

Bootcamp IA — Atelier Guidé : Classification supervisée (Scikit-learn)

Département IA & Ingénierie des Données
Institut Supérieur d'Informatique (ISI Dakar)

Samedi 20 Septembre 2025
14h15 – 16h00

Objectif de l'atelier

Cet atelier vise à pratiquer **trois algorithmes de classification** avec Scikit-learn : **KNN**, **Régression Logistique** et **Arbre de Décision**. Les participants apprendront à :

- Charger et explorer rapidement des jeux de données standards.
- Construire, entraîner et évaluer des modèles de classification.
- Visualiser des matrices de confusion et (pour les tâches binaires) la courbe ROC.

Contexte & Données

Nous utiliserons trois datasets **intégrés à Scikit-learn** :

- **Iris** (3 classes) — pour KNN.
- **Breast Cancer** (binaire) — pour la Régression Logistique.
- **Wine** (3 classes) — pour l'Arbre de Décision.

Mission : entraîner, évaluer, comparer les trois approches et discuter des avantages/limites.

Organisation des participants

- Les 70 participants sont répartis en **14 groupes de 5**.
- Environnement : Python (Anaconda) ou Google Colab.
- Temps imparti : **1h15**.
- Livrables par groupe : code + graphes + mini-interprétation.

Plan de travail

1. **Préparer l'environnement** (imports, fonctions utilitaires).
2. **Partie A (KNN & Iris)** : split, entraînement, évaluation + étude de k .
3. **Partie B (Régression Logistique & Breast Cancer)** : pipeline avec standardisation, ROC/AUC.
4. **Partie C (Arbre de Décision & Wine)** : profondeur, critère (gini/entropy), visualisation.
5. **Synthèse** : tableau récapitulatif & discussion.

Code de départ (commun)

```
1 # ===== Imports communs =====
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5
6 from sklearn.model_selection import train_test_split, GridSearchCV
7 from sklearn.metrics import (
8     accuracy_score, confusion_matrix, ConfusionMatrixDisplay,
9     classification_report, roc_curve, auc
10 )
11
12 # Pour des graphiques plus nets
13 plt.rcParams["figure.dpi"] = 120
14
15 def plot_conf_mat(y_true, y_pred, labels=None, title="Matrice de
16     confusion"):
17     cm = confusion_matrix(y_true, y_pred, labels=labels)
18     disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=
19         labels)
20     disp.plot(values_format='d')
21     plt.title(title)
22     plt.tight_layout()
23     plt.show()
```

Partie A — KNN sur *Iris*

```
1 from sklearn.datasets import load_iris
2 from sklearn.neighbors import KNeighborsClassifier
3
4 # Charger les données
5 iris = load_iris(as_frame=True)
6 X, y = iris.data, iris.target
7 labels = iris.target_names
8
9 # Split train/test
10 X_train, X_test, y_train, y_test = train_test_split(
11     X, y, test_size=0.30, random_state=42, stratify=y
12 )
13
14 # Entraînement KNN avec k=3
15 knn = KNeighborsClassifier(n_neighbors=3)
16 knn.fit(X_train, y_train)
17
18 # Évaluation
19 y_pred = knn.predict(X_test)
20 acc = accuracy_score(y_test, y_pred)
21 print(f"Accuracy (k=3): {acc:.3f}")
22 plot_conf_mat(y_test, y_pred, labels=labels, title="Iris - KNN (k=3)")
23
24 # Étude de k (1 à 15)
25 ks = list(range(1, 16))
26 accs = []
27 for k in ks:
28     clf = KNeighborsClassifier(n_neighbors=k)
```

```

29     clf.fit(X_train, y_train)
30     accs.append(accuracy_score(y_test, clf.predict(X_test)))
31
32 plt.plot(ks, accs, marker="o")
33 plt.xlabel("k (nombre de voisins)")
34 plt.ylabel("Accuracy test")
35 plt.title("Iris - Impact de k sur la performance KNN")
36 plt.grid(True, alpha=0.3)
37 plt.show()

```

Partie B — Régression Logistique sur *Breast Cancer*

```

1  from sklearn.datasets import load_breast_cancer
2  from sklearn.preprocessing import StandardScaler
3  from sklearn.linear_model import LogisticRegression
4  from sklearn.pipeline import Pipeline
5
6  # Charger les données
7  bc = load_breast_cancer(as_frame=True)
8  X, y = bc.data, bc.target
9  target_names = bc.target_names
10
11 # Split train/test
12 X_train, X_test, y_train, y_test = train_test_split(
13     X, y, test_size=0.30, random_state=42, stratify=y
14 )
15
16 # Pipeline : StandardScaler + LogisticRegression
17 pipe_lr = Pipeline(steps=[
18     ("scaler", StandardScaler()),
19     ("lr", LogisticRegression(max_iter=500))
20 ])
21 pipe_lr.fit(X_train, y_train)
22
23 # Évaluation
24 y_pred = pipe_lr.predict(X_test)
25 print("Accuracy:", accuracy_score(y_test, y_pred))
26 print("Classification report:\n", classification_report(y_test, y_pred,
27     target_names=target_names))
28 plot_conf_mat(y_test, y_pred, labels=target_names, title="Breast Cancer
29     - Régression Logistique")
30
31 # Courbe ROC / AUC
32 y_proba = pipe_lr.predict_proba(X_test)[:, 1]
33 fpr, tpr, thresholds = roc_curve(y_test, y_proba)
34 roc_auc = auc(fpr, tpr)
35
36 plt.plot(fpr, tpr, lw=2, label=f"AUC = {roc_auc:.3f}")
37 plt.plot([0, 1], [0, 1], linestyle="--", lw=1)
38 plt.xlabel("Faux positifs (FPR)")
39 plt.ylabel("Vrais positifs (TPR)")
40 plt.title("Breast Cancer - Courbe ROC (Régression Logistique)")
41 plt.legend(loc="lower right")
42 plt.show()

```

Partie C — Arbre de Décision sur *Wine*

```
1 from sklearn.datasets import load_wine
2 from sklearn.tree import DecisionTreeClassifier, plot_tree
3
4 # Charger les données
5 wine = load_wine(as_frame=True)
6 X, y = wine.data, wine.target
7 labels = wine.target_names
8
9 # Split train/test
10 X_train, X_test, y_train, y_test = train_test_split(
11     X, y, test_size=0.30, random_state=42, stratify=y
12 )
13
14 # Arbre profondeur max = 3
15 tree_clf = DecisionTreeClassifier(max_depth=3, criterion="gini",
16     random_state=42)
17 tree_clf.fit(X_train, y_train)
18
19 y_pred = tree_clf.predict(X_test)
20 print(f"Accuracy (depth=3, gini): {accuracy_score(y_test, y_pred):.3f}")
21 plot_conf_mat(y_test, y_pred, labels=labels, title="Wine - Arbre de Dé
22     cision (gini, depth=3)")
23
24 # Visualisation
25 plt.figure(figsize=(10, 6))
26 plot_tree(tree_clf, feature_names=X.columns, class_names=labels, filled=
27     True, rounded=True)
28 plt.title("Wine - Arbre de Décision (gini, depth=3)")
29 plt.show()
```

À retenir

Points clés

- **KNN** : simple, efficace, mais sensible au choix de k et à l'échelle.
- **Régression Logistique** : interprétable, robuste sur des problèmes binaires.
- **Arbre de Décision** : lisible, puissant mais sujet au surapprentissage.