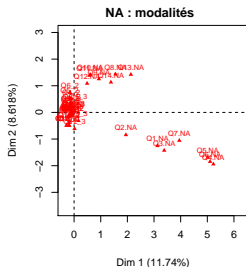
## Traitement d'un questionnaire

- 1232 répondants, 14 questions, 35 modalités, 9% de données manquantes pour 42% des répondants



## Traitement d'un questionnaire

- 1232 répondants, 14 questions, 35 modalités, 9% de données manquantes pour 42% des répondants



# Plan

- 1 Introduction
- 2 Imputation simple pour variables quantitatives
- 3 Imputation simple pour variables qualitatives
- 4 Imputation simple pour données mixtes**
- 5 Imputation multiple
- 6 Conclusion

## Données mixtes

⇒ Modèle joint :

- General location model (Schafer, 1997)  $\implies$  problème quand beaucoup de modalités
- Transformer les variables qualitatives en indicatrices et faire comme si les variables étaient continues (**Amelia**)
- Modèle à classes latentes (Vermunt) – modèles Bayésien non paramétrique (Dunson, Reiter, Duke University)

⇒ Modèle conditionnel : linéaire, logistique, multinomial, logit (**mice**)

## Données mixtes

⇒ Modèle joint :

- General location model (Schafer, 1997)  $\implies$  problème quand beaucoup de modalités
- Transformer les variables qualitatives en indicatrices et faire comme si les variables étaient continues (**Amelia**)
- Modèle à classes latentes (Vermunt) – modèles Bayésien non paramétrique (Dunson, Reiter, Duke University)

⇒ Modèle conditionnel : linéaire, logistique, multinomial, logit (**mice**)

⇒ Forêts aléatoires (Stekhoven & Bühlmann, 2012, **missForest**)

⇒ Analyse factorielle de données mixtes (Audigier, Husson & Josse, 2014, **missMDA**)

# Imputation itérative par forêts aléatoires

- ① Imputation initiale : moyenne - modalité au hasard
- ② Ajuster une forêt aléa  $\mathbf{X}_j^{obs}$  en fct de  $\mathbf{X}_{-j}^{obs}$  puis prédire  $\mathbf{X}_j^{miss}$
- ③ Boucler sur les variables jusqu'à un critère d'arrêt

## Imputation itérative par forêts aléatoires

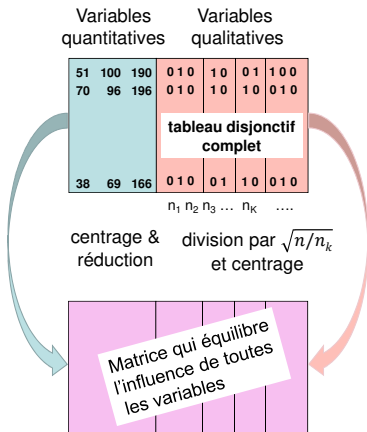
- ➊ Imputation initiale : moyenne - modalité au hasard
- ➋ Ajuster une forêt aléa  $\mathbf{X}_j^{obs}$  en fct de  $\mathbf{X}_{-j}^{obs}$  puis prédire  $\mathbf{X}_j^{miss}$
- ➌ Boucler sur les variables jusqu'à un critère d'arrêt

⇒ Propriétés :

- Relations non-linéaires, interactions complexes
- $n \ll p$
- erreur out-of-bag : approximation de l'erreur d'imputation

⇒ Meilleur que plus proches voisins et mice

# Analyse Factorielle de Données Mixtes (cas complet)

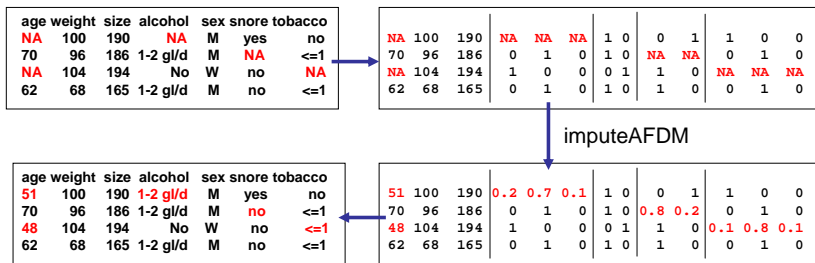


ACP sur une matrice pondérée



## Algorithme d'AFDM itératif régularisé

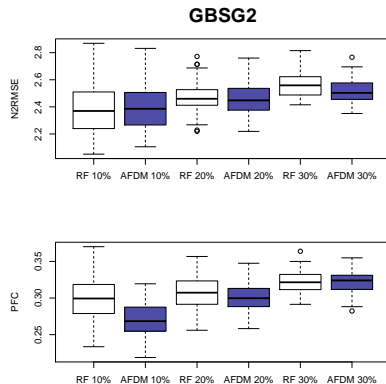
- 1 Initialisation : imputation par la moyenne (quanti) et la proportion (quali)
- 2 Itérer jusqu'à convergence
  - (a) estimation : AFDM sur le jeu complété  $\Rightarrow \mathbf{U}, \mathbf{D}, \mathbf{V}$
  - (b) imputation des valeurs manquantes avec le modèle de reconstitution
  - (c) moyennes,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



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

## Comparaison avec forêts aléatoires

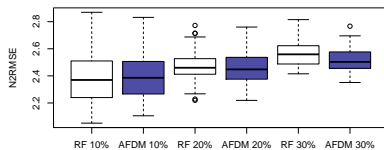
Imputations obtenues par forêts aléatoires & ACP itérative



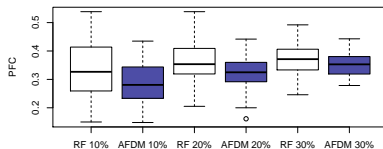
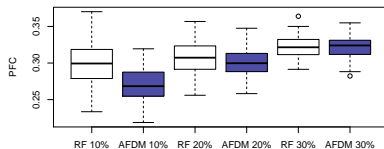
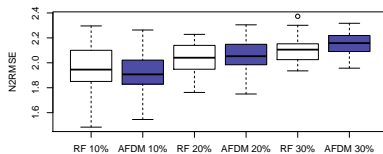
# Comparaison avec forêts aléatoires

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 (complémentaire de FactoMineR)

⇒ 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)

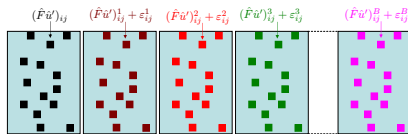
# Plan

- 1 Introduction
- 2 Imputation simple pour variables quantitatives
- 3 Imputation simple pour variables qualitatives
- 4 Imputation simple pour données mixtes
- 5 Imputation multiple**
- 6 Conclusion

## Imputation multiple

Imputation simple : une valeur unique ne peut pas refléter l'incertitude sur la prédiction  $\Rightarrow$  sous-estimation de l'écart-type

- 1 Générer  $M$  valeurs possibles pour chaque valeur manquante



- 2 Faire l'analyse sur chaque tableau imputé :  $\hat{\theta}_m, \widehat{Var}(\hat{\theta}_m)$
- 3 Combiner les résultats :  $\hat{\beta} = \frac{1}{M} \sum_{m=1}^M \hat{\beta}_m$

$$T = \frac{1}{M} \sum_m \widehat{Var}(\hat{\beta}_m) + \left(1 + \frac{1}{M}\right) \frac{1}{M-1} \sum_m (\hat{\beta}_m - \hat{\beta})^2$$

$\Rightarrow$  Objectif : fournir une estimation des paramètres et de leur variabilité (prendre en compte la variabilité due aux données manquantes)



## Imputation multiple : principe

- 1 Créer des jeux de données bootstrap (autre possibilité régression Bayésienne)
- 2 Estimer sur chaque jeu de données les paramètres du modèle :  $(\hat{\beta})^1, \dots, (\hat{\beta})^M \implies$  variabilité sur le modèle
- 3 Ajouter du bruit en imputant pour  $m = 1, \dots, M$  valeurs manquantes  $y_i^m$  en tirant dans la distribution prédictive  $\mathcal{N}(x_i \hat{\beta}^m, (\hat{\sigma}^2)^m)$

2 sources de variabilité : dans les paramètres du modèle & dans le bruit ajouté

Variance de prédiction = variance d'estimation + bruit

## Modèle joint

⇒ Hypothèse  $\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

Algorithme :

- ① Bootstrap des lignes :  $\mathbf{X}^1, \dots, \mathbf{X}^M$   
Algorithme EM :  $(\hat{\boldsymbol{\mu}}^1, \hat{\boldsymbol{\Sigma}}^1), \dots, (\hat{\boldsymbol{\mu}}^M, \hat{\boldsymbol{\Sigma}}^M)$
- ② Imputation :  $x_{ij}^m$  tirée depuis  $\mathcal{N}(\hat{\boldsymbol{\mu}}^m, \hat{\boldsymbol{\Sigma}}^m)$

Facile à paralléliser

Implémenté dans **Amelia** ([website](#))



Amelia Earhart



James Honaker



Gary King



Matt Blackwell

## Modèle conditionnel

⇒ Hypothèse : un modèle par variable

Algorithme :

- ① Imputation initiale : imputation par la moyenne
- ② Pour la variable  $j$ 
  - 2.1  $(\beta^{-j}, \sigma^{-j})$  tirés d'une distribution Bootstrap ou a posteriori
  - 2.2 Imputation : régression aléatoire  $x_{ij}$  tiré dans  $\mathcal{N}(\mathbf{X}_{-j}\beta^{-j}, \sigma^{-j})$
- ③ Boucler sur les variables
- ④ Répéter  $M$  fois les étapes 2 et 3

Implémenté dans `mice` ([website](#))

*“There is no clear-cut method for determining whether the MICE algorithm has converged”*



Stef van Buuren

## Modèle joint versus modèle conditionnel

⇒ Modèle conditionnel prend le leadership ?

- Flexible : un modèle par variable. Facile de gérer les interactions et les variables de natures différentes (binaire, ordinale, quali...)
- Beaucoup de modèles statistiques sont des modèles conditionnels !
- Fonctionne bien en pratique

⇒ Inconvénients : 1 modèle/variable... fastidieux...

## Modèle joint versus modèle conditionnel

⇒ Modèle conditionnel prend le leadership ?

- Flexible : un modèle par variable. Facile de gérer les interactions et les variables de natures différentes (binaire, ordinale, quali...)
- Beaucoup de modèles statistiques sont des modèles conditionnels !
- Fonctionne bien en pratique

⇒ Inconvénients : 1 modèle/variable... fastidieux...

⇒ Que faire avec fortes corrélations ou quand  $n < p$  ?

- modèle joint régularise la covariance  $\Sigma + k\mathbb{I}$  (choix de  $k$  ?)
- modèle conditionnel : régression ridge ou sélection de variables  
⇒ beaucoup de paramètres de réglage ... pas facile ...

## Imputation multiple avec ACP et Bootstrap

$$\begin{aligned}x_{ij} &= \tilde{x}_{ij} + \varepsilon_{ij}, \varepsilon_{ij} \sim \mathcal{N}(0, \sigma^2) \\ &= \sum_{s=1}^S d_s u_{is} v_{js} + \varepsilon_{ij}\end{aligned}$$

- 1 Variabilité des paramètres,  $M$  jeux possibles :  $(\hat{x}_{ij})^1, \dots, (\hat{x}_{ij})^M$   
Bootstrap des résidus :  $\mathbf{X}^1 = \hat{\mathbf{X}} + \varepsilon^1, \dots, \mathbf{X}^M = \hat{\mathbf{X}} + \varepsilon^M$   
ACP itérative :  $\hat{\mathbf{X}}^1 = \mathbf{U}^1 \mathbf{D}^1 \mathbf{V}^1, \dots, \hat{\mathbf{X}}^M = \mathbf{U}^M \mathbf{D}^M \mathbf{V}^M$
- 2 Bruit : pour  $m = 1, \dots, M$ , valeurs manquantes  $x_{ij}^m$  sont imputées en choisissant depuis une distribution prédictive  $\mathcal{N}(\hat{x}_{ij}^m, \hat{\sigma}^2)$

Implémenté dans `missMDA` ([website](#))

# Imputation multiple en pratique

⇒ Etape 1 : Générer  $M$  jeux de données imputés

```
> library(Amelia)
> res.amelia <- amelia(don,m=100)  ## avec package zelig

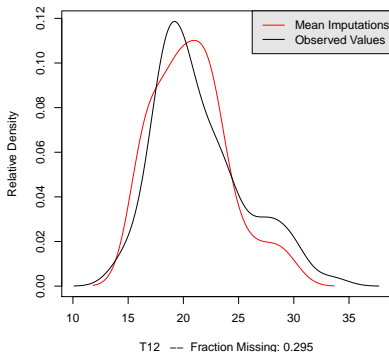
> library(mice)
> res.mice <- mice(don,m=100,defaultMethod="norm.boot")

> library(missMDA)
> res.MIPCA <- MIPCA(don,ncp=2,nboot=100)
> res.MIPCA$resMI
```

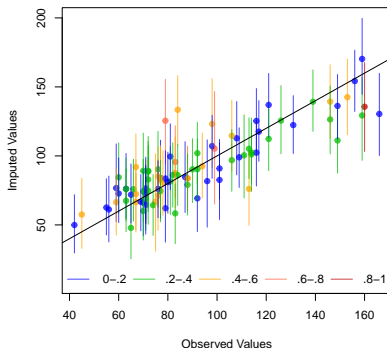
# Imputation multiple en pratique

## Etape 2 : visualisation

Observed and Imputed values of T12



Observed versus Imputed Values of maxO3



```
> library(Amelia)
> res.amelia <- amelia(don,m=100)
> compare.density(res.amelia, var="T12")
> overimpute(res.amelia, var="maxO3")
```

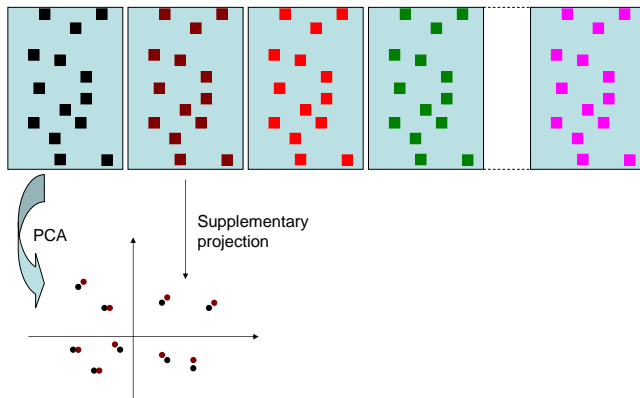
fonction stripplot dans mice



## Imputation multiple en pratique

Etape 2 : visualisation de l'incertitude liée aux NA

Quelle confiance accorder aux représentations ? Notion de variance

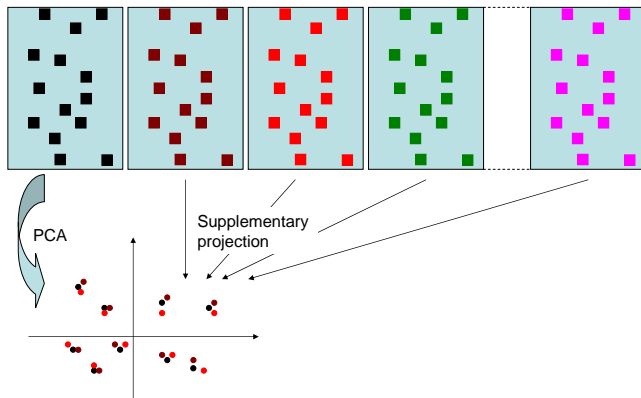


ACP itérative régularisée  
⇒ configuration de référence

## Imputation multiple en pratique

Etape 2 : visualisation de l'incertitude liée aux NA

Quelle confiance accorder aux représentations ? Notion de variance

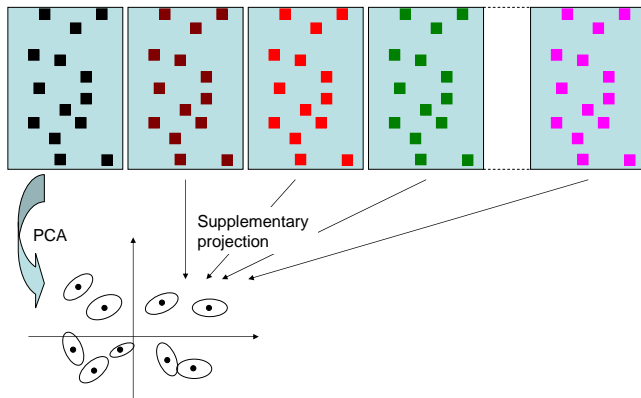


ACP itérative régularisée  
⇒ configuration de référence

## Imputation multiple en pratique

Etape 2 : visualisation de l'incertitude liée aux NA

Quelle confiance accorder aux représentations ? Notion de variance



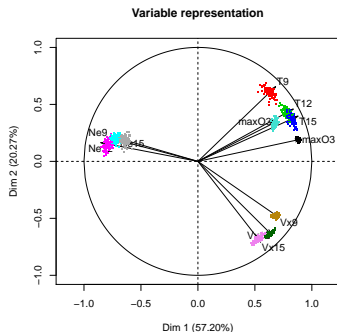
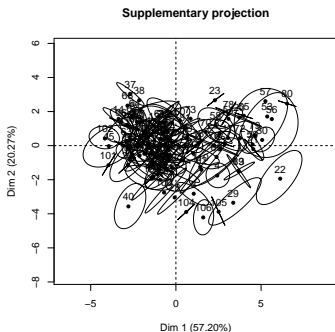
ACP itérative régularisée  
⇒ configuration de référence

# Imputation multiple en pratique

⇒ Etape 2 : visualisation de l'incertitude liée aux NA

```
> res.MIPCA <- MIPCA(don,ncp=2)
```

```
> plot(res.MIPCA,choice= "ind.supp"); plot(res.MIPCA,choice= "var ")
```



## Imputation multiple en pratique

⇒ Etape 3. Régression par tableau et combinaison des résultats

$$\hat{\beta} = \frac{1}{M} \sum_{m=1}^M \hat{\beta}_m$$

$$T = \frac{1}{M} \sum_m \widehat{Var}(\hat{\beta}_m) + \left(1 + \frac{1}{M}\right) \frac{1}{M-1} \sum_m (\hat{\beta}_m - \hat{\beta})^2$$

```
> require(mice)
> imp<-prelim(res.mi=res.MIPCA,X=ozone[,1:11])
> fit <- with(data=imp,exp=lm(maxO3~T9+T12+T15+Ne9+Ne12+...+Vx15+maxO3v))
> res.pool<-pool(fit)
> summary(res.pool)
```

	est	se	t	df	Pr(> t )	lo 95	hi 95	nmis	fmi	lambda
(Intercept)	19.31	16.30	1.18	50.48	0.24	-13.43	52.05	NA	0.46	0.44
T9	-0.88	2.25	-0.39	26.43	0.70	-5.50	3.75	37	0.71	0.69
T12	3.29	2.38	1.38	27.54	0.18	-1.59	8.18	33	0.70	0.68
....										
Vx15	0.23	1.33	0.17	39.00	0.87	-2.47	2.93	21	0.57	0.55
maxO3v	0.36	0.10	3.65	46.03	0.00	0.16	0.56	12	0.50	0.48

# Imputation multiple pour variables qualitatives ou mixtes

⇒ Modèle joint :

- Modèle log-linear (Schafer, 1997) (**cat**) : pb si bcp de modalités
- Modèles à classes latentes (Vermunt, 2014) - Bayésien non-paramétrique (Si & Reiter, 2014, Murray & Reiter, 2016)  
(**MixedDataImpute**, **NPBayesImpute**, **NestedCategBayesImpute**)

⇒ Modèle conditionnel : logistique, multinomial, forêts (**mice**)

⇒ MIMCA ((**MIMCA**) de **missMDA**) fournit des inférences valides (ex. régression logistique avec NA) pour des jeux de données avec bcp de modalités et/ou modalités rares

**Imputation multiple pour données mixtes** : **MIFAMD** sur le même principe, et modèles joints et conditionnels

# Plan

- 1 Introduction
- 2 Imputation simple pour variables quantitatives
- 3 Imputation simple pour variables qualitatives
- 4 Imputation simple pour données mixtes
- 5 Imputation multiple
- 6 Conclusion**

## Remarques

- ***“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 and Rubin, 1983)
- Théorie de l'IM : bonne pour la régression. Autres méthodes ?
- Modèle d'imputation doit être aussi complexe que le modèle d'analyse (interaction)



## Remarques

- ***"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 and Rubin, 1983)
- Théorie de l'IM : bonne pour la régression. Autres méthodes ?
- Modèle d'imputation doit être aussi complexe que le modèle d'analyse (interaction)
- Quelques problèmes pratiques :
  - Imputation de  $X$  et  $X^2$
  - Problèmes de bornes ( $> 0$ )  $\Rightarrow$  tronquer ?
  - Comment faire avec des données de grandes dimensions ?

# Une page Web et des didacticiels

[http://factominer.free.fr/missMDA/index\\_fr.html](http://factominer.free.fr/missMDA/index_fr.html)



## Le package missMDA

Le package **missMDA** est complémentaire de FactoMineR. Il permet de gérer les données manquantes pour les méthodes d'analyses factorielles (ACP, AFC, ACM, AFDM, AFM). Il permet de faire de l'imputation simple et multiple.

L'imputation simple consiste à remplacer les valeurs manquantes par des valeurs plausibles. Cela revient à compléter le jeu de données qui peut ensuite être analysé par n'importe quelle méthode d'analyse factorielle.

**missMDA** impute les valeurs manquantes de sorte que les valeurs imputées n'ont aucune influence sur les résultats de l'analyse factorielle (pas d'influence dans le sens où les valeurs imputées n'ont aucun poids, et donc les résultats de l'analyse factorielle sont obtenues uniquement avec les valeurs observées).

**missMDA** utilise des méthodes de réduction de données, ce qui lui permet d'imputer de façon satisfaisante de gros jeux de données contenant des variables quantitatives et/ou qualitatives. En effet, il impute par ACP (ou ACM, ou AFDM ou AFM) en prenant en compte à la fois les similarités entre individus et les liens entre variables.

Voir cette vidéo si vous voulez comprendre le principe de missMDA quelque soit les jeux de données (quantitatifs et/ou qualitatifs).

Les imputations sont très bonnes comparées aux méthodes classiques permettant d'imputer des tableaux incomplets (forêts aléatoires par exemple).

- **missMDA** gère les données manquantes dans:
  - les jeux de données avec variables quantitatives grâce à l'ACP (Voir la vidéo)
  - les jeux de données avec variables qualitatives grâce à l'ACM (Voir la vidéo)
  - les tableaux de contingence grâce à l'AFC
  - les données mixtes grâce à l'AFDM
  - les jeux de données où les variables sont structurées par groupe grâce à l'AFM
- **missMDA** permet de faire de l'imputation multiple:
  - pour les variables quantitatives grâce à l'ACP: Voir la vidéo
  - pour les variables qualitatives grâce à l'ACM

## Menu sur les données manquantes

### Le package missMDA

ACP avec données manquantes

ACM avec données manquantes

Imputation multiple

Peut-on croire dans les valeurs imputées ?

Références - Conférences

## Les auteurs de missMDA

François Husson

Julie Josse

## Ressources

⇒ Logiciels :

- R CRAN task View: Missing Data
- R-miss-tastic

⇒ Articles :

- Imbert, A., & Vialaneix, N. (2018). Décrire, prendre en compte, imputer et évaluer les valeurs manquantes dans les études statistiques : une revue des approches existantes. *Journal de la SFdS*, **159(2)**, 1-55.
- Josse J, Husson F. & Pagès J (2009) Gestion des données manquantes en Analyse en Composantes Principales. *Journal de la SFdS*. **150 (2)**, 28-51.