

IA : prétraitement des données

Françoise Bouvet
francoise.bouvet@ijclab.in2p3.fr

Les différentes étapes du prétraitement

1. Visualiser et examiner les données pour détecter des corrélations
2. Isoler la(les) variable(s) cible(s) et l'encoder
3. Détecter les données manquantes et y pallier
4. Sélectionner les bonnes variables d'entraînement ou en créer d'autres
5. Normaliser les variables d'entraînement continues
6. Séparer l'ensemble en sous-ensembles (entraînement/validation/test)

Préambule : représentation des données

Entrée X

$$\mathbf{X} = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & x_3^{(1)} & \dots & x_n^{(1)} \\ x_1^{(2)} & x_2^{(2)} & x_3^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_1^{(m)} & x_2^{(m)} & x_3^{(m)} & \dots & x_n^{(m)} \end{bmatrix}$$

Sortie Y

$$\mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

Exemple

Prédire le risque d'attaque cardiaque à partir de l'âge, du sexe, de la pression artérielle, de la glycémie, du taux de cholestérol et du rythme cardiaque.

- L'ensemble de données contient 200 patients $\rightarrow m = 200$
- $n = 6$
- y prend une des 5 valeurs suivantes : 1, 2, 3, 4, 5 (de risque faible à élevé)

m : nombre d'échantillons

n : nombre de paramètres qui caractérisent chaque échantillon

Ensemble de données : m échantillons

Données catégorielles

- Encodage ordinal

- Attribution d'une valeur entière à chaque classe
- Génère un ordre artificiel entre les catégories

Ex : encoder 3 couleurs : bleu -> 1 ; vert -> 2 ; rouge -> 3

- Encodage « One-hot »

- La catégorie est représentée par un vecteur de 0, sauf une valeur à 1
- La taille du vecteur est égale au nombre de classes

Ex : bleu -> $\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$; vert -> $\begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$; rouge -> $\begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$

Données continues (1)

Beaucoup d'algorithmes de ML sont très sensibles aux valeurs absolues des valeurs d'entrée alors que l'information cruciale se situe plutôt au niveau des valeurs relatives

➤ **normalisation** des données en entrée

Centrage et normalisation

$$x_i^m \leftarrow \frac{x_i^m - \mu_i}{\sigma_i}$$

où :

- $X = \{x^m\}, m \in \{1, M\}, x^m \in R^N$
- $\mu_i = \frac{1}{M} \sum_{m=1}^M x_i^m$
- $\sigma_i^2 = \frac{1}{M} \sum_{m=1}^M (x_i^m - \mu_i)^2$

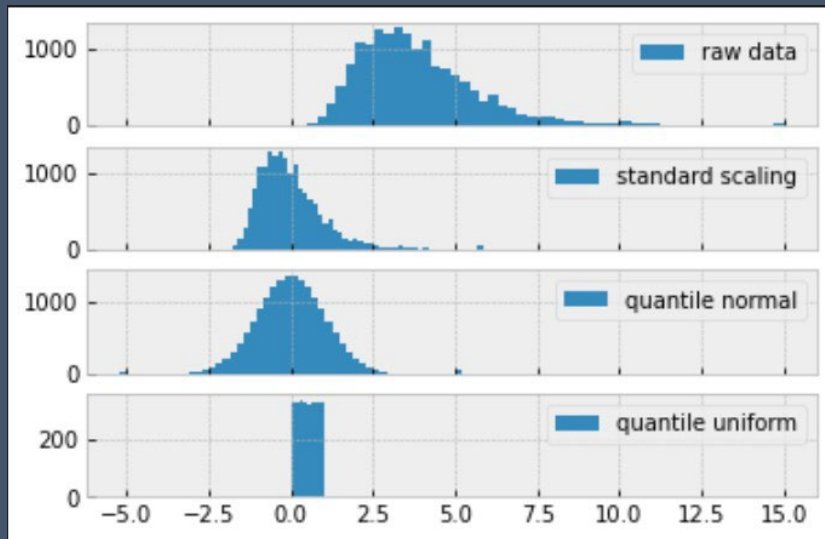
Cas particulier des données « **sparse** » (beaucoup de 0)

➤ Normaliser sans centrer pour ne pas les rendre « denses » (mémoire)

Données continues (2)

Quantile -> float

Transformation non linéaire des données pour obtenir une distribution spécifique



Discretisation (binning) -> classe

Regroupe les valeurs en k classes (bins)

Techniques classiques :

- **Uniforme** : chaque bin a la même étendue en valeur
- **Quantile** : basé sur les percentiles, chaque bin a le même nombre de valeurs
- **Clustering** : des groupes sont identifiés et des exemples sont assignés à chaque groupe

Problèmes courants : manque de données

Solution usuelle : **augmentation de données**

- Transformations géométriques
 - Rotation
 - Translation
 - Distorsion
 - Modification du contraste
- Modèles génératifs
 - GAN
 - Transformers
 - Variational Auto Encoder

Problèmes courants : données non balancées

Imbalanced data : grande différence entre la quantité de données dans chaque classe

- Problème classique en **classification**



- Attention à la **métrique**
- Solution usuelle :
 - **Sous-échantillonnage** des classes majoritaires par suppression
 - **Sur-échantillonnage** des classes minoritaires par duplication / synthèse

Problèmes courants : données manquantes

Données manquantes : NaN

- Ignorer l'échantillon
- Faire de l'imputation de donnée
 - Technique statistique de remplissage de données manquantes à partir des autres valeurs
 - Ex : catégorie la plus présente, moyenne, médiane, valeur de l'échantillon précédent, etc

```
[[ 6.  2.9  4.5  1.5]
 [ 5.9  3.  5.1  1.8]
 [ 4.4  3.  1.3  0.2]
 [ 5.1  3.3  nan  nan]
 [ 5.  3.5  1.6  0.6]
 [ 5.4  3.4  nan  nan]
 [ 5.7  3.8  nan  0.3]
 [ 5.6  2.5  3.9  nan]
 [ 7.7  2.6  6.9  2.3]
 [ 5.8  2.7  5.1  1.9]
 [ 6.7  3.1  5.6  2.4]
 [ 4.8  3.4  1.9  nan]
 [ 7.2  3.2  6.  1.8]
 [ 4.4  2.9  nan  nan]
 [ 6.9  3.2  5.7  2.3]
 [ 5.5  4.2  1.4  nan]
 [ 6.3  2.3  4.4  1.3]
 [ 7.  3.2  4.7  1.4]
 [ 5.8  2.7  nan  nan]
 [ 6.8  2.8  4.8  1.4]
 [ 5.4  3.9  1.7  nan]
 [ 7.6  3.  6.6  2.1]
 [ 7.7  2.8  6.7  2. ]
 [ 5.  3.3  nan  0.2]
 [ 5.9  3.  4.2  1.5]
 [ 6.1  2.8  4.  1.3]
 [ 5.  3.6  1.4  0.2]
 [ 7.4  2.8  6.1  1.9]
 [ 6.3  2.5  5.  1.9]
 [ 6.7  3.3  5.7  2.5]]
```

```
from sklearn.impute import SimpleImputer
imp = SimpleImputer(strategy="median").fit(X_train)
X_median_imp = imp.transform(X_train)
```

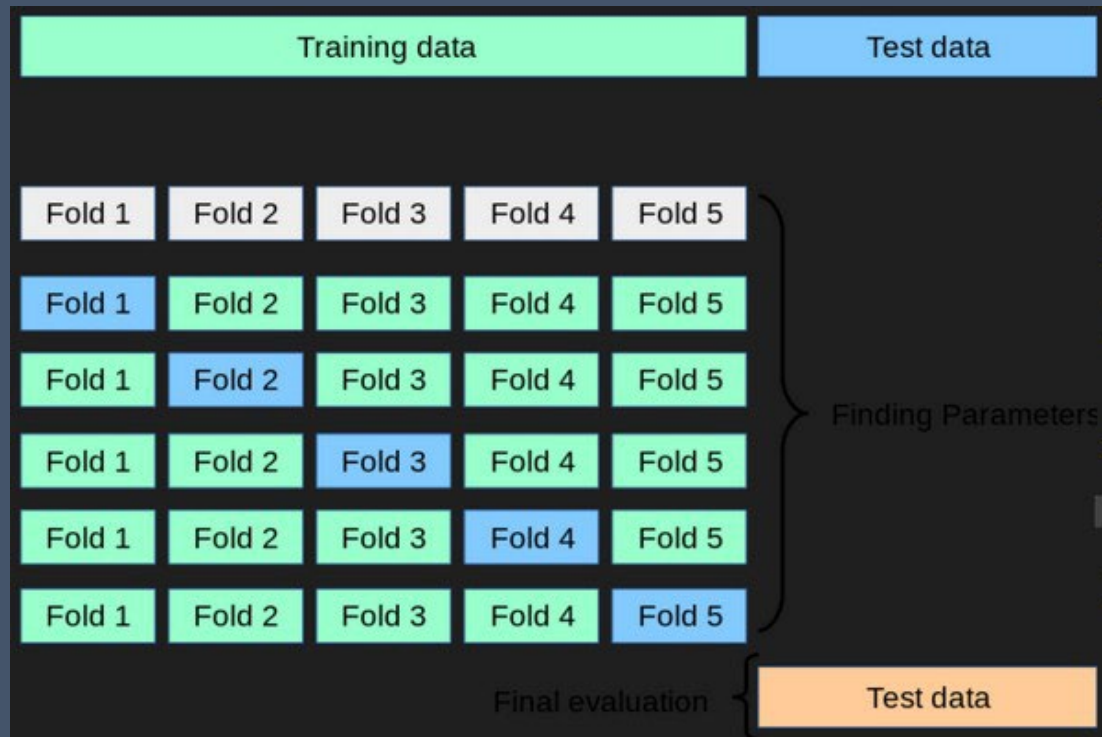
Imputation

```
array([[ 6.  ,  2.9 ,  4.5 ,  1.5 ],
 [ 5.9 ,  3.  ,  5.1 ,  1.8 ],
 [ 4.4 ,  3.  ,  1.3 ,  0.2 ],
 [ 5.1 ,  3.3 ,  4.116,  1.462],
 [ 5.  ,  3.5 ,  1.6 ,  0.6 ],
 [ 5.4 ,  3.4 ,  4.116,  1.462],
 [ 5.7 ,  3.8 ,  4.116,  0.3 ],
 [ 5.6 ,  2.5 ,  3.9 ,  1.462],
 [ 7.7 ,  2.6 ,  6.9 ,  2.3 ],
 [ 5.8 ,  2.7 ,  5.1 ,  1.9 ],
 [ 6.7 ,  3.1 ,  5.6 ,  2.4 ],
 [ 4.8 ,  3.4 ,  1.9 ,  1.462],
 [ 7.2 ,  3.2 ,  6.  ,  1.8 ],
 [ 4.4 ,  2.9 ,  4.116,  1.462],
 [ 6.9 ,  3.2 ,  5.7 ,  2.3 ],
 [ 5.5 ,  4.2 ,  1.4 ,  1.462],
 [ 6.3 ,  2.3 ,  4.4 ,  1.3 ],
 [ 7.  ,  3.2 ,  4.7 ,  1.4 ],
 [ 5.8 ,  2.7 ,  4.116,  1.462],
 [ 6.8 ,  2.8 ,  4.8 ,  1.4 ],
 [ 5.4 ,  3.9 ,  1.7 ,  1.462],
 [ 7.6 ,  3.  ,  6.6 ,  2.1 ],
 [ 7.7 ,  2.8 ,  6.7 ,  2. ],
 [ 5.  ,  3.3 ,  4.116,  0.2 ],
 [ 5.9 ,  3.  ,  4.2 ,  1.5 ],
 [ 6.1 ,  2.8 ,  4.  ,  1.3 ],
 [ 5.  ,  3.6 ,  1.4 ,  0.2 ],
 [ 7.4 ,  2.8 ,  6.1 ,  1.9 ],
 [ 6.3 ,  2.5 ,  5.  ,  1.9 ],
 [ 6.7 ,  3.3 ,  5.7 ,  2.5 ]])
```

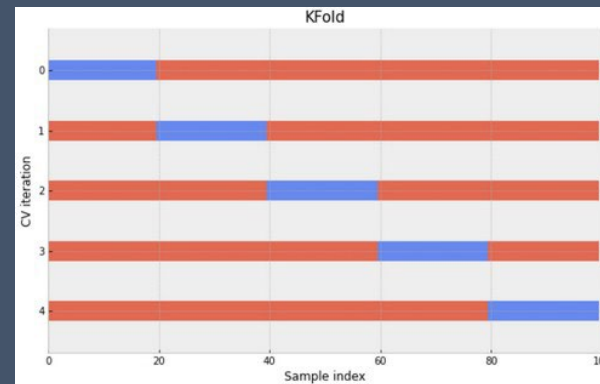
Problèmes courants : données bruitées

- Discrétisation (binning)
- Clustering
- Supprimer les valeurs aberrantes (outliers)

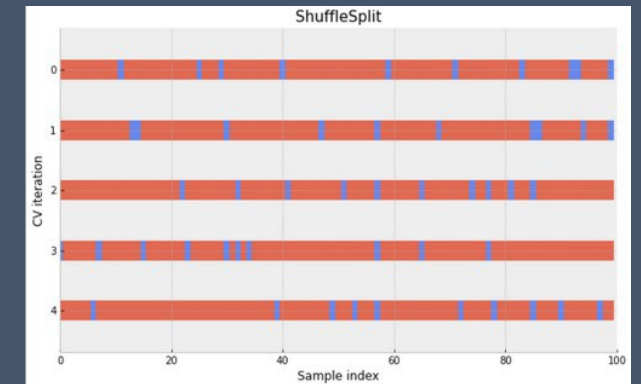
Cross-validation



1. Entraîner K modèles avec des jeux de données différents
2. Au choix :
 - Prendre le meilleur des K modèles
 - Prendre le meilleur des K modèles et le ré-entraîner sur l'ensemble des données
 - Garder les K modèles et se fier à l'avis du plus grand nombre



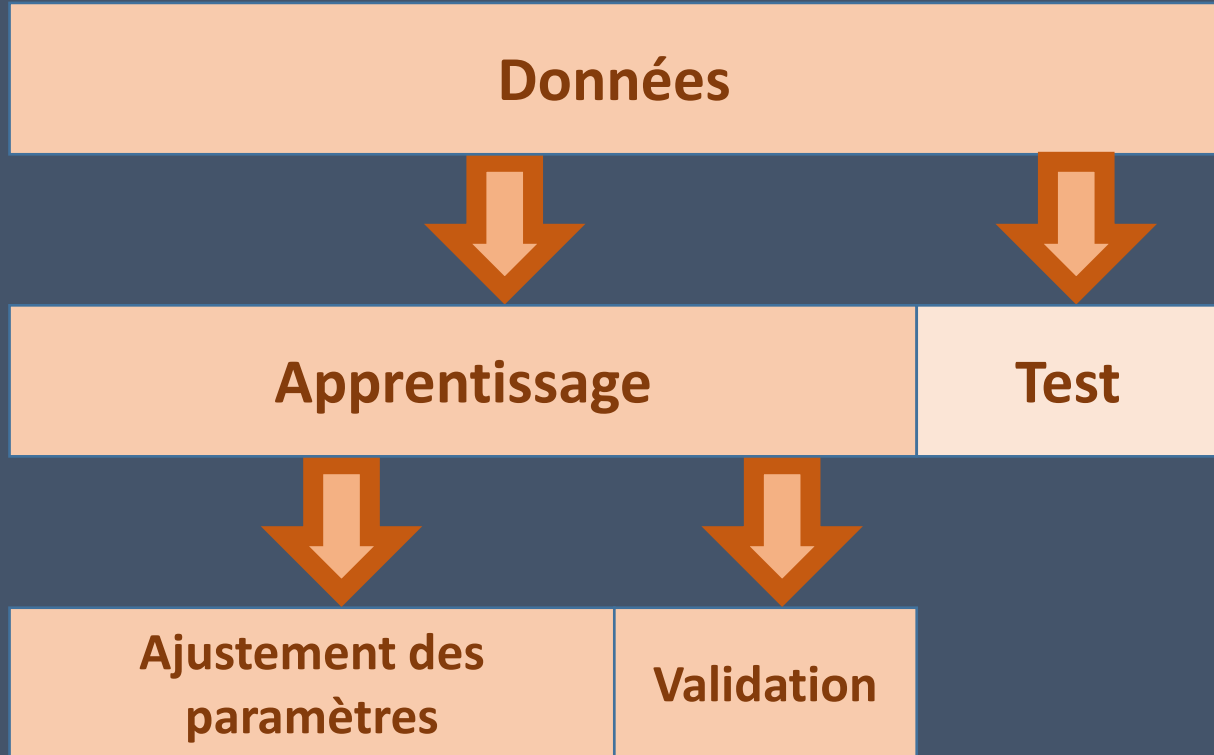
Séparation en KFold



Séparation aléatoire

➤ Méthode utilisée en **Machine Learning**

Séparation des données : test / apprentissage / validation



- L'échantillon d'apprentissage est divisé en 2 : une partie pour ajuster les paramètres du réseau ; une partie « validation » pour évaluer les performances du réseau à chaque itération.
- L'échantillon test sert uniquement à donner les performances finales du modèle.

Typiquement : 60%/20%/20% ; 70%/20%/10%

➤ Méthode utilisée en **Deep Learning**