

Machine Learning

Séance 10

L'objectif de cette session est de (i) rappeler le fonctionnement des arbres de décision et (ii) d'appliquer nos connaissances sur python.

1 Importation, description et visualisation des données

1. Importez la base de données

— Première possibilité (la plus simple)

```
1 from sklearn.datasets import load_iris
2 iris = load_iris()
3 X = iris.data[:, 2:] # petal length and width
4 y = iris.target
```

— Deuxième possibilité (moins simple)

```
1 import pandas
2 X = pandas.read_excel("/Users/data/data_X.xlsx")
3 y = pandas.read_excel("/Users/data/data_y.xlsx")
4 X = X.to_numpy()
5 y = y.to_numpy()
```

2. Décrivez la structure des données :

— Pour la base X : `type(X)`, `X[:10]`, `X.shape`

```
1 print(type(X))
2 print(X[:10])
3 print(X.shape)
```

— Quelle est la structure de la base X? Combien possède-t-elle de lignes? Combien de colonnes? Que représentent les colonnes?

3. Visualisez les données :

```
1
2 import matplotlib.pyplot as plt
3 plt.scatter(X[:, 0], X[:, 1], c=y, cmap=plt.cm.Set1, edgecolor="k")
4 plt.xlabel("Petal_length")
5 plt.ylabel("Petal_width")
```

4. Que représente Y?

2 Arbre de décision

5. Implémentez un arbre de décision

```
1 from sklearn import tree
2 tree_clf = DecisionTreeClassifier(max_depth=2)
3 tree_clf.fit(X, y)
4 predictions = tree_clf.predict(X)
5 predictions
6 print(f"%Classification_Correct: {tree_clf.score(X, y):.3f}")
```

6. Visualisez les deux arbres de décision suivants. Quelle est la différence principale entre ces deux arbres?

```
1 from sklearn import tree
2 clf = tree.DecisionTreeClassifier(max_depth=2)
3 clf.fit(X, y)
4 plot = tree.plot_tree(clf, feature_names=iris.feature_names[2:], class_names=iris.target_names,
5 proportion = True, rounded=True, filled=True)
6
```

```

7 from sklearn import tree
8 clf = tree.DecisionTreeClassifier(max_depth=3)
9 clf.fit(X, y)
10 tree.plot_tree(clf, feature_names=iris.feature_names[2:], class_names=iris.target_names,
11 proportion = True, rounded=True, filled=True)

```

7. En utilisant un arbre de décision, prédir le type de fleur ayant des pétales de 5cm de long et de 1,5cm de largeur :

```

1 tree_clf = tree.DecisionTreeClassifier(max_depth=2)
2 tree_clf.fit(X, y)
3 tree_clf.predict_proba([[5, 1.5]])
4 tree_clf.predict([[5, 1.5]])

```

8. Questions sur les arbres de décision :

- (a) Vrai ou faux : les arbres de décision désignent une méthode linéaire et non-paramétrique
- (b) Vrai ou faux : le fonctionnement des arbres de décision revient à segmenter l'espace de décision
- (c) Vrai ou faux : pour obtenir les meilleures performances, il vaut toujours mieux utiliser des arbres plus profonds
- (d) Quel est le problème principal des arbres de décision simples ?
- (e) A quoi correspond le *Tree Pruning* ?

3 Préparation des données

9. Divisez la base de données en deux sous-bases : une base d'apprentissage et une base de test — Utilisez le code suivant :

```

1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=2)

```

- Combien d'observations contient la base de test ? Quelle proportion de la base de données cela représente ?

10. Vrai ou faux :

- (a) La base d'apprentissage est utilisée pour tester la performance du modèle
- (b) Il suffit de classer aléatoirement les observations avant la division de la base de données en deux pour obtenir toujours les mêmes performances pour un modèle donné
- (c) Il faut toujours mettre 20% des observations dans la base de test et 80% dans la base d'apprentissage

4 Comparaison de plusieurs modèles

4.1 Entraînement

11. Utilisez le code suivant :

```

1 from sklearn.ensemble import BaggingClassifier
2
3 clf = BaggingClassifier(base_estimator=tree.DecisionTreeClassifier(), n_estimators = 10, random_state=10)
4 clf.fit(X_train, y_train)
5 print(f'Test : {clf.score(X_test, y_test):.3f}')
6 print(f'Train : {clf.score(X_train, y_train):.3f}')
7
8 from sklearn.ensemble import RandomForestClassifier
9
10 RF_Model = RandomForestClassifier()
11 RF_Model = RandomForestClassifier(random_state=10, n_estimators = 10)

```

```

12 RF_Model.fit(X_train, y_train)
13 print(f'Test : {RF_Model.score(X_test, y_test):.3f}')
14 print(f'Train : {RF_Model.score(X_train, y_train):.3f}')
15
16
17 from sklearn.ensemble import GradientBoostingClassifier
18
19 clf = GradientBoostingClassifier(n_estimators=10, learning_rate=0.1,
20                                max_depth=1, random_state=0).fit(X_train, y_train)
21
22 clf.fit(X_train, y_train)
23 print(f'Test : {clf.score(X_test, y_test):.3f}')
24 print(f'Train : {clf.score(X_train, y_train):.3f}')

```

12. Questions sur le code :

- (a) Quel(s) modèle(s) a-t-on utilisé ?
- (b) Quel modèle donne la meilleure performance ?
- (c) Quels sont les caractéristiques de chaque modèle ?

13. Questions sur le bagging :

- (a) Quelle est la différence entre un arbre de décision simple et le bagging ?
- (b) Vrai ou faux : le bagging n'utilise pas le bootstrap
- (c) Vrai ou faux : le bagging utilise un sous-ensemble des prédicteurs sélectionné aléatoirement
- (d) Vrai ou faux : le bagging fait "pousser" des arbres de manière séquentielle

14. Questions sur le random forest :

- (a) Quelle est la différence entre le bagging et le random forest ?
- (b) Vrai ou faux : le random forest n'utilise pas le bootstrap
- (c) Vrai ou faux : le random forest n'est jamais équivalent au bagging
- (d) Vrai ou faux : le random forest utilise un sous-ensemble des prédicteurs sélectionné aléatoirement
- (e) Vrai ou faux : le random forest fait "pousser" des arbres de manière séquentielle

15. Questions sur le boosting :

- (a) Quelle est la différence entre le boosting et le random forest ?
- (b) Vrai ou faux : le boosting n'utilise pas le bootstrap
- (c) Vrai ou faux : le boosting fait "pousser" des arbres de manière séquentielle
- (d) Vrai ou faux : le boosting utilise un sous-ensemble des prédicteurs sélectionné aléatoirement
- (e) Vrai ou faux : la performance du boosting est proportionnelle à la profondeur des arbres utilisés
- (f) Dans le cadre du boosting, peut-on dire que chaque arbre essaie de corriger les faiblesses du précédent ?

4.2 Hyperparameter Tuning

16. Que désigne l'expression *Hyperparameter Tuning* ?

4.2.1 Random Forest

17. Lancez le code suivant :

```
1 n_arbre = [1,2,3,4,5]
2 train_results = []
3 test_results = []
4 for i in n_arbre:
5     rf = RF_Model = RandomForestClassifier(bootstrap=True,
6                                           n_estimators=i,
7                                           random_state=100)
8     rf.fit(X_train, y_train)
9     train_pred = rf.predict(X_train)
10    test_r = rf.score(X_test, y_test)
11    test_results.append(test_r)
12    train_r = rf.score(X_train, y_train)
13    train_results.append(train_r)
14 from matplotlib.legend_handler import HandlerLine2D
15 line1, = plt.plot(n_arbre, train_results, label='Erreur_d_apprentissage')
16 line2, = plt.plot(n_arbre, test_results, label='Erreur_de_test')
17 plt.legend(handler_map={line1: HandlerLine2D(numpoints=2)})
18 plt.ylabel('Taux_de_Classification_Correct')
19 plt.xlabel('Nombre_d_arbres')
20 plt.show()
```

- (a) Quel hyperparamètre est étudié ici ?
- (b) Comment varie la performance du modèle en fonction de cet hyperparamètre ?
- (c) Si vous deviez choisir une valeur pour cet hyperparamètre afin d'obtenir la meilleure performance, quelle valeur choisiriez-vous ?

18. Lancez le code suivant :

```
1 max_depths = [1, 2, 3, 4, 5]
2 train_results = []
3 test_results = []
4 for i in max_depths:
5     rf = RF_Model = RandomForestClassifier(bootstrap=True, max_depth=i, random_state=100)
6     rf.fit(X_train, y_train)
7     train_pred = rf.predict(X_train)
8     test_r = rf.score(X_test, y_test)
9     test_results.append(test_r)
10    train_r = rf.score(X_train, y_train)
11    train_results.append(train_r)
12 from matplotlib.legend_handler import HandlerLine2D
13 line1, = plt.plot(max_depths, train_results, label='Erreur_d_apprentissage')
14 line2, = plt.plot(max_depths, test_results, label='Erreur_de_test')
15 plt.legend(handler_map={line1: HandlerLine2D(numpoints=2)})
16 plt.ylabel('Taux_de_Classification_Correct')
17 plt.xlabel('Profondeur_d_un_arbre')
18 plt.show()
```

- (a) Quel hyperparamètre est étudié ici ?
- (b) Comment varie la performance du modèle en fonction de cet hyperparamètre ?
- (c) Si vous deviez choisir une valeur pour cet hyperparamètre afin d'obtenir la meilleure performance, quelle valeur choisiriez-vous ?

4.2.2 Boosting

19. Lancez le code suivant :

```
1 #####
2 from sklearn.ensemble import GradientBoostingClassifier
3 #####
4 n_estimators_n = [1, 2, 5, 10, 20]
5 train_results = []
6 test_results = []
7 for i in n_estimators_n:
8     rf = GradientBoostingClassifier(n_estimators=i, learning_rate=0.1,
9                                     max_depth=1, random_state=0).fit(X_train, y_train)
10    rf.fit(X_train, y_train)
11    train_pred = rf.predict(X_train)
12    test_r = rf.score(X_test, y_test)
13    test_results.append(test_r)
14    train_r = rf.score(X_train, y_train)
15    train_results.append(train_r)
16 from matplotlib.legend_handler import HandlerLine2D
```

```

17 line1, = plt.plot(max_depths, train_results, label='Erreur_d_apprentissage')
18 line2, = plt.plot(max_depths, test_results, label='Erreur_de_test')
19 plt.legend(handler_map={line1: HandlerLine2D(numpoints=2)})
20 plt.ylabel('Taux_de_Classification_Correct')
21 plt.xlabel('Nombre_d_arbres')
22 plt.show()

```

- Quel hyperparamètre est étudié ici ?
- Comment varie la performance du modèle en fonction de cet hyperparamètre ?
- Si vous deviez choisir une valeur pour cet hyperparamètre afin d'obtenir la meilleure performance, quelle valeur choisiriez-vous ?

20. Lancez le code suivant :

```

1 max_depths = [1, 2, 3, 4]
2 train_results = []
3 test_results = []
4 for i in max_depths:
5     rf = GradientBoostingClassifier(n_estimators=20, learning_rate=0.1,
6         max_depth=i, random_state=0).fit(X_train, y_train)
7     rf.fit(X_train, y_train)
8     train_pred = rf.predict(X_train)
9     test_r = rf.score(X_test, y_test)
10    test_results.append(test_r)
11    train_r = rf.score(X_train, y_train)
12    train_results.append(train_r)
13 from matplotlib.legend_handler import HandlerLine2D
14 line1, = plt.plot(max_depths, train_results, label='Erreur_d_apprentissage')
15 line2, = plt.plot(max_depths, test_results, label='Erreur_de_test')
16 plt.legend(handler_map={line1: HandlerLine2D(numpoints=2)})
17 plt.ylabel('Taux_de_Classification_Correct')
18 plt.xlabel('Profondeur_d_un_arbre')
19 fig1 = plt.gcf()
20 fig1.savefig('/Users/maitre/Dropbox/Assas/cours/machine_learning/slides/lecture10/erreur_classification_profondeur_arbre')
21 plt.show()

```

- Quel hyperparamètre est étudié ici ?
- Comment varie la performance du modèle en fonction de cet hyperparamètre ?
- Si vous deviez choisir une valeur pour cet hyperparamètre afin d'obtenir la meilleure performance, quelle valeur choisiriez-vous ?