

Imputation multiple à l'aide des méthodes d'analyse factorielle

Vincent Audigier & Julie Josse & François Husson

Agrocampus Rennes

45e journées de Statistique, Toulouse, 28 Mai 2013

Données manquantes

	X				Y		
NA					NA	NA	
			NA				
	NA						NA
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
		NA	NA	NA			

Données manquantes

					X		Y	
NA					NA	NA		
			NA					
	NA						NA	
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	
		NA	NA	NA				

- Suppression des individus : cas complet

Données manquantes

X					Y		
NA					NA	NA	
			NA				
	NA						NA
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
		NA	NA	NA			

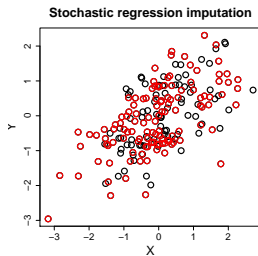
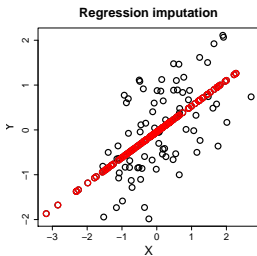
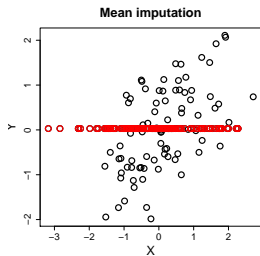
- Suppression des individus : cas complet
- Adaptation de la méthode au cas incomplet

Données manquantes

	X				Y		
NA					NA	NA	
			NA				
	NA						NA
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
		NA	NA	NA			

- Suppression des individus : cas complet
- Adaptation de la méthode au cas incomplet
- Imputation

Méthodes d'imputation simple



$\mu_y = 0$
 $\sigma_y^2 = 1$
 $\rho = 0.6$
 $IC_{\mu_y, 95\%}$

0.01
0.5
0.30
39.4

0.01
0.72
0.78
61.6

0.01
0.99
0.59
70.8

Principe

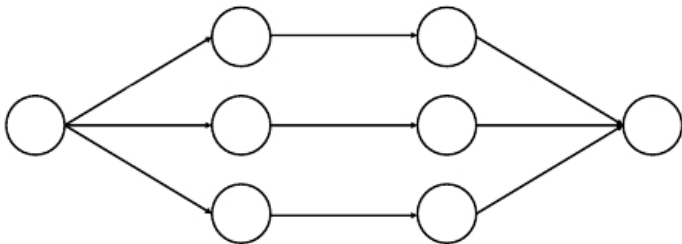
IM comme solution au problème de sous-estimation de la variance

Tableau incomplet

Imputation

Analyse

Agrégation



Principe

Estimation du paramètre d'intérêt

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

Décomposition de la variance

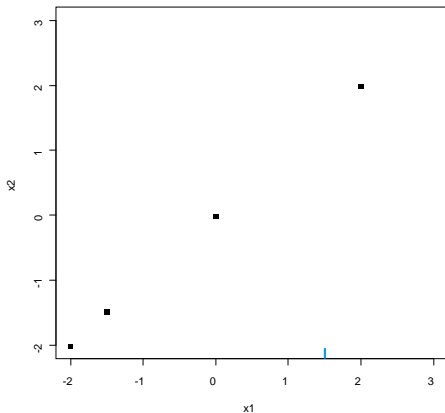
$$T = \frac{1}{M} \sum_{m=1}^M \widehat{Var}(\hat{\theta}_m) + \left(1 + \frac{1}{M}\right) \frac{1}{M-1} \sum_{m=1}^M (\hat{\theta}_m - \hat{\theta})^2$$

Une méthode d'imputation multiple repose sur une
méthode d'imputation simple

Algorithme : ACP itérative

Le jeu de données

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

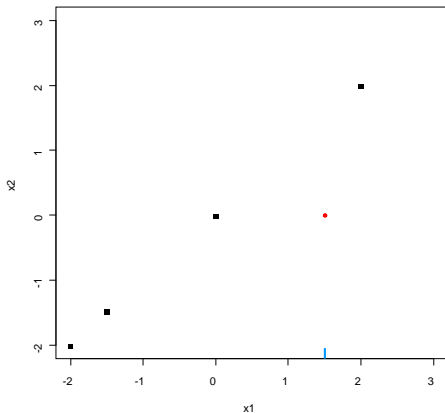


Algorithme : ACP itérative

Initialisation : imputation par la moyenne

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



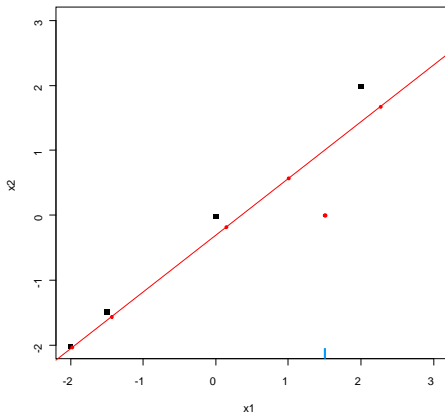
Algorithme : ACP itérative

ACP sur le jeu complété à partir d'une dimension

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



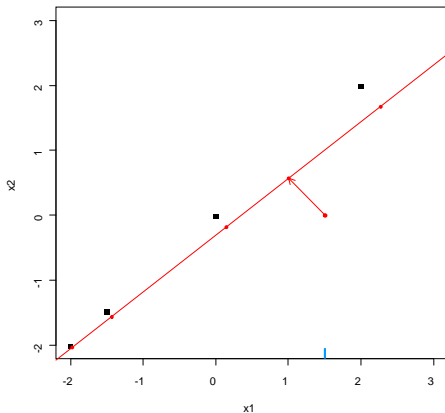
Algorithme : ACP itérative

Détermination de la valeur prédite par le modèle

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



Algorithme : ACP itérative

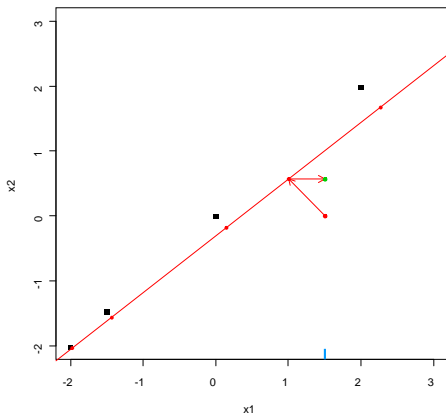
Etape d'imputation : $X^\ell = W * X + (1 - W) * \hat{X}^\ell$

```
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
```

```
  x1  x2
  x1  x2
-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
```



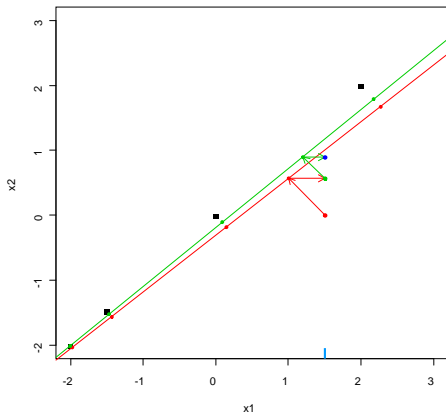
Algorithme : ACP itérative

ACP sur le jeu complété à partir d'une dimension

```
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
```



Algorithme : ACP itérative

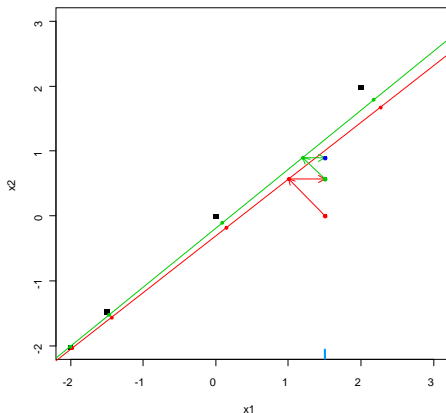
Etape d'imputation : $X^\ell = W * X + (1 - W) * \hat{X}^\ell$

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



Algorithme : ACP itérative

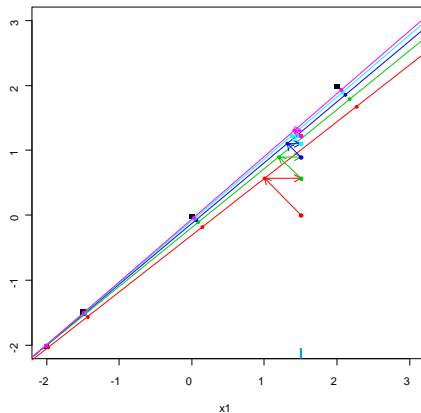
Processus itéré jusqu'à convergence

```
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
```

```
^      ^
x1     x2
-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
```



Algorithme : ACP itérative

- 1 initialisation $\ell = 0$: X^0 (imputation par la moyenne)
- 2 itération ℓ :
 - (a) recherche de $\hat{F}_{I \times S}^\ell$ et $\hat{U}_{K \times S}^\ell$
 - (b) imputation par $\hat{x}_{ik}^\ell = \sum_{s=1}^S \hat{f}_{is}^\ell \hat{u}_{ks}^\ell = \sum_{s=1}^S \sqrt{\hat{\lambda}_s} \hat{v}_{is}^\ell \hat{u}_{ks}^\ell$
- 3 (a) et (b) sont alternées jusqu'à convergence

Algorithme : ACP itérative régularisée

- 1 initialisation $\ell = 0$: X^0 (imputation par la moyenne)
- 2 itération ℓ :
 - (a) recherche de $\hat{F}_{I \times S}^\ell$ et $\hat{U}_{K \times S}^\ell$
 - (b) imputation par $\hat{x}_{ik}^\ell = \sum_{s=1}^S \left(\sqrt{\hat{\lambda}_s} - \frac{\hat{\sigma}^2}{\hat{\lambda}_s} \right) \hat{v}_{is}^\ell \hat{u}_{ks}^\ell$
- 3 (a) et (b) sont alternées jusqu'à convergence

Algorithme MIPCA

1 Initialisation

- Estimation des paramètres du modèle d'ACP

$$X_{I \times K} = F_{I \times S} U_{K \times S}^T + E_{I \times K} \text{ avec } E = (\epsilon_{ik})_{ik} \text{ et}$$
$$\epsilon_{ik} \sim \mathcal{N}(0, \sigma^2) \Rightarrow \hat{F}, \hat{U}, \hat{\sigma}^2$$

Algorithme MIPCA

1 Initialisation

- Estimation des paramètres du modèle d'ACP

$$X_{I \times K} = F_{I \times S} U_{K \times S}^\top + E_{I \times K} \text{ avec } E = (\epsilon_{ik})_{ik} \text{ et} \\ \epsilon_{ik} \sim \mathcal{N}(0, \sigma^2) \Rightarrow \hat{F}, \hat{U}, \hat{\sigma}^2$$

2 Répéter M fois :

- Propagation de l'incertitude sur \hat{F} et \hat{U}
 - Création du tableau $X_m^* = \hat{F} \hat{U}^\top + E_m^*$ avec $E_m^* \sim \mathcal{N}(0, \hat{\sigma}^2)$
 - ACP sur $X_m^* \Rightarrow \hat{F}_m^*, \hat{U}_m^*, \hat{\sigma}_m^{*2}$

Algorithme MIPCA

1 Initialisation

- Estimation des paramètres du modèle d'ACP

$$X_{I \times K} = F_{I \times S} U_{K \times S}^\top + E_{I \times K} \text{ avec } E = (\epsilon_{ik})_{ik} \text{ et} \\ \epsilon_{ik} \sim \mathcal{N}(0, \sigma^2) \Rightarrow \hat{F}, \hat{U}, \hat{\sigma}^2$$

2 Répéter M fois :

- Propagation de l'incertitude sur \hat{F} et \hat{U}
 - Création du tableau $X_m^* = \hat{F} \hat{U}^\top + E_m^*$ avec $E_m^* \sim \mathcal{N}(0, \hat{\sigma}^2)$
 - ACP sur $X_m^* \Rightarrow \hat{F}_m^*, \hat{U}_m^*, \hat{\sigma}_m^{*2}$
- Respecter la distribution du jeu de données
 - Imputation simple : $X_m \leftarrow W * X + (1 - W) * \hat{X}_m^*$
 - Ajout du résidu : $X_m \leftarrow X_m + \tilde{E}$ tiré selon $\mathcal{N}(0, \hat{\sigma}^{*2})$

Algorithme Expectation Maximization Bootstrap

- Hypothèse : $X_{I \times K} \sim \mathcal{N}(\mu, \Sigma)$
- Algorithme :
 - 1 Bootstrap des individus
 - 2 Estimation de μ et Σ
 - 3 Imputation par régression aléatoire
- Implémenté dans le package Amelia du logiciel R

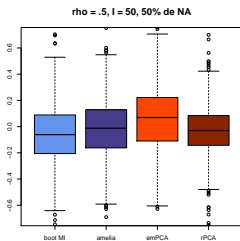
Imputation par équations enchaînées (MICE)

- Hypothèse sur les distributions conditionnelles
- Algorithme :
 - ① Initialisation :
 - Spécifier un modèle pour chaque variable
 - Imputation des données manquantes par tirage dans les données observées
 - ② Pour chaque variable
 - Estimation des paramètres du modèle
 - Imputation par régression aléatoire
 - ③ Répéter (2) un nombre de fois fixé à l'avance
- Implémenté dans le package MICE du logiciel R

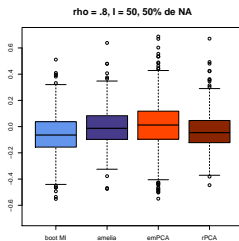
Protocole de simulation

- Paramètres d'intérêt : $\theta_1 = \beta_1$, $\theta_2 = \mathbb{E}[Y]$
- 1000 simulations
 - Simulation de X selon une loi normale pour 3 matrices de covariance différentes
 - Génération 50% de NA complètement au hasard
 - Imputation de $M = 20$ tableaux par EMB, MICE, MIPCA
- Evaluation par
 - le biais sur θ_i
 - le coverage

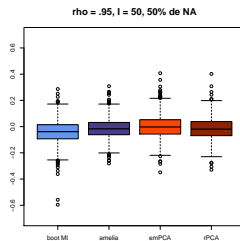
Résultats $\theta_1 = \beta_1$



collinearite faible

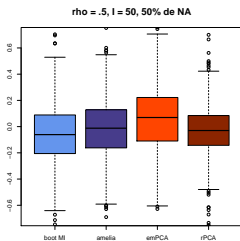


collinearite moyenne

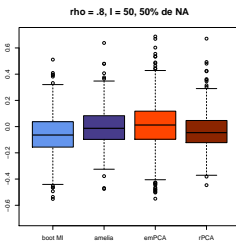


collinearite forte

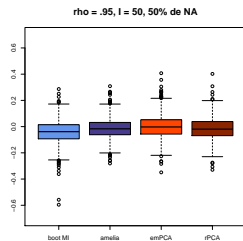
Résultats $\theta_1 = \beta_1$



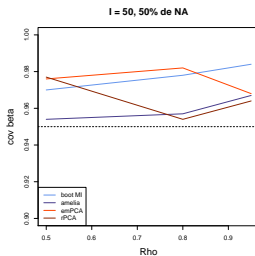
collinearite faible



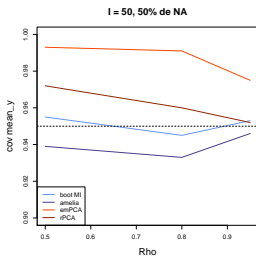
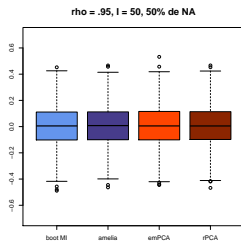
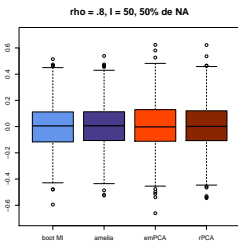
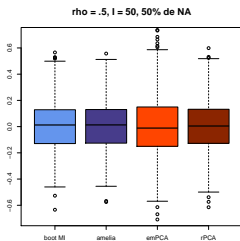
collinearite moyenne



collinearite forte



Résultats $\theta_2 = \mathbb{E}[Y]$



Conclusion

- Une méthode nouvelle pour effectuer de l'imputation multiple qui
- permet d'imputer des jeux de données quantitatifs, y compris quand $I \ll K$.
 - nécessite le réglage d'un paramètre (nombre de dimensions)
 - est disponible dans le package `missMDA` du logiciel R
 - peut s'étendre aux données mixtes