

# 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  $\mathbf{Z} = \begin{matrix} z_{ij} & i & n \\ & j & p \end{matrix}$  ensemble de données numériques

**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  $c_k \in \mathcal{M}_p$

**Objectif** Identifier  $C_i$  pour chaque individu  $i$  à partir de  $\{z_{i1}^0, \dots, z_{ip}^0\}$

## 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...  $\mathbf{z}_i = \mathbf{z}_i^{obs}, \mathbf{z}_i^{misso}$

## 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)

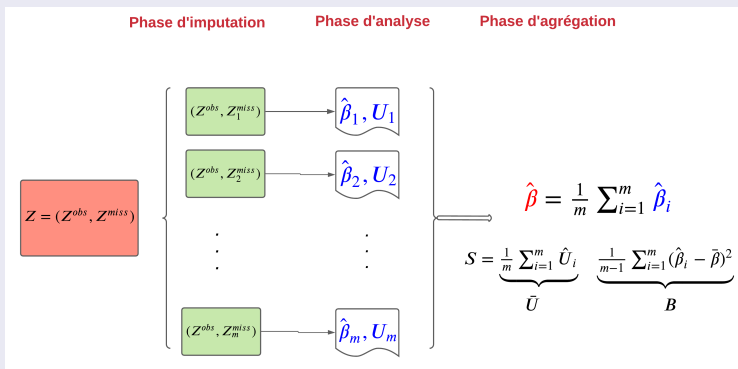
Optimal Completion Strategy of fuzzy c-means (Hathaway et Bezdek, 2001)

Simulations

## ④ Conclusion

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

## Principe



Estimation unique de  $\hat{\beta}$  et sa variance :  $S = \bar{U} + B$

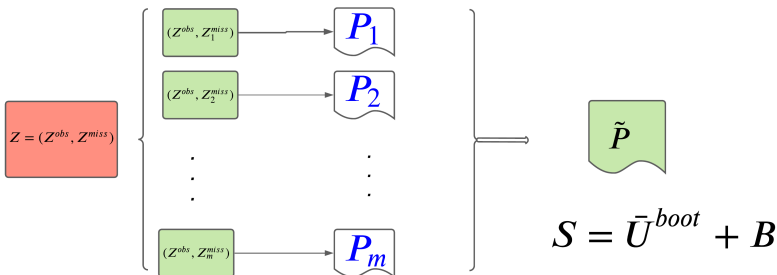
Qualité de  $\hat{\beta}$  dépend de l'indépendance des tableaux imputés conditionnellement aux valeurs observées

# 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)

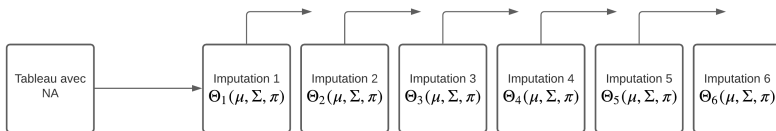
## Approche d'imputation

On suppose que  $Z_{ij} C_i$ ,  $N^1 Z_{ij} C_i, C_i^0$

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  $^1, , ^0$ .

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



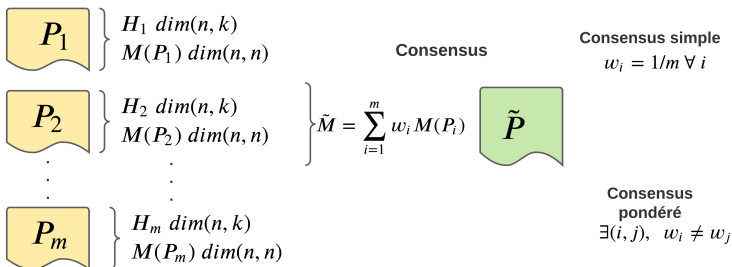


# Phase d'agrégation : consensus de partitions

## 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 :

$$\min_H \| \tilde{M} - HH^0 \|_F \quad (1)$$

$H = f, g^{1nk^0}$  : matrice d'indicatrices

$M = f, g^{1nr^0}$  : matrice de connectivité.

En posant  $\tilde{H} = H^1 H^0 H^0$  et  $D = \text{diag}^1 H^0 H^0$  (1) devient :

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

$$\tilde{H} = \tilde{H} \frac{S}{\frac{1}{\tilde{H} \tilde{H}^0 \tilde{M} \tilde{H}^0} \frac{1}{\tilde{M} \tilde{H} D^0}} \quad (3)$$

$$D = D \frac{S}{\frac{1}{\tilde{H}^0 \tilde{M} \tilde{H}^0} \frac{1}{\tilde{H}^0 \tilde{H} D \tilde{H}^0 \tilde{H}^0}} \quad (4)$$

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

Introduire  $w$  tel que :

$$w = (w_1, w_2, \dots, w_m)^T, \quad w_t \geq 0, \quad \sum_{t=1}^m w_t = 1 \quad (5)$$

On définit  $\tilde{M} = \sum_{t=1}^m w_t M^1 P_t^0$  et le problème (1) devient :

$$\min_{w, \tilde{H}} \|\tilde{M} - \tilde{H}\tilde{H}^0\| \quad (6)$$

- ① Chercher  $\tilde{H}$  en fixant  $w$  en utilisant NMF
- ② Chercher  $w$  en fixant  $\tilde{H}$  en résolvant  $\min_w w^0 A w = b^0$ , 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

$n =$  observations,  $p =$  variables

$K =$  clusters :  $n_w =$  ,  $w \in \{1, \dots, g\}$

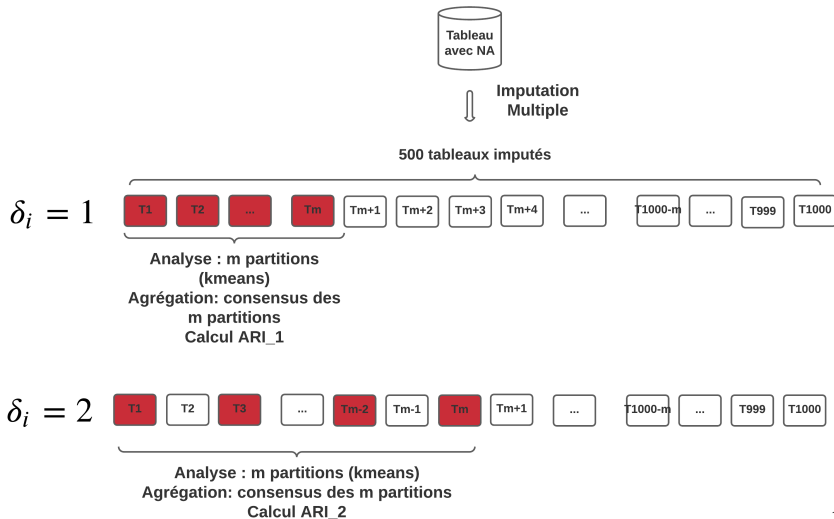
moyennes  $\mu_w, \Sigma_w, \gamma_w$  :  $\mu_w = 1, \dots, 0$   
 $\Sigma_w = 1, \dots, 0$

Matrice variance covariance =  $\begin{pmatrix} \sigma^2 & & & & \\ & \sigma^2 & & & \\ & & \dots & & \\ & & & \dots & \\ & & & & \sigma^2 \end{pmatrix}$

30% de données manquantes selon un mécanisme MAR.

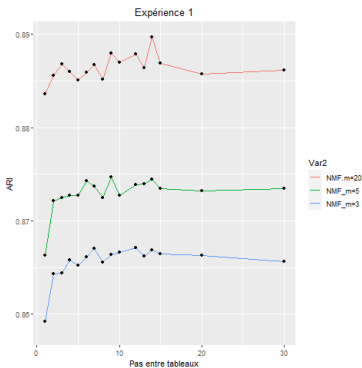
## Plan

Plan d'expérimentation  
 $m \geq 2$ ,  $f$ ,  $g$  : nombre de tableaux

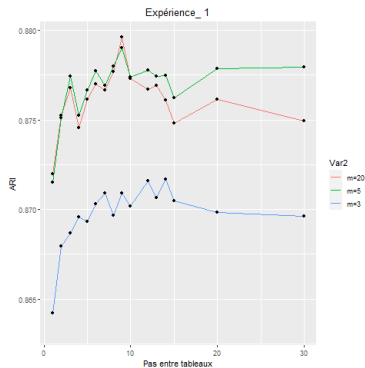


## Résultats

## Résultats NMF



## Résultats WNMF



## Résultats

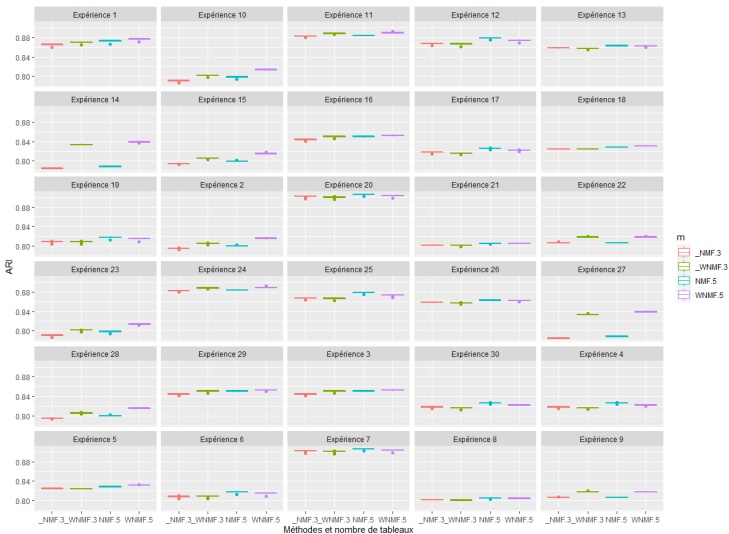
Moins sensible à l'indépendance entre tableaux.

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

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

( )

$$\min_{H, C} \sum_{j \in Z} \|H C\|_F$$

## k-POD

$$\min_{c_1, \dots, c_K, C_1, \dots, C_K} \sum_{k=1}^K \sum_{i \in C_k} \sum_{j \in C_k} \|Z_{ij} - c_k\|^2$$

( )

$$\min_{H, C} \sum_{j \in P} \|Z^0 - P^{-1} H C\|_F$$

**Pas de solution explicite**

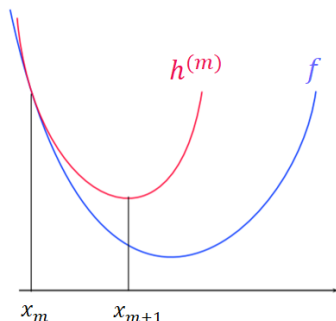
$$H = [h_{ik}]_{i \in M, k \in K} \quad h_{ik} = \begin{cases} s_i & i \in C_k \\ \text{sinon} & \end{cases}$$

$f_1, \dots, f_m$   $f_1, \dots, f_p$  ! sous-ensemble des indices des valeurs observées

$$P^{-1} Z^0_{ij} = \begin{cases} Z_{ij} & \text{si } i \in f_1 \\ s_i & \text{si } i \in f_2 \end{cases} \dots$$

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

## Algorithme de majoration-minimisation



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

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

Fonction de perte  $f^1(H, C^0)$   
 $\|j_j P^{-1} Z^0 - P^{-1} H C^0\|_F$

Fonction majorante  $g^{1,q^0}$  en  ${}^1H^{1,q^0}, C^{1,q^0}$   
 $\|j_j P^{-1} Z^0 - P^{-1} H C^0\|_F \leq \|j_j P_c^{-1} Z^0 - P_c^{-1} H^{1,q^0} C^{1,q^0}\|_F$   
 =  
 $\|j_j Z^{1,q^0} - H C\|_F$

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

Pour le nombre  $K$  de clusters fixé :

Initialiser les valeurs manquantes et obtenir la base complétée  $Z^{1,0}$ , puis calculer la partition initiale  $H^{1,0}, C^{1,0}$  par k-means ;

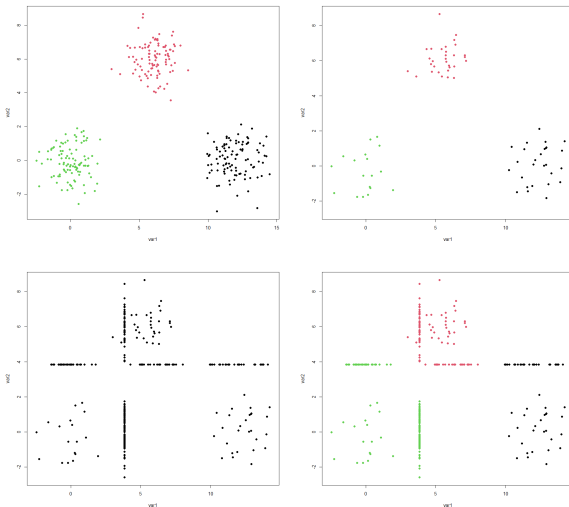
Répéter pour les itérations  $q = 1, 2, \dots$  jusqu'à convergence :

remplacer les valeurs manquantes par celles des centroïdes de leur groupe respectif pour obtenir  $Z^{1,q}$  ;

estimer le minimum  $H^{1,q}, C^{1,q}$  de la fonction majorante  $g^{1,q}(H, C) = J(H^{1,q}, C^{1,q})$  par k-means sur  $Z^{1,q}$ .

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

## Initialisation des valeurs manquantes





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

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

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

GMM :

$$Z_{ij}c_i = k \quad g_k^{1:j} \quad k^0 = \begin{matrix} 1; & k, & k^0 \end{matrix}$$

Ignorable-GMM :

$$Z_{ij}c_i = k \quad \begin{matrix} 1; & k, & k^0 = \\ & \bigcirc & \\ & \sum_{j \in O_i} & \end{matrix} \quad \begin{matrix} 1Z_{ij}; & k_j, & k_j^0 \end{matrix}$$

Log-vraisemblance ignorable-GMM :

$$l^1 Z, \quad \circ = \sum_{i=1}^{\tilde{n}} \log \sum_{k=1}^{\tilde{K}} \begin{matrix} \bigcirc \\ \sum_{j \in O_i} \end{matrix} \quad \begin{matrix} 1Z_{ij}; & k_j, & k_j^0 \end{matrix}$$

# Fuzzy c-means

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

$$\underset{c_1, \dots, c_K}{\operatorname{argmin}} \sum_{k=1}^K \sum_{i=1}^n d_{ik}^{2m} x_i c_k$$

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

Attribuer aléatoirement des coefficients  $c_{ik}$  à chaque individu ;

Répéter jusqu'à convergence :

Calculer les centroïdes  $c_{kj} = \frac{\sum_{i=1}^n d_{ik}^{2m} x_{ij}}{\sum_{i=1}^n d_{ik}^{2m}}$

Recalculer les coefficients  $c_{ik} = \frac{d_{ik}^{2m} x_i}{\sum_{l=1}^K d_{il}^{2m}}$

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

Pour le nombre de clusters  $K$  et le paramètre de flou 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  $w_{ik}^{1q} = \frac{d^{\alpha}(X_i, c_k^{1q})^{-\alpha}}{\sum_{l=1}^K d^{\alpha}(X_i, c_l^{1q})^{-\alpha}}$

Calculer les nouveaux centroïdes de chacun des  $K$  clusters

$$c_{kj}^{1q} = \frac{\sum_{i=1}^n w_{ik}^{1q} X_{ij}}{\sum_{i=1}^n w_{ik}^{1q}}$$

Calculer les valeurs manquantes  $X_{ij}^{1q} = \frac{\sum_{k=1}^K w_{ik}^{1q} c_{kj}^{1q}}{\sum_{k=1}^K w_{ik}^{1q}}$

# Simulations

## Configurations de données simulées

Configuration de référence (modèle I)

$n$  = observations,  $p$  = variables

$K$  = clusters :  $n_w = w_1, w_2, \dots, w_g$

moyennes ( $w_1, w_2, \dots, w_g$ ) :

$$= \begin{matrix} 1 & & & & & & & & 0 \\ & 1 & & & & & & & \\ & & \ddots & & & & & & \\ & & & \ddots & & & & & \\ & & & & 1 & & & & \\ & & & & & \ddots & & & \\ & & & & & & 1 & & \\ & & & & & & & \ddots & \\ & & & & & & & & 1 & & 0 \end{matrix}$$

Matrices de variance-covariance

$w_1, w_2, \dots, w_g = 1, 0 = \begin{matrix} \textcircled{0} & / & \textcircled{0} \\ - & & \\ - & & \\ - & & \\ - & & \\ - & & \\ - & & \\ - & & \\ - & & \\ - & & \\ \ll & & \neg \end{matrix}$  avec  $= \dots$

# Simulations

Pourcentages de valeurs manquantes : = %, %, %

Mécanismes de valeurs manquantes :

$$\text{MCAR}^a : \text{Prob}^1 r_{ij} = 0 =$$

$$\text{MAR}^b 1 : \text{Prob}^1 r_{ij} = 0 = 1 a_{,j} x_i^0$$

$$\text{MAR} 2 : \text{Prob}^1 r_{ij} = 0 = 1 a_{,j} x_i^0$$

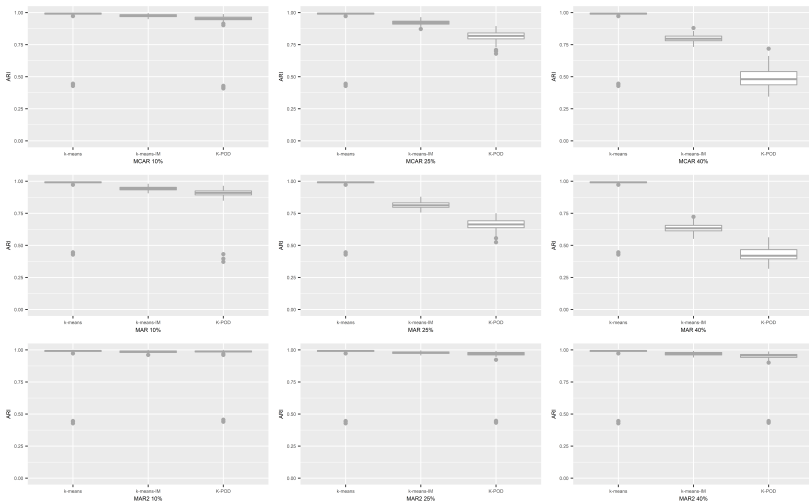
- 
- a. Missing Completely At Random  
b. Missing At Random

la fonction de répartition de la loi normale standard  
une constante contrôlant la proportion de valeurs  
manquantes attendue

$R = 1 r_{ij}^0$   $i = 1, \dots, n$ ,  $j = 1, \dots, p$  le dispositif des valeurs manquantes,  $r_{ij} =$   
si  $x_{ij}$  est observé et  $r_{ij} =$  si  $x_{ij}$  est manquant

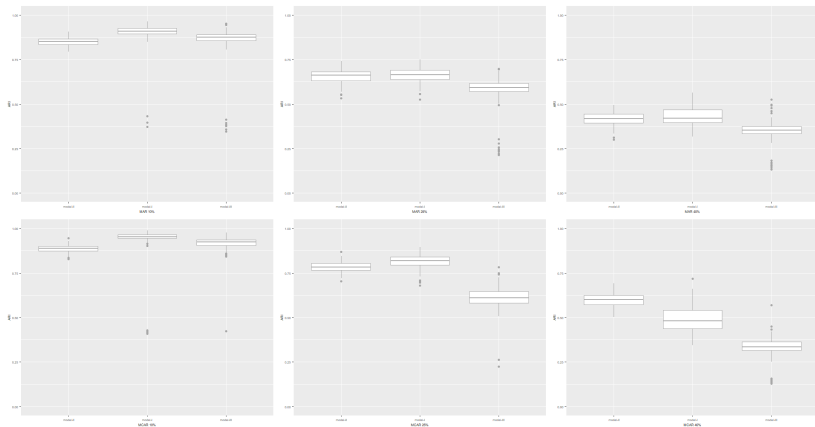
## Résultats

## 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.