





# Engénierie des caracteristiques (Feature Engineering)

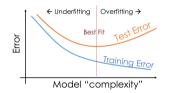
- Processus de sélection, modification ou création de caracteristiques d'éntrée
- Rappel des points clés
  - Sélection des caracteristiques identifier les caractéristiques les plus pertinantes
    - Lasso, Ridge, Elastic Net, RFE, PCA, ...
  - Transformation des caractéristiques
    - <u>Standardisation</u> (standardscaler), <u>one-hot</u> encoding, label encoding, normalisation (minmaxscaler), ...
  - Création de caractéristiques générer des variables à partir de la connaissance métier ou en combinant des variables existantes
    - Expansion polynomiale, (polynomialfeatures), ...





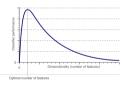
## Concepts liés

• Sur/sous apprentissage



Split train-test

• Fléau de la surdimensionalité



Sélection de caracteristiques (Lasso, Ridge, ...)

- Biais dans les échantillons
  - Les données ne représentent pas fidèlement la population réelle

Stratified train/test split, ...







# Validation Croisée (Cross Validation)

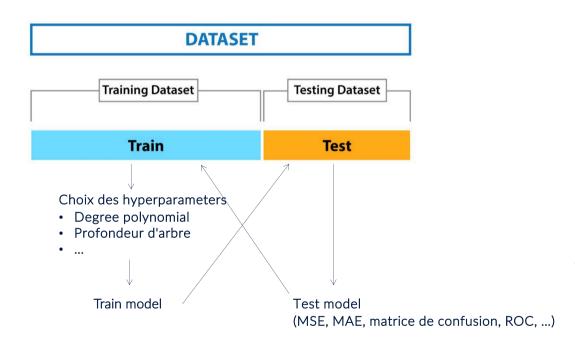
- Définition
  - Évaluer la capacité d'un modèle à généraliser sur des données non vues à l'entrainement
  - En divisant le jeu de données en sous-ensembles d'entraînement et de test <u>plusieurs</u> fois

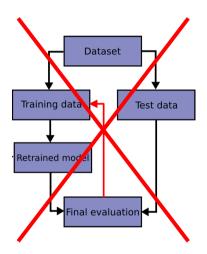




### Validation Croisée

• Example d'une approche <u>naïve et inexacte</u> d'entrainement machine





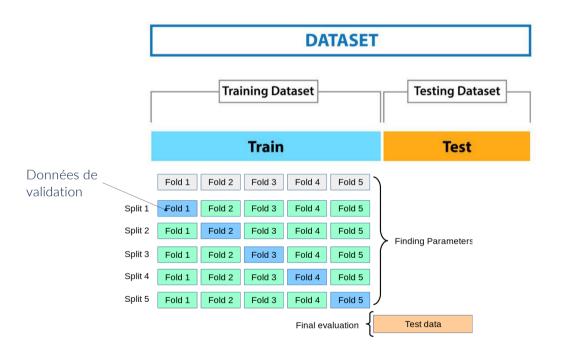
Fuite d'information des données de test dans les données d'entrainement! - modèle trop optimiste par choix des hyperparamètres

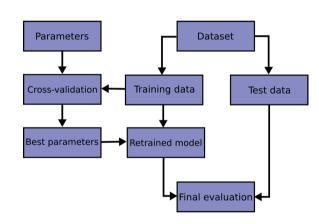




### Validation Croisée

• Méthode de validation croisée





- Validation croisée en k blocs (k-Fold Cross-Validation)
- Validation croisée "Leave-One-Out" (LOOCV)
- k-Fold stratifié
- Division temporelle (Time-Series Split)





## Validation Croisée

- Importance de la validation croisée
  - Réduit le risque de sur/sous-estimer les performances d'un modèle
  - Pour ajuster les <u>hyperparamètres</u> et évaluer les performances, en gardant le jeu de test totalement séparé pour l'évaluation finale





#### Fonctions scikit-learn utiles

- sklearn.pipeline.Pipeline
- sklearn.model\_selection.cross\_val\_score methode de validation croisée
- sklearn.model\_selection.GridSearchCV validation croisée avec selection des hyperparamètres optimums
- sklearn.feature\_selection.RFE Recursive Feature Elimination (RFECV), pour réduire l'espace des caractéristiques

Ref: https://scikit-learn.org/stable/modules/cross\_validation.html





