

# Gestion des données manquantes en clustering

Présenté par :

Fadela SADOU ZOULEYA  
Mouhamadou Lamine NDAO

Encadré par :

Ndèye Niang  
Vincent Audigier

29 Octobre 2021

# Clustering

**Données**  $\mathbf{Z} = (z_{ij})$   $1 \leq i \leq n$  ensemble de données numériques  
 $1 \leq j \leq p$

**P** partition en  $K$  classes où chaque  $i$  appartient à une seule classe  
 $C_i \in \{1, \dots, K\}$ . Le centroïde de la classe  $k$  est noté  $c_k \in \mathcal{M}_{1 \times p}$

**Objectif** Identifier  $C_i$  pour chaque individu  $i$  à partir de  $(\mathbf{z}_i)_{(1 \leq i \leq n)}$

## Approches existantes

basées sur distances

- k-means
- fuzzy C-means
- hierarchical clustering
- pam

basées sur modèles

- gaussian mixture models
- mixture of multivariate  $t$ -distributions

# Gestion des données manquantes en clustering

**Contraintes**,  $\mathbf{Z}$  souvent **incomplet**...  $\mathbf{z}_i = (\mathbf{z}_i^{obs}, \mathbf{z}_i^{miss})$

## Méthodes "Ad-hoc"

- Analyse des cas complets
- Imputation simple

## Méthodes avancées

- Imputation Multiple (IM)
- Méthodes directes (fuzzy C-means, k-Pod, ...)

## ① Introduction

## ② Imputation multiple

Principe général : exemple de la régression

Imputation multiple en clustering

Simulations

## ③ Méthodes directes

Méthodes

k-POD (Chi et al., 2016)

Ignorable-GMM (Marbac et al., 2019)

Optimal Completion Strategy of fuzzy c-means (Hathaway et Bezdek, 2001)

Simulations

## ④ Conclusion

## ① Introduction

## ② Imputation multiple

Principe général : exemple de la régression

Imputation multiple en clustering

Simulations

## ③ Méthodes directes

Méthodes

k-POD (Chi et al., 2016)

Ignorable-GMM (Marbac et al., 2019)

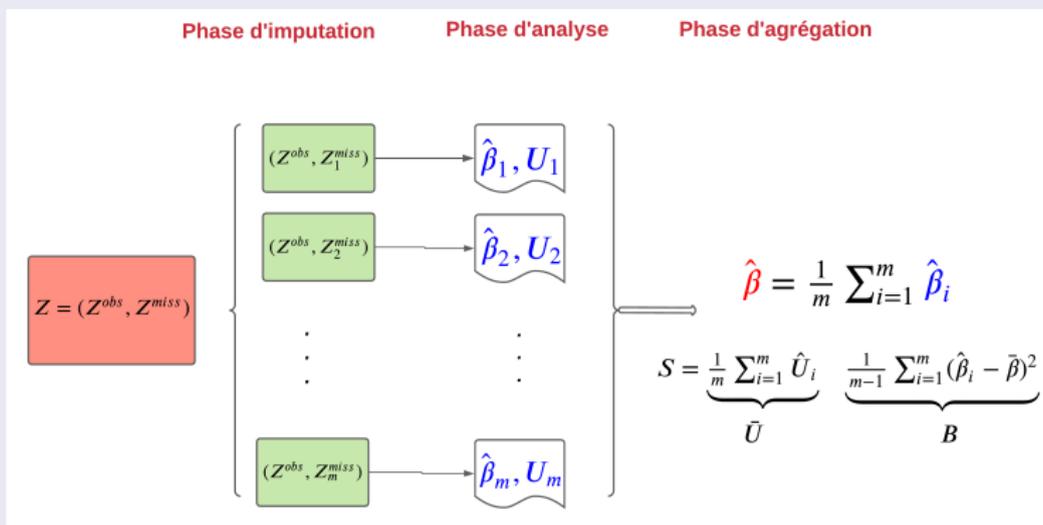
Optimal Completion Strategy of fuzzy c-means (Hathaway et Bezdek, 2001)

Simulations

## ④ Conclusion

# Imputation multiple (Rubin, 1987) : cas de la régression

## Principe



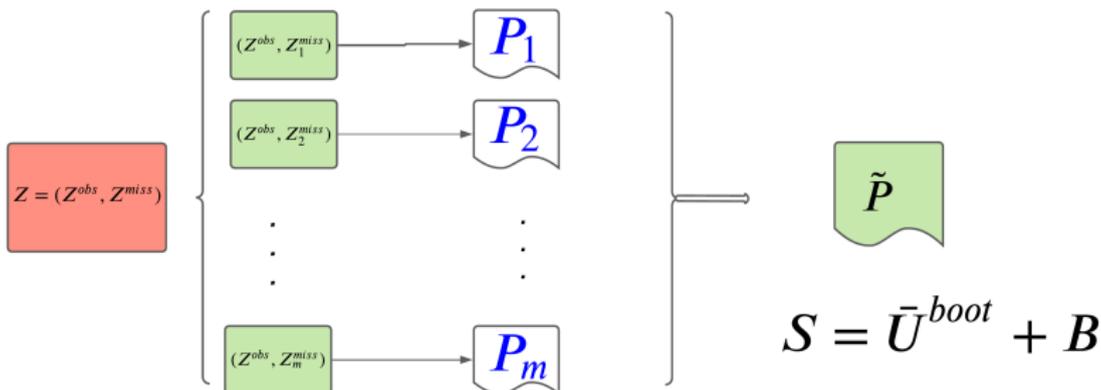
- Estimation unique de  $\beta$  et sa variance :  $S = \bar{U} + B$
- Qualité de  $\hat{\beta}$  dépend de l'indépendance des tableaux imputés conditionnellement aux valeurs observées

# IM pour le clustering (Audigier and Niang, 2020)

**Phase d'imputation**  
Prise en compte de la structure  
en groupes (JM-DP)

**Phase d'analyse**  
Appliquer un  
clustering (kmeans)

**Phase d'agrégation**  
Consensus de partitions  
(extension des règles de  
Rubin)



# Phase d'imputation : JM-DP (Kim et al., 2014)

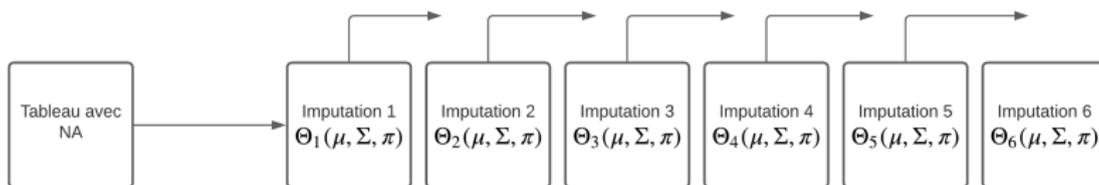
## Approche d'imputation

On suppose que  $Z_i | C_i, \mu, \Sigma \sim N(Z_i | \mu_{C_i}, \Sigma_{C_i})$

### Principe :

- 1 Alterner le tirage des paramètres dans leur distribution a posteriori et l'imputation des NA selon les paramètres simulés.
- 2 Répéter jusqu'à convergence vers la distribution a posteriori des paramètres  $\Theta(\mu, \Sigma, \pi)$ .

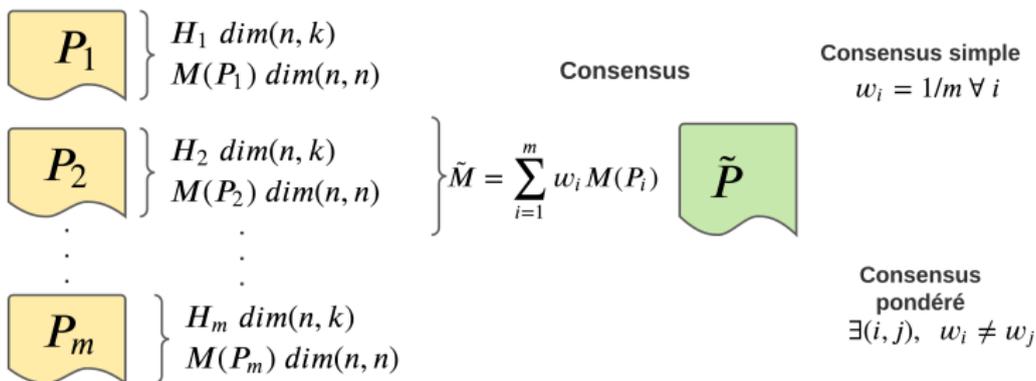
**Imputation multiple** : Répéter à nouveau (1) et (2) de façon à récupérer  $m$  réalisations des paramètres issus de la loi a posteriori.



# Phase d'agrégation : consensus de partitions

## Consensus de partitions

- **Objectif** : Identifier une partition compromis d'un ensemble de partitions obtenues sur le même ensemble d'observations.
- **Principe** :



# Approche simple : NMF (Tao Li and al., 2007)

- Formulation :

$$\min_H \|\tilde{M} - HH'\|^2 \quad (1)$$

- $H = \{0, 1\}^{(nk)}$  : matrice d'indicatrices  
 $M = \{0, 1\}^{(nn)}$  : matrice de connectivité.
- En posant  $\tilde{H} = H(H'H)^{-1/2}$  et  $D = \text{diag}(H'H)$  (1) devient :

$$\min_{\tilde{H}'\tilde{H}=I, \tilde{H} \geq 0} \|\tilde{M} - \tilde{H}D\tilde{H}'\|^2 \quad (2)$$

$$\tilde{H} \leftarrow \tilde{H} \sqrt{\frac{(\tilde{M}'\tilde{H}D)}{(\tilde{H}\tilde{H}'\tilde{M}'\tilde{H}D)}} \quad (3)$$

$$D \leftarrow D \sqrt{\frac{(\tilde{H}'\tilde{M}\tilde{H})}{(\tilde{H}'\tilde{H}D\tilde{H}'\tilde{H})}} \quad (4)$$

## Approche pondérée : Weighted-NMF (Tao Li and al., 2008)

- Introduire  $w$  tel que :

$$w = (w_1, w_2, \dots, w_m)', \quad w_t \geq 0, \quad \|w\|_1 = \sum_{t=1}^m w_t = 1 \quad (5)$$

- On définit  $\tilde{M} = \sum_{t=1}^m w_t M(P_t)$  et le problème (1) devient :

$$\min_{w, \tilde{H}} \|\tilde{M} - \tilde{H}\tilde{H}'\|^2 \quad (6)$$

- 1 Chercher  $\tilde{H}$  en fixant  $w$  en utilisant NMF
- 2 Chercher  $w$  en fixant  $\tilde{H}$  en résolvant  $\min w'Aw - 2b' + cte$

# Problématique

- Les travaux sur l'agrégation des partitions après IM n'effectue que du consensus simple
- Problème : par construction, les tableaux imputés ne sont jamais indépendants.
- Idée : utiliser le consensus pondéré pour corriger les redondances des partitions.

# Simulations

# Données

**30 expériences** ont été simulées. Pour une expérience, nous avons :

## Données

- $n = 100$  observations,  $p = 10$  variables
- $K = 2$  clusters :  $n_w = 50$ ,  $w \in \{1, 2\}$
- moyennes  $\mu_w$ ,  $w \in \{1, 2\}$  :  $\mu_1 = (0, 0, 0, 0, 0, 0, 0, 0, 0, 0)$   
 $\mu_2 = (0, 0, 0, 0, 0, 2, 2, 2, 2, 2)$
- Matrice variance covariance  $\Sigma = \begin{pmatrix} I_5 & \mathbf{0} \\ \mathbf{0} & \begin{matrix} 1 & 0.3 & 0.3 & 0.3 & 0.3 \\ 0.3 & 1 & 0.3 & 0.3 & 0.3 \\ 0.3 & 0.3 & 1 & 0.3 & 0.3 \\ 0.3 & 0.3 & 0.3 & 1 & 0.3 \\ 0.3 & 0.3 & 0.3 & 0.3 & 1 \end{matrix} \end{pmatrix}$
- **30%** de données manquantes selon un mécanisme MAR.

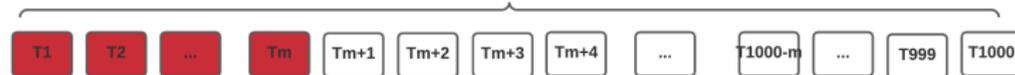
## Plan

## Plan d'expérimentation

 $m \in \{3, 5, 20\}$  : nombre de tableaux


Imputation  
Multiple

500 tableaux imputés

 $\delta_i = 1$ 


Analyse : m partitions  
(kmeans)

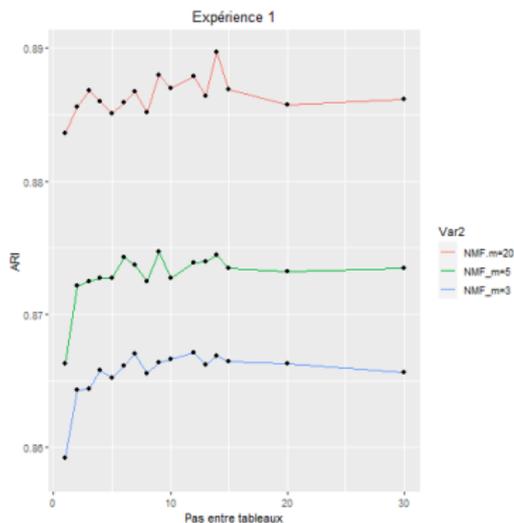
Agrégation: consensus des  
m partitions  
Calcul ARI\_1

 $\delta_i = 2$ 

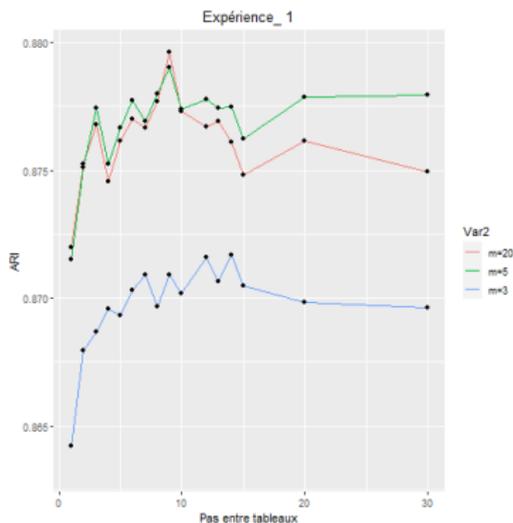

Analyse : m partitions (kmeans)  
Agrégation: consensus des m partitions  
Calcul ARI\_2

## Résultats

## Résultats NMF



## Résultats WNMF

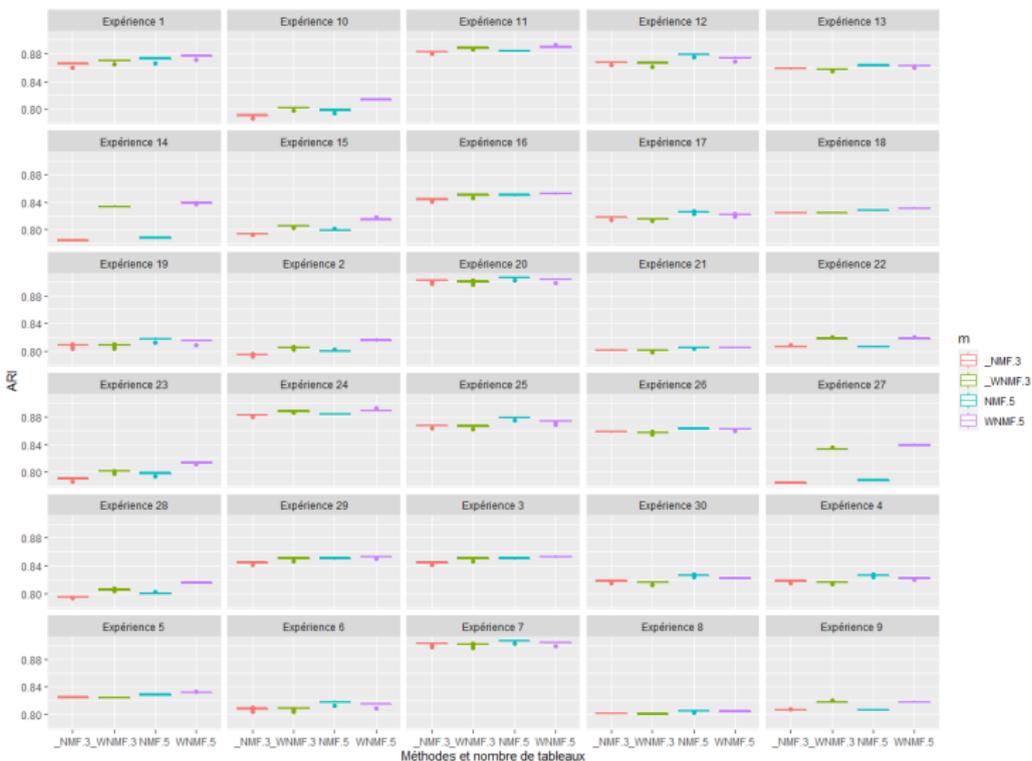


## Résultats

- Moins sensible à l'indépendance entre tableaux.
- Le nombre de tableaux  $\tilde{P}$  considéré améliore la qualité  $\tilde{P}$ .

# Résultats

## Comparaison des approches NMF et WNMF



# Synthèse

- Clustering après IM est beaucoup moins sensible à **l'indépendance**
- En clustering après IM, en prenant un nombre relativement élevé de tableaux, il n'est pas nécessaire de s'assurer de l'indépendance entre les tableaux.
- Pas de **différences significatives**, en termes de performances, entre les approches NMF et WNMF après IM.

## Perspectives

- L'apport du consensus pondéré dans une stratégie de stacking.
- Évaluation de la qualité d'une partition consensus après IM

## ① Introduction

## ② Imputation multiple

Principe général : exemple de la régression

Imputation multiple en clustering

Simulations

## ③ Méthodes directes

Méthodes

k-POD (Chi et al., 2016)

Ignorable-GMM (Marbac et al., 2019)

Optimal Completion Strategy of fuzzy c-means (Hathaway et Bezdek, 2001)

Simulations

## ④ Conclusion

# K-POD (Chi et al., 2016)

## k-means

$$\min_{c_1, \dots, c_K, C_1, \dots, C_K} \sum_{k=1}^K \sum_{i \in C_k} d(Z_i, c_k)^2$$

$$\iff$$

$$\min_{H, C} \|Z - HC\|_F^2$$

## k-POD

$$\min_{c_1, \dots, c_K, C_1, \dots, C_K} \sum_{k=1}^K \sum_{ij, i \in C_k \text{ et } ij \in \Omega} (z_{ij} - c_{kj})^2$$

$$\iff$$

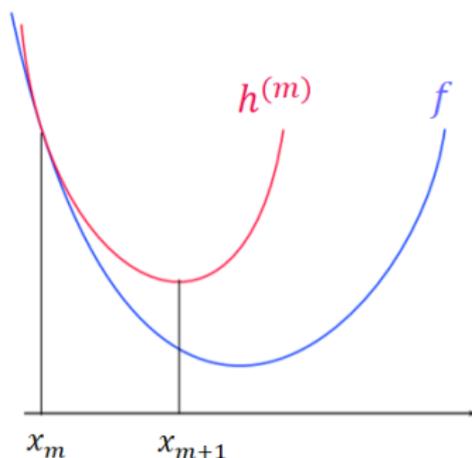
$$\min_{H, C} \|P_\Omega(Z) - P_\Omega(HC)\|_F^2$$

**Pas de solution explicite**

- $H = [h_{ik}] \in M_{n \times K} \rightarrow h_{ik} = \begin{cases} 1 & \text{si } i \in C_k \\ 0 & \text{sinon} \end{cases}$
- $\Omega \subset \{1, \dots, n\} \times \{1, \dots, p\} \rightarrow$  sous-ensemble des indices des valeurs observées
- $P_\Omega \rightarrow [P_\Omega(Z)]_{ij} = \begin{cases} z_{ij} & \text{si } (i, j) \in \Omega \\ 0 & \text{si } (i, j) \in \Omega^c \end{cases}$

## K-POD (Chi et al., 2016)

## Algorithme de majoration-minimisation



La fonction de perte  $f : x \rightarrow f(x)$  est majorée à l'itération  $m$  par la fonction  $h^{(m)}$ .

## K-POD (Chi et al., 2016)

**Fonction de perte  $f(H, C)$**

$$\| P_{\Omega}(Z) - P_{\Omega}(HC) \|_F^2$$

**Fonction majorante  $g^{(q)}$  en  $(H^{(q)}, C^{(q)})$**

$$\| P_{\Omega}(Z) - P_{\Omega}(HC) \|_F^2 + \| P_{\Omega^c}(Z) - P_{\Omega^c}(H^{(q)}C^{(q)}) \|_F^2$$

=

$$\| Z^{(q)} - HC \|_F^2$$

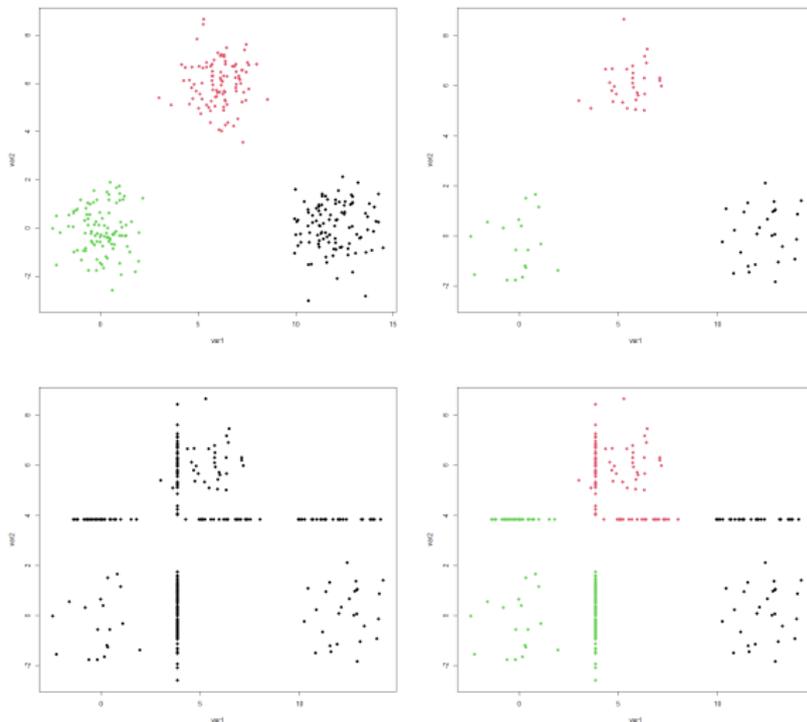
# K-POD (Chi et al., 2016)

Pour le nombre  $K$  de clusters fixé :

- **Initialiser** les valeurs manquantes et obtenir la base complétée  $Z^{(0)}$ , puis calculer la **partition initiale**  $(H^{(0)}, C^{(0)})$  par **k-means** ;
- **Répéter** pour les itérations  $q = 0, 1, 2, \dots$  jusqu'à convergence :
  - **remplacer** les valeurs manquantes par celles des **centroïdes** de leur groupe respectif pour obtenir  $Z^{(q)}$  ;
  - estimer le minimum  $(H^{(q+1)}, C^{(q+1)})$  de la fonction majorante  $g^{(q)}(H, C | H^{(q)}, C^{(q)})$  par **k-means** sur  $Z^{(q)}$ .

# K-POD (Chi et al., 2016)

## Initialisation des valeurs manquantes



# Ignorable-GMM (Marbac et al., 2019)

**Notation**  $O_i \subseteq 1, \dots, p$  l'ensemble des indices des variables pour lesquelles l'individu  $i$  est observé.

Si l'individu  $i$  appartient au cluster  $C_k$ , alors :

- GMM :

$$Z_i | c_i = k \sim g_k(\cdot | \theta_k) = \mathcal{N}(\cdot; \mu_k, \Sigma_k)$$

- Ignorable-GMM :

$$Z_i | c_i = k \sim \mathcal{N}(\cdot; \mu_k, \Sigma_k) = \prod_{j \in O_i} \mathcal{N}(z_{ij}; \mu_{kj}, \Sigma_{kj})$$

- Log-vraisemblance ignorable-GMM :

$$l(Z, \Psi) = \sum_{i=1}^n \log \sum_{k=1}^K \pi_k \prod_{j \in O_i} \mathcal{N}(z_{ij}; \mu_{kj}, \Sigma_{kj})$$

# Fuzzy c-means

Degré d'appartenance/probabilité d'appartenir à chaque cluster

$$\operatorname{argmin}_{c_1, \dots, c_K, \gamma_{11}, \dots, \gamma_{nK}} \sum_{k=1}^K \sum_{i=1}^n \gamma_{ik}^m \|x_i - c_k\|_2^2$$

Pour le nombre de clusters  $K$  et le paramètre de flou  $\alpha$  fixés :

- Attribuer aléatoirement des coefficients  $\gamma_{ik}$  à chaque individu ;
- Répéter jusqu'à convergence :

- Calculer les centroïdes  $c_{kj} = \frac{\sum_{i=1}^n \gamma_{ik}^\alpha \cdot x_{ij}}{\sum_{i=1}^n \gamma_{ik}^\alpha}$
- Recalculer les coefficients  $\gamma_{ik} = \frac{d(X_i, c_k)^{\frac{1}{\alpha-1}}}{\sum_{l=1}^K d(X_i, c_l)^{\frac{1}{\alpha-1}}}$

# Fuzzy c-means : Optimal Completion Strategy (Hathaway et Bezdek, 2001)

Pour le nombre de clusters  $K$  et le paramètre de flou  $\alpha$  fixés :

- **Initialiser** les centroïdes  $c_{k,k=1,\dots,K}$  et **initialiser** les valeurs manquantes par celles du centroïde le plus proche au sens de la **distance locale** ;
- Répéter pour les itérations  $q = 1, 2, \dots$  jusqu'à convergence :
  - Calculer les coefficients  $\gamma_{ik}^{(q)} = \frac{d(X_i, c_k^{(q-1)})^{\frac{1}{\alpha-1}}}{\sum_{l=1}^K d(X_i, c_l^{(q-1)})^{\frac{1}{\alpha-1}}}$
  - Calculer les nouveaux centroïdes de chacun des  $K$  clusters

$$c_{kj}^{(q)} = \frac{\sum_{i=1}^n \gamma_{ik}^{\alpha} \cdot x_{ij}}{\sum_{i=1}^n \gamma_{ik}^{\alpha}}$$

- Calculer les valeurs manquantes  $x_{ij}^{(q+1)} = \frac{\sum_{k=1}^K \gamma_{ik}^{(q)\alpha} \cdot c_{kj}^{(q)}}{\sum_{k=1}^K \gamma_{ik}^{(q)\alpha}}$

# Simulations

## Configurations de données simulées

### Configuration de référence (**modèle I**)

- $n = 750$  observations,  $p = 8$  variables
- $K = 3$  clusters :  $n_w = 250, w \in \{1, 2, 3\}$
- moyennes  $\mu_w, w \in \{1, 2, 3\}$  :  
 $\mu_1 = (0, 0, 0, 0, \Delta, \Delta, 0, \Delta^2)$   
 $\mu_2 = (0, 0, 0, 0, -\Delta, -\Delta, -\Delta, 0)$   
 $\mu_3 = (0, 0, 0, 0, -\Delta, \Delta, \Delta, -\Delta^2)$
- Matrices de variance-covariance

$$\Sigma_{w,w=1,2,3} = \Sigma(\rho) = \begin{pmatrix} I_4 & \mathbf{0} \\ \mathbf{0} & \begin{matrix} 1 & \rho & \rho & \rho \\ \rho & 1 & \rho & \rho \\ \rho & \rho & 1 & \rho \\ \rho & \rho & \rho & 1 \end{matrix} \end{pmatrix} \text{ avec } \rho = 0.3.$$

# Simulations

Pourcentages de valeurs manquantes :  $\tau = 10\%$ ,  $25\%$ ,  $40\%$

Mécanismes de valeurs manquantes :

- MCAR<sup>a</sup> :  $Prob(r_{ij} = 0) = \tau$
- MAR<sup>b</sup> 1 :  $Prob(r_{ij} = 0) = \Phi(a_\tau + x_{i1})$
- MAR 2 :  $Prob(r_{ij} = 0) = \Phi(a_\tau + x_{i8})$

---

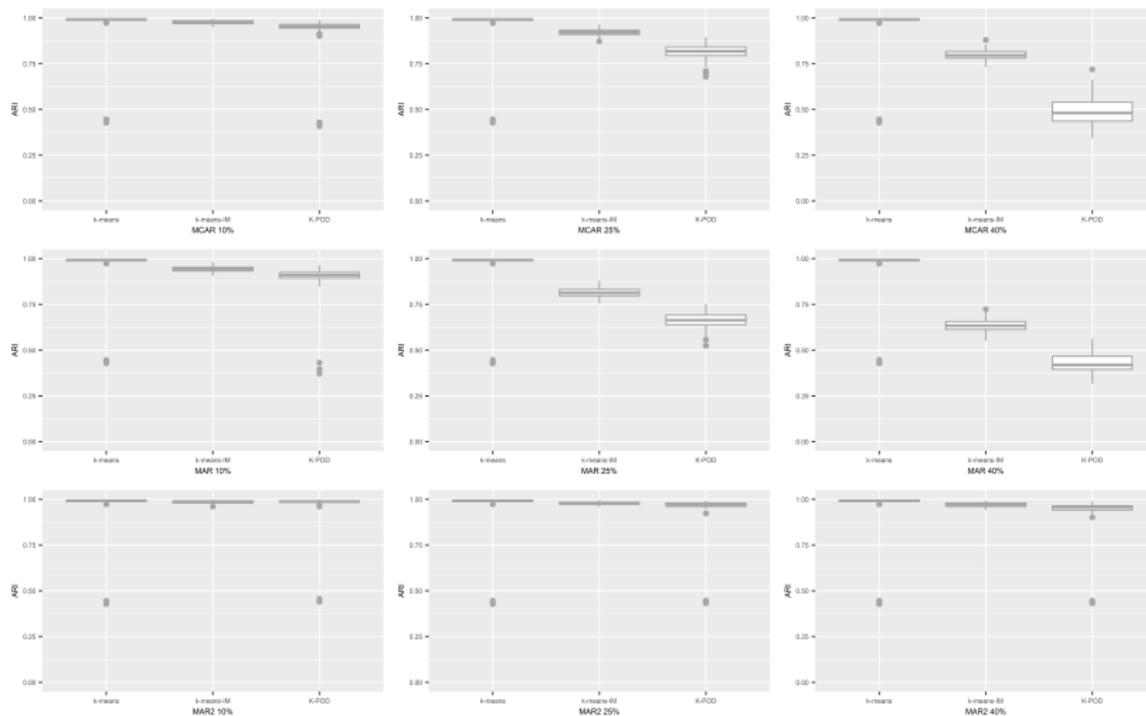
a. Missing Completely At Random

b. Missing At Random

- $\Phi$  la fonction de répartition de la loi normale standard
- $\tau$  une constante contrôlant la proportion de valeurs manquantes attendue
- $R = (r_{ij})_{1 \leq i \leq n, 1 \leq j \leq p}$  le dispositif des valeurs manquantes,  $r_{ij} = 1$  si  $x_{ij}$  est observé et  $r_{ij} = 0$  si  $x_{ij}$  est manquant

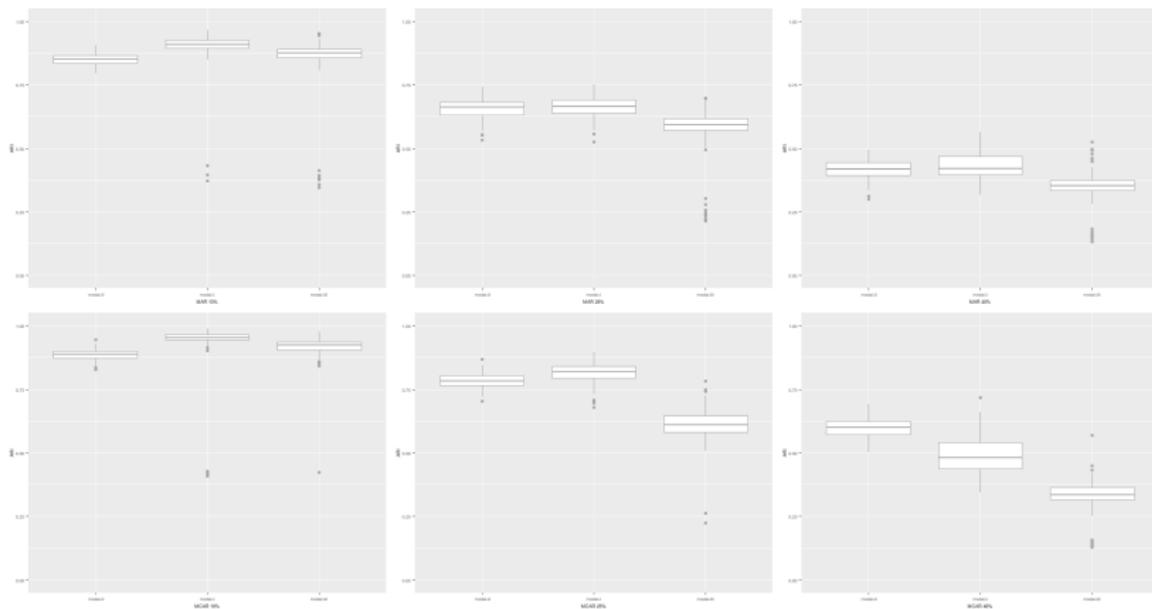
# Résultats

## K-means



## Résultats

## K-POD



- Le consensus pondéré n'est pas meilleur que le consensus simple pour l'imputation multiple en clustering ;
- L'indépendance des tableaux imputés n'est pas très importante pour effectuer une imputation multiple en clustering ;
- L'imputation multiple est meilleure que les méthodes directes de clustering sur données incomplètes.