

SY09

TD/TP — Arbres de décision

numpy=1.26.4; seaborn=0.13.2; matplotlib=3.8.3; pandas=2.2.0

Une implémentation des arbres de décision est disponible dans `scikit-learn`. Il faut importer la classe `DecisionTreeClassifier`

```
from sklearn.tree import DecisionTreeClassifier
```

Les arguments intéressants lors de l'instanciation de la classe sont les suivants :

- `criterion` : le critère utilisé pour évaluer la qualité d'une séparation,
- `max_leaf_nodes` : le nombre maximum de feuilles autorisée lors de la construction de l'arbre,
- `max_depth` : la profondeur maximale autorisée lors de la construction de l'arbre,
- `max_features` : le nombre de caractéristiques à retenir de manière aléatoire pour la recherche de la meilleure séparation,
- `ccp_alpha` : paramètre λ utilisé pour l'élagage coût-complexité.

1 Arbres de décision

- [1] Retrouver l'arbre construit dans le polycopié de cours en utilisant le jeu de données `data/toy2.csv`. Pour visualiser l'arbre construit, on pourra utiliser la fonction `plot_tree` avec

```
from sklearn.tree import plot_tree
```

- [2] Retrouver le critère de Gini du nœud racine (sans séparation).
- [3] Montrer que le gain tel que vu en cours dû à la première séparation et avec le critère de Gini, vaut $\simeq 1.93773$.

Les effectifs et les critères de Gini de tous les nœuds sont accessibles via l'attribut `tree_`, en spécifiant l'indice du nœud :

```
<gini> = cls.tree_.impurity[<indice d'un nœud>]  
<effectif nœud> = cls.tree_.n_node_samples[<indice d'un nœud>]
```

L'indice du nœud racine est 0, et les indices des autres nœuds peuvent être calculés comme suit :

```
<indice du nœud fils gauche> = cls.tree_.children_left[<indice d'un nœud>]  
<indice du nœud fils droit> = cls.tree_.children_right[<indice d'un nœud>]
```

2 Modèles basés sur les arbres

2.1 *Bagging*

Dans cette section, nous allons recréer manuellement les différentes étapes de la technique du *bagging* pour les arbres de décision. La première étape consiste à créer plusieurs jeux de données avec la technique du *bootstrap*.

- 4 Créer trois jeux de données échantillonnés par *bootstrap* grâce à la fonction `resample` disponible avec

```
from sklearn.utils import resample
```

On utilisera le jeu de données `Synth1-2000.csv` du TP07.

- 5 Apprendre un arbre de décision limité à 50 feuilles maximum sur chaque réplication du jeu de données. Afficher leurs régions de décision respectives.
- 6 En utilisant les trois modèles précédents, compléter la fonction suivante qui réalise l'agrégation des trois fonctions de décision. On pourra utiliser la fonction `np.where`.

```
def aggregating(X):  
    y1 = ...  
    y2 = ...  
    y3 = ...  
  
    return ...
```

Utiliser ensuite cette fonction avec `add_decision_boundary` en spécifiant `model_classes` pour afficher les frontières de décision du modèle agrégé.

- 7 Réaliser le même travail en utilisant cette fois la classe `BaggingClassifier`.

2.2 Forêt aléatoire

Une implémentation des forêts aléatoires est disponible dans `scikit-learn`. Il faut importer la classe `RandomForestClassifier`

```
from sklearn.ensemble import RandomForestClassifier
```

En plus des arguments des arbres de décision, on trouve l'argument `n_estimators` qui fixe le nombre d'arbres utilisés dans la forêt. De plus, l'argument `max_features` est cette fois fixé par défaut à \sqrt{p} (avec p le nombre de caractéristiques).

- 8 Apprendre une forêt aléatoire sur le jeu de données précédent et visualiser le résultat.

2.3 Élagage coût-complexité et validation croisée

Dans cette section, on se propose d'implémenter la recherche du paramètre λ optimal par validation croisée telle que décrite aux pages 127 à 128 du polycopié de cours.

- 9 On commence par déterminer les valeurs λ_k pour la séquence $(\mathcal{A}_k, \lambda_k)$ d'arbres appris sur l'ensemble d'apprentissage total \mathcal{T} .

La série des λ_k est disponible grâce à la fonction `cost_complexity_pruning_path`.

- 10 Calculer les moyennes géométriques $\bar{\lambda}_k = \sqrt{\lambda_k \lambda_{k+1}}$.

- 11 Compléter le générateur suivant qui génère les erreurs $\hat{\varepsilon}_v(\mathcal{A}^v(\bar{\lambda}_k))$ pour chaque k et chaque v .

```
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import KFold
from sklearn.metrics import accuracy_score
from sklearn.utils import check_X_y

def decision_tree_cross_validation_accuracies(X, y, n_folds, lambdas):
    X, y = check_X_y(X, y)

    # Création d'un objet `KFold` pour la validation croisée
    kf = ...

    for train_index, val_index in kf:
        # Création de `X_train`, `y_train`, `X_val` et `y_val`
        X_train = ...
        y_train = ...
        X_val = ...
        y_val = ...

        for k, lmb in enumerate(lambdas):
            # Création d'un arbre avec un coefficient coût-complexité
            # égal à `lmb`
            clf = ...

            # Apprentissage sur l'ensemble d'apprentissage et calcul
            # du taux de bonne classification sur l'ensemble de
            # validation
            ...
            yield k, lmb, acc
```

12 Créer le jeu de données issu du générateur précédent. Agréger les erreurs de tous les plis ensemble et afficher les erreurs de validation en fonction des $\bar{\lambda}_k$.

13 En déduire le sous-arbre optimal obtenu par cette procédure et afficher sa frontière de décision.



2.4 Comparaison de la diversité des arbres

Dans cette section, on se propose d'illustrer le bénéfice du paramètre `max_features` en terme de diversité des arbres appris.

Pour évaluer la diversité des arbres, on choisit de comparer les ensembles de caractéristiques sélectionnées pour effectuer des séparations jusqu'à une certaine profondeur.

On pourra utiliser la fonction suivante qui génère les indices des caractéristiques utilisées dans les séparations jusqu'à une certaine profondeur.

```
def features_depth(model, depth, acc=False):
    """Génère les indices des caractéristiques utilisées dans un arbre.

    L'argument `model` est l'arbre. Les indices sont générés
    uniquement à la profondeur `depth` sauf si `acc` est vrai. Dans ce
    cas, toutes les caractéristiques jusqu'à la profondeur `depth`
    sont générées.

    """

    tree = model.tree_
```

```
def gen_id(i, depth):  
    if tree.feature[i] >= 0:  
        if acc or depth == 0:  
            yield tree.feature[i]  
    if depth != 0:  
        yield from gen_id(tree.children_left[i], depth - 1)  
        yield from gen_id(tree.children_right[i], depth - 1)  
  
yield from gen_id(0, depth)
```

14 Afficher la distribution des caractéristiques retenues à la profondeur 0, 1, 2 pour une forêt de 100 arbres en faisant varier `max_features`. Que peut-on constater ? Quelle est la distribution lorsque `max_features` vaut 1 ?

On utilisera le jeu de données `spams` disponible avec le TP09.