

Apprentissage Machine / Statistique

Imputation de Données Manquantes

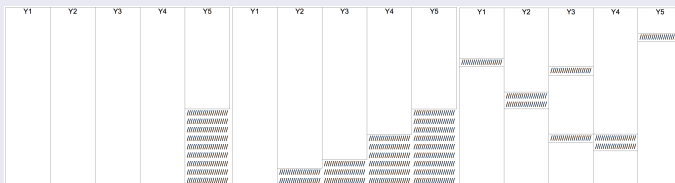
PHILIPPE BESSE

INSA – Dpt GMM
Institut de Mathématiques – ESP

Types de données manquantes

- Missing completely at random (MCAR)
- Missing at random (MAR) : absence liée à une ou plusieurs autres variables
- Missing not at random (MNAR) : absence dépend de la variable
- Type des variables : quantitatives ou qualitatives

Répartition univariée, monotone, arbitraire



Analyse sans imputation

- **Suppression**
 - Lignes : perte de précision, biais si non MCAR
 - Colonnes
- **Méthodes tolérantes**
 - CART : divisions de substitution
 - NIPALS en PLS (imputation implicite)
 - XGBoost

Imputations unidimensionnelles

- Complétion stationnaire : plus fréquente, dernière valeur, *Last Observation Carry Forward (LOCF)*
- Combinaison linéaire : moyenne, médiane

Imputations par régression

- k plus proches voisins
- Régression, régression locale (LOESS)
- NIPALS (PLS, PCA) ou SVD

Imputation par SVD

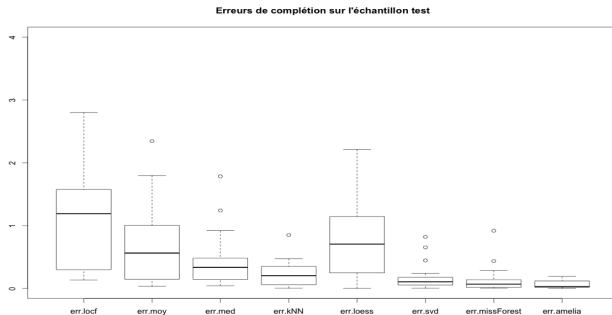
- Matrice \mathbf{Y}^0 complétée de \mathbf{Y} par des moyennes
- **Tant que** $\| \mathbf{Y}^i - \mathbf{Y}^{i+1} \| / \| \mathbf{Y}^i \| < \epsilon$
 - $\hat{\mathbf{Y}}$: SVD tronquée (rang J) de \mathbf{Y}^i centrée
 - \mathbf{Y}^{i+1} est \mathbf{Y} complétée par les éléments de $\hat{\mathbf{Y}}$

Imputation multiple : Amelia II

- Imputer plusieurs valeurs
- Combinaison des résultats
- Réduction du bruit, mesure d'incertitude
- Hypothèse de normalité
- Algorithme EMB (Expectation Minimisation et Bootstrap)

Imputation par forêts aléatoires

- ❶ Complétion “naïve” des valeurs manquantes.
- ❷ k indices des colonnes de Y triées par quantité croissante de valeurs manquantes ;
- ❸ **Tant que** γ n'est pas atteint **faire**
 - ❶ Y_{imp}^{old} = matrice précédemment imputée
 - ❷ **Pour** s dans k **faire**
 - ❶ Ajuster $y_{obs}^{(s)} \sim x_{obs}^{(s)}$ par forêt aléatoire
 - ❷ Prédire $y_{mis}^{(s)}$ avec les régresseurs $x_{mis}^{(s)}$
 - ❸ Y_{imp}^{new} : nouvelle matrice complétée par les $y_{mis}^{(s)}$
 - ❸ mettre à jour le critère γ



EBP - Erreurs de complétion sur un échantillon test de 10%