



IFT712 – Techniques d'apprentissage

Projet de Session : Classification de Feuilles d'Arbres

Équipe :

Antoine Debin deba0971

Antoine Sabatier saba5483

gabriel Généreux geng0394

Lien GitHub : <https://github.com/TonUser/Leaf-Classification>

26 novembre 2025

Table des matières

1	Introduction	3
2	Description des Données	3
3	Méthodologie	3
3.1	Protocole de Validation (Hold-out & Cross-Validation)	3
3.2	Prétraitement (Preprocessing)	3
3.3	Recherche d'Hyper-paramètres	4
4	Architecture et Design Logiciel	4
5	Résultats et Analyse	4
5.1	Comparaison des Performances	4
5.2	Analyse Critique	5
5.3	Matrice de Confusion	5
6	Gestion de Projet	6
7	Conclusion	6

1 Introduction

Dans le cadre du cours IFT712, ce projet vise à concevoir et comparer des systèmes de classification automatique pour identifier des espèces végétales. Nous nous appuyons sur le jeu de données *Leaf Classification* issu de la plateforme Kaggle.

L'objectif principal est de mettre en œuvre une démarche scientifique rigoureuse pour comparer six algorithmes d'apprentissage supervisé (KNN, SVM, Random Forest, MLP, Gradient Boosting et Régression Logistique). Une attention particulière est portée à la méthodologie de validation pour éviter le sur-apprentissage (*overfitting*), un risque majeur compte tenu du faible nombre d'exemples par classe (environ 10 images pour 99 espèces).

2 Description des Données

Le jeu de données est constitué de feuilles préalablement traitées. Contrairement à une approche de vision par ordinateur brute (CNN sur pixels), nous travaillons sur des données tabulaires.

- **Volume** : 990 échantillons d'entraînement.
- **Classes** : 99 espèces d'arbres différentes.
- **Caractéristiques (Features)** : 192 variables numériques extraites mathématiquement décrivant :
 - La Forme (*Shape*)
 - La Marge / le contour (*Margin*)
 - La Texture (*Texture*)

3 Méthodologie

Cette section détaille le protocole expérimental mis en place pour garantir la fiabilité des résultats.

3.1 Protocole de Validation (Hold-out & Cross-Validation)

Pour respecter les bonnes pratiques scientifiques et éviter toute *data leakage* :

1. **Séparation Train / Test** : Avant tout traitement, nous avons isolé **20%** des données (198 échantillons) pour constituer un ensemble de test final (Hold-out set). Ce jeu de données n'a jamais été vu par les modèles durant la phase d'entraînement.
2. **Validation Croisée Stratifiée** : Sur les 80% restants, nous avons utilisé une *Stratified K-Fold Cross-Validation* (k=5). L'aspect stratifié est crucial ici pour assurer que chaque pli de validation contienne une proportion représentative des 99 classes.

3.2 Prétraitement (Preprocessing)

Nous avons encapsulé le prétraitement dans un Pipeline Scikit-Learn.

- **StandardScaler** : Toutes les données ont été normalisées (moyenne = 0, variance = 1). Cette étape est indispensable pour la convergence du MLP et pour le calcul de distance du KNN et du SVM.
- **LabelEncoder** : Encodage des noms d'espèces en entiers.

3.3 Recherche d'Hyper-paramètres

Chaque modèle a été optimisé via un `GridSearchCV` testant diverses combinaisons (ex : noyau *rbf* vs *linear* pour le SVM, nombre de neurones pour le MLP).

4 Architecture et Design Logiciel

Le projet respecte une structure modulaire professionnelle (Package Python) séparant clairement le code, les données et les tests.

- `src/data_loader.py` : Classe responsable du chargement robuste des CSV.
- `src/model_trainer.py` : Classe gérant la logique d'entraînement et les Pipelines.
- `src/models.py` : Centralisation de la configuration des modèles (Pattern *Factory*).

Cette approche permet une grande extensibilité : ajouter un nouveau modèle ne nécessite de modifier qu'un seul fichier de configuration sans toucher à la logique d'entraînement.

5 Résultats et Analyse

5.1 Comparaison des Performances

Le tableau ci-dessous résume les performances obtenues après optimisation. Le **Score Test** est la métrique la plus importante car calculée sur des données inconnues.

Modèle	Score CV (Train)	Score Test (Réal)	Meilleurs Paramètres
MLP (Réseau de Neurones)	0.977	0.995	tanh, (100,)
Régression Logistique	0.986	0.990	C=1, solver=lbfgs
SVM	0.983	0.990	Kernel Linear, C=0.1
Random Forest	0.980	0.985	100 arbres
KNN	0.968	0.970	k=5, Manhattan
Gradient Boosting	0.607	0.611	lr=0.1, depth=3

TABLE 1 – Résultats finaux (Classés par performance sur le Test Set)

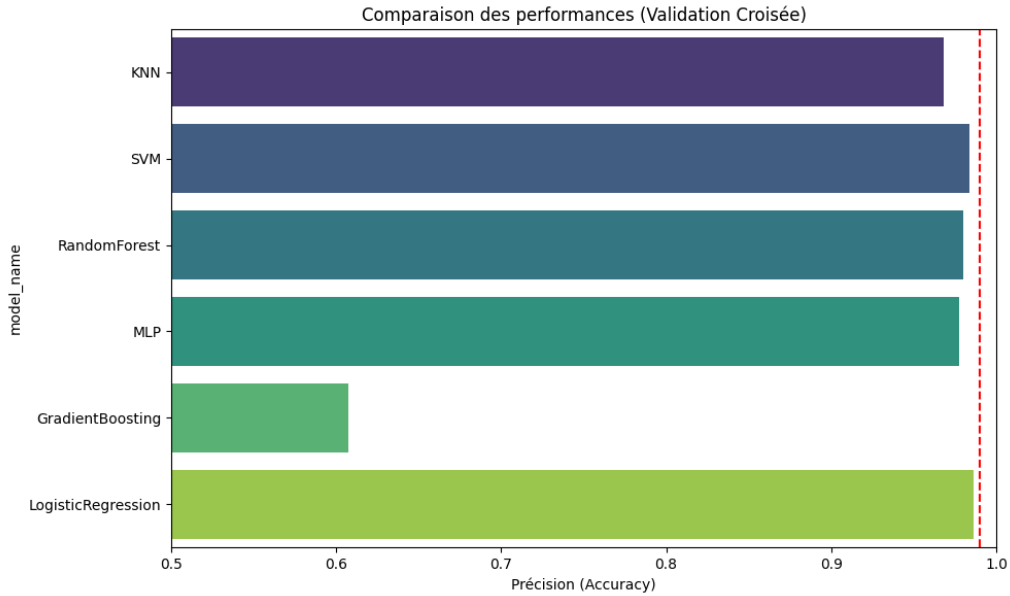


FIGURE 1 – Comparaison graphique des performances (Accuracy)

5.2 Analyse Critique

1. **La domination du MLP (99.5%)** : Le Perceptron Multicouche s'impose comme le meilleur modèle. Avec une seule couche cachée de 100 neurones et une fonction d'activation *tanh*, il a su capturer parfaitement les subtilités géométriques des feuilles.
2. **Linéarité du problème** : Les excellents scores du SVM (noyau linéaire) et de la Régression Logistique (99%) indiquent que les classes sont **linéairement séparables** dans l'espace des 192 caractéristiques. Les features extraites (Marge, Forme, Texture) sont donc de très haute qualité.
3. **L'échec du Gradient Boosting (61%)** : Ce résultat, bien qu'inférieur, est très instructif. Le Boosting est un algorithme complexe qui requiert un grand volume de données pour généraliser. Avec seulement ≈ 8 images par classe dans le set d'entraînement, le modèle a souffert d'un sur-apprentissage massif, n'arrivant pas à converger vers une solution généralisable.

5.3 Matrice de Confusion

La matrice de confusion du modèle champion (MLP) montre une diagonale quasi-parfaite, confirmant la robustesse de la solution.

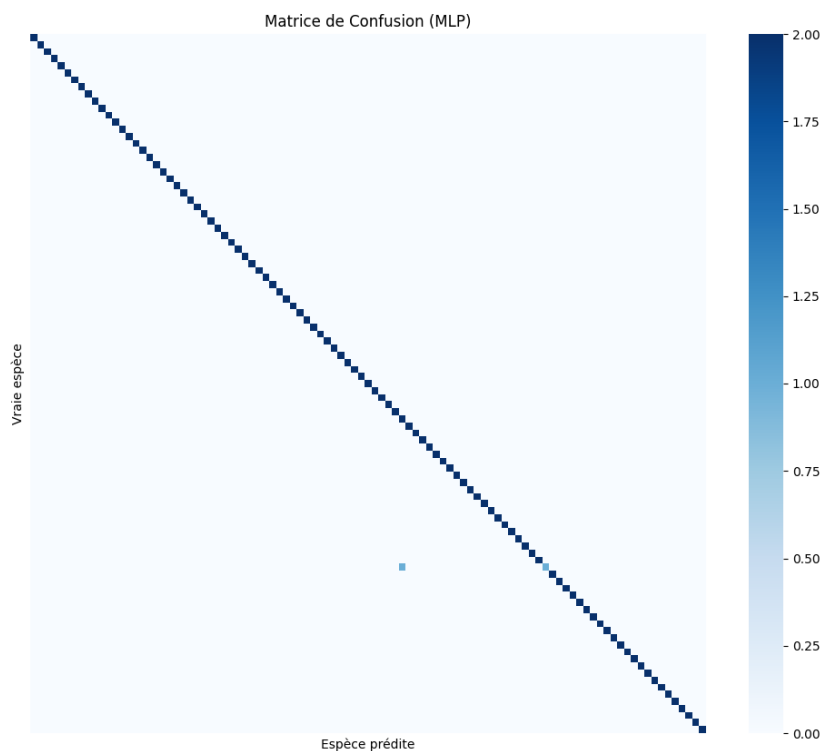


FIGURE 2 – Matrice de Confusion du modèle MLP sur l'ensemble de Test

6 Gestion de Projet

Le développement a suivi les bonnes pratiques logicielles :

- Utilisation de **Git** et **GitHub** pour le versionnage.
- Flux de travail basé sur les branches (*Feature Branch Workflow*).
- Exclusion des fichiers de données et binaires via `.gitignore`.
- Reproductibilité assurée par un fichier `requirements.txt`.

7 Conclusion

Ce projet a permis de développer un classifieur de feuilles d'arbres atteignant une précision de **99.5%** sur des données inédites. L'analyse comparative a mis en évidence que pour ce type de données tabulaires à haute dimensionnalité mais faible volumétrie, les approches comme le MLP ou le SVM sont nettement supérieures aux méthodes d'ensemble complexes comme le Gradient Boosting. La structure modulaire du code et la rigueur du protocole de validation garantissent la fiabilité scientifique de ces conclusions.