

Rapport TP2 Arbres

Mathieu Le Séac'h

Classification avec les arbres

```
import os
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from matplotlib import rc

from sklearn import tree, datasets, model_selection

from itertools import product
from collections import defaultdict

from tp_arbres_source import (rand_gauss, rand_bi_gauss, rand_tri_gauss,
                              rand_checkers, rand_clown,
                              plot_2d, frontiere)
```

Q1)

Dans le cadre d'une régression, on peut utiliser TODO comme mesure d'homogénéité.

Q2)

On cherche à tracer les courbes donnant le pourcentage d'erreurs commises en fonction de la profondeur maximale de l'arbre pour

```
np.random.seed(10)

criteria = ["gini", "entropy"]
```

```

depths = range(1, 14)

n = 456
data = rand_checkers(n//4, n//4, n//4, n//4)
X_fit = data[:, :2]
y_fit = data[:, 2]

def fit_all_classifiers(criteria, depths, X_fit, y_fit):
    classifiers = defaultdict(list)

    for (criterion, max_depth) in product(criteria, depths):
        classifier = tree.DecisionTreeClassifier(
            criterion=criterion,
            max_depth=max_depth,
        ).fit(
            X_fit, y_fit
        )

        classifiers[criterion].append(classifier)

    return classifiers

classifiers = fit_all_classifiers(criteria, depths, X_fit, y_fit)

def compute_all_scores(classifiers, X_test, y_test):
    scores = dict()

    for criterion in classifiers.keys():
        scores[criterion] = list(map(
            lambda classifier: classifier.score(X_test, y_test),
            classifiers[criterion]
        ))

    return scores

scores = compute_all_scores(classifiers, X_fit, y_fit)

plt.figure()

```

```

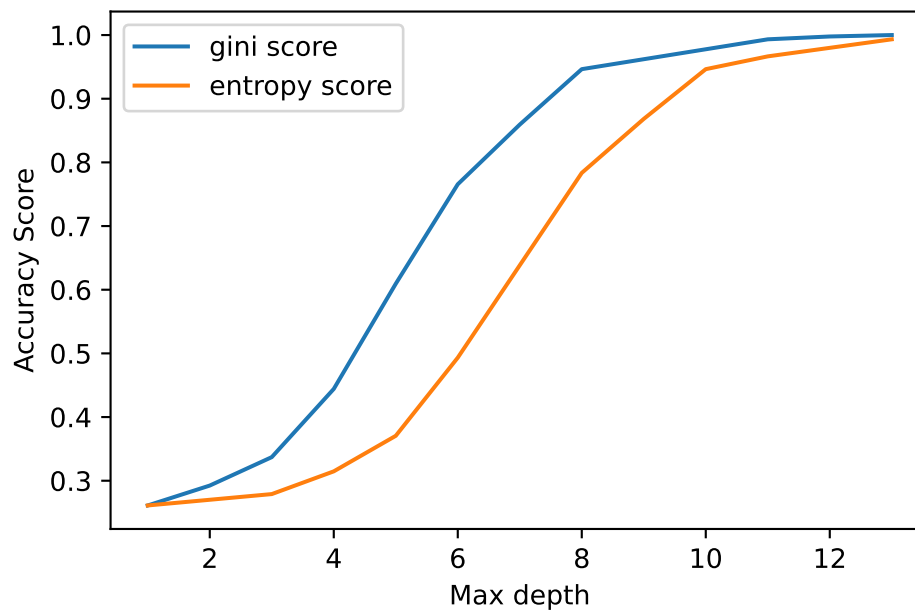
for criterion in scores.keys():
    plt.plot(depths, scores[criterion], label=f"{criterion} score")
    print(f"Scores with {criterion} criterion: {scores[criterion]}")

plt.xlabel('Max depth')
plt.ylabel('Accuracy Score')
plt.legend()
plt.draw()

```

Scores with gini criterion: [0.2611607142857143, 0.2924107142857143, 0.33705357142857145, 0.44285714285714285, 0.5714285714285714, 0.7142857142857143, 0.8571428571428571, 0.9285714285714286, 0.9571428571428571, 0.9714285714285714, 0.9857142857142857, 0.9928571428571428, 0.9957142857142857, 0.9971428571428571, 0.9985714285714286, 0.9992857142857143, 0.9995714285714286, 0.9997142857142857, 0.9998571428571428, 0.9999285714285714, 0.9999571428571428, 0.9999714285714286, 0.9999857142857143, 0.9999928571428571, 0.9999957142857143, 0.9999971428571428, 0.9999985714285714, 0.9999992857142857, 0.9999995714285714, 0.9999997142857143, 0.9999998571428571, 0.9999999285714286, 0.9999999571428571, 0.9999999714285714, 0.9999999857142857, 0.9999999928571428, 0.9999999957142857, 0.9999999971428571, 0.9999999985714286, 0.9999999992857143, 0.9999999995714286, 0.9999999997142857, 0.9999999998571428, 0.9999999999285714, 0.9999999999571428, 0.9999999999714286, 0.9999999999857143, 0.9999999999928571, 0.9999999999957143, 0.9999999999971428, 0.9999999999985714, 0.9999999999992857, 0.9999999999995714, 0.9999999999997143, 0.9999999999998571, 0.9999999999999286, 0.9999999999999571, 0.9999999999999714, 0.9999999999999857, 0.9999999999999929, 0.9999999999999957, 0.9999999999999971, 0.9999999999999986, 0.9999999999999993, 0.9999999999999996, 0.9999999999999998, 0.9999999999999999, 1.0]

Scores with entropy criterion: [0.2611607142857143, 0.2700892857142857, 0.27901785714285715, 0.33705357142857145, 0.44285714285714285, 0.5714285714285714, 0.7142857142857143, 0.8571428571428571, 0.9285714285714286, 0.9571428571428571, 0.9714285714285714, 0.9857142857142857, 0.9928571428571428, 0.9957142857142857, 0.9971428571428571, 0.9985714285714286, 0.9992857142857143, 0.9995714285714286, 0.9997142857142857, 0.9998571428571428, 0.9999285714285714, 0.9999571428571428, 0.9999714285714286, 0.9999857142857143, 0.9999928571428571, 0.9999957142857143, 0.9999971428571428, 0.9999985714285714, 0.9999992857142857, 0.9999995714285714, 0.9999997142857143, 0.9999998571428571, 0.9999999285714286, 0.9999999571428571, 0.9999999714285714, 0.9999999857142857, 0.9999999928571428, 0.9999999957142857, 0.9999999971428571, 0.9999999985714286, 0.9999999992857143, 0.9999999995714286, 0.9999999997142857, 0.9999999998571428, 0.9999999999285714, 0.9999999999571428, 0.9999999999714286, 0.9999999999857143, 0.9999999999928571, 0.9999999999957143, 0.9999999999971428, 0.9999999999985714, 0.9999999999992857, 0.9999999999995714, 0.9999999999997143, 0.9999999999998571, 0.9999999999999286, 0.9999999999999571, 0.9999999999999714, 0.9999999999999857, 0.9999999999999929, 0.9999999999999957, 0.9999999999999971, 0.9999999999999986, 0.9999999999999993, 0.9999999999999996, 0.9999999999999998, 0.9999999999999999, 1.0]



Comme on s'y attendait, la précision des arbres étant testés sur les données d'apprentissage, plus le modèle a de paramètres (ici de profondeur), mieux il va connaître les données d'apprentissage.

Q3)

Comme expliqué à la question précédente, plus la profondeur sera grande, plus l'erreur sera petite. Ainsi avec les profondeurs testées (de 1 à 13), l'arbre de décision qui minimise l'erreur est celui de profondeur 13, cependant dans notre cas le classificateur atteint déjà une précision

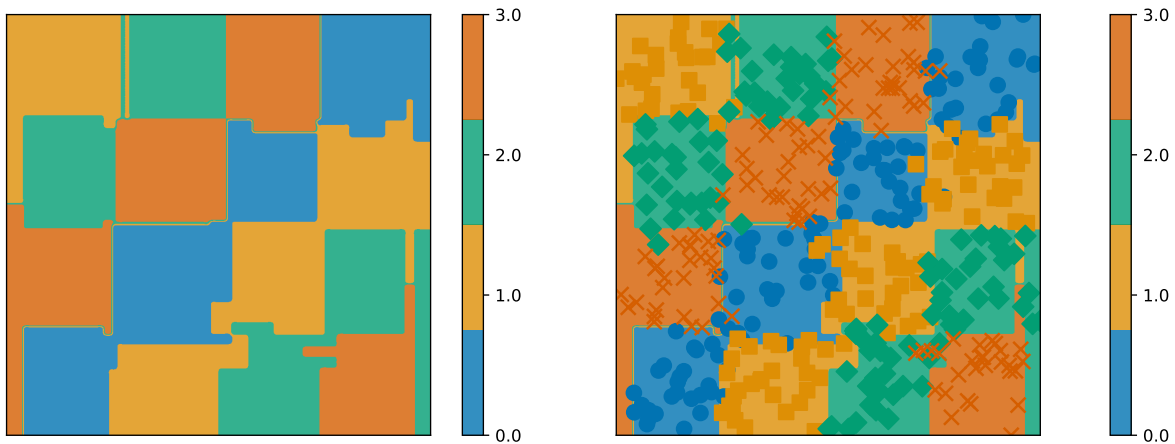
parfaite à partir d'une profondeur de 16. Ainsi avec le critère d'entropie et une profondeur de 13 on obtient la classification suivante:

```
best_depth = np.argmax(scores["entropy"]) + 1
assert(best_depth == 13)

best_classifier = classifiers["entropy"][best_depth - 1]

plt.figure(figsize=(15, 10))
plt.subplot(2, 2, 1)
frontiere(
    lambda x: best_classifier.predict(x.reshape((1, -1))),
    X_fit, y_fit,
    step=100, samples=False
)

plt.subplot(2, 1, 1)
plot_2d(X_fit, y_fit)
frontiere(
    lambda x: best_classifier.predict(x.reshape((1, -1))),
    X_fit, y_fit,
    step=100, samples=False
)
```



Dans l'ensemble le damier ressemble bien à celui des données, cependant on remarque certains rectangles

Q4)

TODO: exporter l'arbre

Q5)

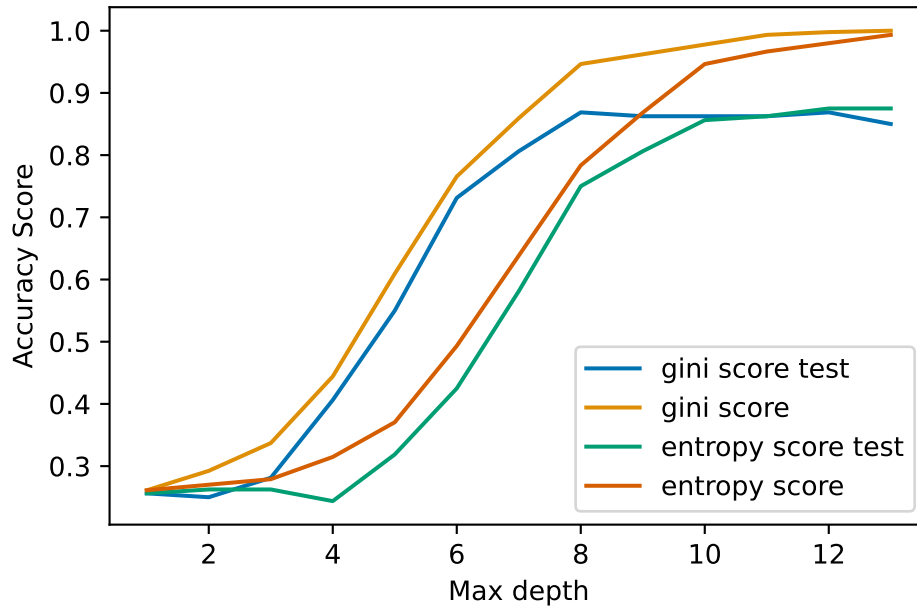
```
n_test = 160
data_test = rand_checkers(n_test//4, n_test//4, n_test//4, n_test//4)
X_test = data_test[:, :2]
y_test = data_test[:, 2]

scores_test = compute_all_scores(classifiers, X_test, y_test)

plt.figure()
for criterion in criterions:
    plt.plot(depths, scores_test[criterion], label=f"{criterion} score test")
    print(f"Scores on test data with {criterion} criterion: {scores_test[criterion]}")
    plt.plot(depths, scores[criterion], label=f"{criterion} score ")
    print(f"Scores with {criterion} criterion: {scores[criterion]}")

plt.xlabel('Max depth')
plt.ylabel('Accuracy Score')
plt.legend()
plt.draw()
```

```
Scores on test data with gini criterion: [0.25625, 0.25, 0.28125, 0.40625, 0.55, 0.73125, 0.8]
Scores with gini criterion: [0.2611607142857143, 0.2924107142857143, 0.33705357142857145, 0.4]
Scores on test data with entropy criterion: [0.25625, 0.2625, 0.2625, 0.24375, 0.31875, 0.42]
Scores with entropy criterion: [0.2611607142857143, 0.2700892857142857, 0.27901785714285715,
```



Q6)

```

digits = datasets.load_digits()
X = digits.data
y = digits.target

X_fit, X_test, y_fit, y_test = model_selection.train_test_split(X, y, test_size=0.2, random_state=42)

criteria = ["gini", "entropy"]
depths = range(1, 14)

classifiers = fit_all_classifiers(criteria, depths, X_fit, y_fit)

scores = compute_all_scores(classifiers, X_fit, y_fit)
scores_test = compute_all_scores(classifiers, X_test, y_test)

plt.figure()

for criterion in scores.keys():
    plt.plot(depths, scores_test[criterion], label=f"{criterion} score test")
    print(f"Scores on test data with {criterion} criterion: {scores_test[criterion]}")
    plt.plot(depths, scores[criterion], label=f"{criterion} score ")

```

