

Sklearn, Généralisation & Overfitting

Jour 2 — Après-midi

Julien Rolland

Formation M2 Développement Fullstack

Jour 2

- 1 Le Vrai Objectif : Généraliser
- 2 L'Ennemi N°1 : l'Overfitting
- 3 Biais-Variance
- 4 Méthodologie : Train / Validation / Test
- 5 Cross-Validation
- 6 Scikit-Learn

Le Piège

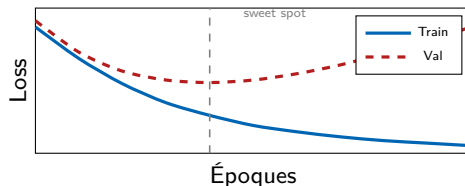
- Minimiser $J(\Theta, X_{\text{train}})$ = problème résolu
- La performance sur le train **ne compte pas**

Ce qu'on veut vraiment

- Données **non-vues** (unseen data)
- Capter la loi sous-jacente
- Ignorer le bruit du dataset

Définition

Généralisation : capacité à performer sur des exemples jamais vus pendant l'entraînement.



Underfitting vs Overfitting

Underfitting

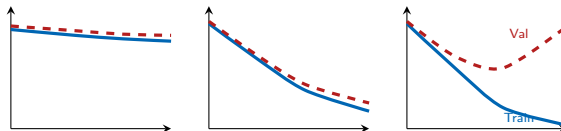
- Modèle trop **simple** — biais élevé
- $\text{Train} \approx \text{Val} \Rightarrow$ toutes deux élevées

Bon fit

- Complexité adaptée aux données
- $\text{Train} \approx \text{Val} \Rightarrow$ toutes deux basses

Overfitting

- Modèle trop **complexe** — variance élevée
- $\text{Train} \ll \text{Val} \Rightarrow$ Val explose

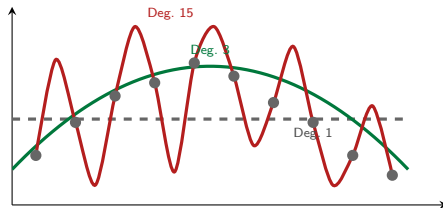


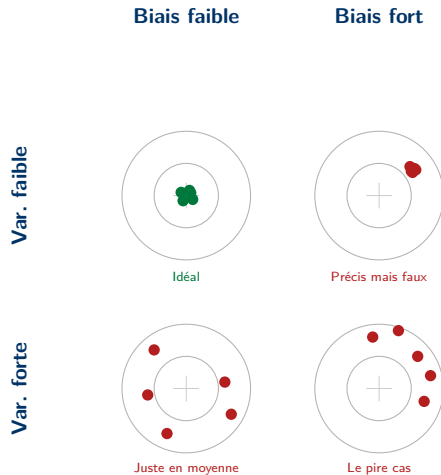
Théorème de Cover

- Dans \mathbb{R}^D , $N \leq D$ points sont toujours linéairement séparables
- Un modèle assez complexe peut **mémoriser** n'importe quel dataset
- Accuracy train = 100% ne signifie rien

Analogie MNIST

- Mémoriser les pixels exacts de chaque «3»
- Incapable de reconnaître un «3» écrit avec un stylo différent
- \Rightarrow 0% de généralisation





En Machine Learning

- **Biais** : erreur systématique
modèle trop simple
- **Variance** : sensibilité aux données
modèle trop complexe

$$\text{Erreur} = \text{Biais}^2 + \text{Variance} + \varepsilon$$

Leviers

- **Complexité** du modèle
- **Régularisation** (λ)
- **Volume** de données



3 ensembles distincts

- 1 **Train** — ajuster les paramètres Θ
- 2 **Validation** — choix du modèle et des hyperparamètres
- 3 **Test** — mesure finale

Règle d'or : Data Leakage

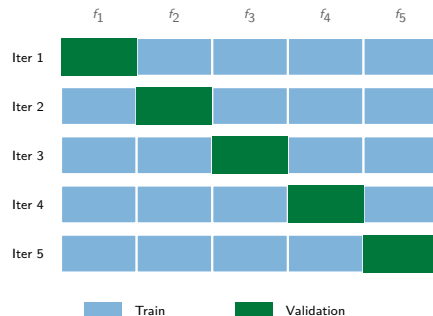
Ne **jamais** entraîner sur le test set.
Ne **jamais** utiliser le test pour choisir les HP.
Le test set doit rester **invisible**.

Le Problème

- Peu de données \Rightarrow split instable
- Un seul val set \Rightarrow variance élevée de l'estimation

Solution : K-Fold

- Découper le train en K morceaux (*folds*)
- K rotations : chaque fold sert de validation
- Score final = moyenne sur K runs
- Utilisation maximale des données disponibles



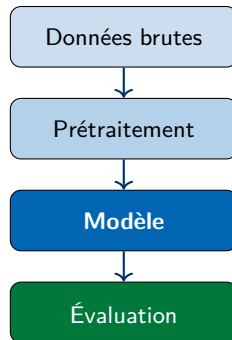
$$s = \frac{1}{K} \sum_{k=1}^K s_k$$

La bibliothèque ML de référence

- Open-source, basée sur NumPy / SciPy
- Standard de l'industrie pour le ML classique
- Interface cohérente pour tous les algorithmes

Ce qu'elle couvre

- **Prétraitement** : scaling, encodage, imputation
- **Modèles** : régression, classification, clustering
- **Validation** : cross-validation, grid search
- **Métriques** : accuracy, F1, AUC, ...



Interface universelle

- `.fit(X, y)` — entraînement
- `.predict(X)` — inférence
- `.score(X, y)` — évaluation
- `.transform(X)` — transformation

Modèles clés

- `LogisticRegression`
- `LinearRegression`
- `StandardScaler`
- `train_test_split`

```
1 from sklearn.linear_model import LogisticRegression
2 from sklearn.preprocessing import StandardScaler
3
4 scaler = StandardScaler()
5 X_tr = scaler.fit_transform(X_train)
6 X_te = scaler.transform(X_test)
7
8 clf = LogisticRegression()
9 clf.fit(X_tr, y_train)
10 print(clf.score(X_te, y_test))
```

Principe

- Enchaîner Scaler → Modèle
- `.fit()` appliqué séquentiellement
- Compatible avec la cross-validation

Sans Pipeline

Scaler fitté sur train+val
⇒ **Data Leakage !**

```
1 from sklearn.pipeline import Pipeline
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.linear_model import LogisticRegression
4
5 pipe = Pipeline([
6     ('scaler', StandardScaler()),
7     ('clf', LogisticRegression()),
8 ])
9
10 pipe.fit(X_train, y_train)
11 print(pipe.score(X_test, y_test))
```

train_test_split

- Découpe aléatoire du dataset
- test_size : proportion du test set
- random_state : reproductibilité
- stratify : préserve les proportions de classes

cross_val_score

- K-Fold intégré, compatible Pipeline
- Retourne un score par fold
- cv : nombre de folds
- scoring : métrique ('accuracy', ...)

```
1 from sklearn.model_selection import train_test_split, cross_val_score
2
3 X_train, X_test, y_train, y_test = train_test_split(
4     X, y, test_size=0.2, random_state=42, stratify=y)
5
6 scores = cross_val_score(pipe, X_train, y_train, cv=5)
7 print(f"CV: {scores.mean():.3f} +/- {scores.std():.3f}")
```