

Machine Learning

Séance 11

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

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

2. Que représente Y ?
3. Que représente X ?

2 Machines à Vecteur de Support

2.1 Deux classes séparables

4. Visualisez les données :

```
1 X = X_base[:100]
2 y = y_base[:100]
3
4 import matplotlib.pyplot as plt
5 plt.scatter(X[:, 0], X[:, 1], c=y, cmap=plt.cm.Set1, edgecolor="k")
6 plt.xlabel("Petal_length")
7 plt.ylabel("Petal_width")
8 plt.xlim(0, 6)
9 plt.ylim(0, 2)
```

5. Tracez l'hyperplan séparateur

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn import svm
4
5 svc = svm.SVC(kernel='linear', C=1).fit(X, y)
6
7 x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
8 y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
9 xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.02),
10                      np.arange(y_min, y_max, 0.02))
11
12 Z = svc.predict(np.c_[xx.ravel(), yy.ravel()])
13 Z = Z.reshape(xx.shape)
14 plt.contourf(xx, yy, Z, cmap=plt.cm.coolwarm, alpha=0.8)
15 plt.scatter(X[:, 0], X[:, 1], c=y, cmap=plt.cm.Set1)
16 plt.xlabel('Petal_length')
17 plt.ylabel('Petal_width')
18 plt.xlim(0, 6)
19 plt.ylim(0, 2)
20 plt.show()
```

6. Vrai ou faux :

- (a) Dans un espace de p -dimension, un hyperplan est un sous-espace de dimension $p+1$
- (b) Dans un espace de 3 dimensions, un hyperplan est simplement une droite
- (c) Un hyperplan séparateur permet de classifier les observations selon leur localisation par rapport à l'hyperplan

- (d) Dans un espace de 2 dimensions, il existe toujours au moins un hyperplan séparateur
- (e) Dans un espace de 3 dimensions, il peut ne pas exister d'hyperplan séparateur
- (f) Si les données sont parfaitement séparables avec un hyperplan, alors il existe une infinité d'hyperplans séparateurs
- (g) La marge correspond à la distance perpendiculaire entre l'hyperplan et l'observation la plus proche de l'hyperplan

2.2 Deux classes non séparables

7. Visualisez les données

```

1 X = X_base[50:]
2 y = y_base[50:]
3
4 import matplotlib.pyplot as plt
5 plt.scatter(X[:, 0], X[:, 1], c=y, cmap=plt.cm.Set1, edgecolor="k")
6 plt.xlabel("Petal_length")
7 plt.ylabel("Petal_width")

```

8. Soft Margin classifier

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn import svm
4
5 svc1 = svm.SVC(kernel='linear', C=1).fit(X, y)
6 svc2 = svm.SVC(kernel='linear', C=10).fit(X, y)
7 svc3 = svm.SVC(kernel='linear', C=100).fit(X, y)
8 svc4 = svm.SVC(kernel='linear', C=1000).fit(X, y)
9
10 x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
11 y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
12 xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.02),
13                      np.arange(y_min, y_max, 0.02))
14
15 titles = ['C_u=1',
16           'C_u=10',
17           'C_u=100',
18           'C_u=1000']
19
20 for i, clf in enumerate((svc1, svc2, svc3, svc4)):
21     plt.subplot(2, 2, i + 1)
22     plt.subplots_adjust(wspace=0.4, hspace=0.4)
23     Z = clf.predict(np.c_[xx.ravel(), yy.ravel()])
24     Z = Z.reshape(xx.shape)
25     plt.contourf(xx, yy, Z, cmap=plt.cm.coolwarm, alpha=0.8)
26     plt.scatter(X[:, 0], X[:, 1], c=y, cmap=plt.cm.Set1)
27     plt.xlabel('Petal_length')
28     plt.ylabel('Petal_width')
29     plt.xlim(xx.min(), xx.max())
30     plt.ylim(yy.min(), yy.max())
31     plt.xticks(())
32     plt.yticks(())
33     plt.title(titles[i])
34 plt.show()

```

9. Vrai ou faux :

- (a) Dans le cadre d'un problème de classification, on peut toujours utiliser le maximal margin classifier
- (b) En présence d'un hyperplan séparateur, le maximal margin classifier produit toujours des meilleures performances que le soft margin classifier
- (c) La différence principale entre le maximal margin classifier et le soft margin classifier repose sur la non-linéarité de l'hyperplan séparateur
- (d) Le paramètre C détermine le nombre et le degré de violation de la marge
- (e) Le paramètre C est toujours compris entre 0 et 1
- (f) Plus C est grand, plus le biais sera faible sur la base d'apprentissage

2.3 Trois classes

10. Visualisez les données :

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

11. Tracez les plans séparateurs

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn import svm
4
5 svc = svm.SVC(kernel='linear', C=1).fit(X, y)
6
7 x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
8 y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
9 xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.02),
10                      np.arange(y_min, y_max, 0.02))
11
12 Z = svc.predict(np.c_[xx.ravel(), yy.ravel()])
13 Z = Z.reshape(xx.shape)
14 plt.contourf(xx, yy, Z, cmap=plt.cm.coolwarm, alpha=0.8)
15 plt.scatter(X[:, 0], X[:, 1], c=y, cmap=plt.cm.Set1)
16 plt.xlabel('Petal_length')
17 plt.ylabel('Petal_width')
18 plt.show()
```

12. Changez le noyau :

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn import svm
4
5 svc = svm.SVC(kernel='linear', C=1).fit(X, y)
6 rbf_svc = svm.SVC(kernel='rbf', gamma=0.7, C=1).fit(X, y)
7 poly_svc = svm.SVC(kernel='poly', degree=3, C=1).fit(X, y)
8
9 x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
10 y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
11 xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.02),
12                      np.arange(y_min, y_max, 0.02))
13
14 titles = ['SVC_avec_noyau_lineaire',
15          'SVC_avec_noyau_RBF',
16          'SVC_avec_noyau_polynomial_(degre_3)']
17
18 for i, clf in enumerate((svc, rbf_svc, poly_svc)):
19     plt.subplot(2, 2, i + 1)
20     plt.subplots_adjust(wspace=0.4, hspace=0.4)
21     Z = clf.predict(np.c_[xx.ravel(), yy.ravel()])
22     Z = Z.reshape(xx.shape)
23     plt.contourf(xx, yy, Z, cmap=plt.cm.coolwarm, alpha=0.8)
24     plt.scatter(X[:, 0], X[:, 1], c=y, cmap=plt.cm.Set1)
25     plt.xlabel('Petal_length')
26     plt.ylabel('Petal_width')
27     plt.xlim(xx.min(), xx.max())
28     plt.ylim(yy.min(), yy.max())
29     plt.xticks(())
30     plt.yticks(())
31     plt.title(titles[i])
32 plt.show()
```

13. Vrai ou faux :

- (a) En général, l'astuce du noyau consiste à transposer les données dans un autre espace de plus petite dimension dans lequel elles sont linéairement séparables
- (b) L'intuition derrière l'astuce du noyau consiste à transposer un problème non-linéaire en problème linéaire
- (c) Dans le cadre des machines à vecteur de support, utiliser un noyau polynomial amène toujours les meilleures performances
- (d) Dans le cadre des machines à vecteur de support, utiliser un noyau gaussien amène toujours les meilleures performances

- (e) Si on souhaite déterminer le risque de crise cardiaque en fonction des caractéristiques d'un individu, on peut utiliser une machine à vecteur de support
- (f) Les machines à vecteur de support peuvent constituer une alternative à la régression logistique
- (g) Les machines à vecteur de support sont souvent utilisées pour des problèmes d'inférence

2.4 Comparaison de plusieurs modèles

14. Divisez la base de données en deux :

```
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=0.15, random_state=6)
```

15. Comparaison de la performance des modèles selon la valeur de C :

```
1 from sklearn import svm
2
3 list_C = [1,10,100,1000]
4 for i in list_C:
5     svc = svm.SVC(kernel='linear', C=i).fit(X_train, y_train)
6     print(f'%_Classification_Correct_-_Lineaire_(Base_Apprentissage):_{svc.score(X_train, y_train):.3f}')
```

16. Comparaison de la performance des modèles selon le noyau utilisé :

```
1 from sklearn import svm
2
3 svc = svm.SVC(kernel='linear', C=10).fit(X_train, y_train)
4 rbf_svc = svm.SVC(kernel='rbf', gamma=0.7, C=10).fit(X_train, y_train)
5 poly_svc = svm.SVC(kernel='poly', degree=3, C=10).fit(X_train, y_train)
6 print(f'%_Classification_Correct_-_Lineaire_(Base_Test):_{svc.score(X_test, y_test):.3f}')
```

17. Questions sur la base des résultats obtenues suite aux questions précédentes :

- (a) Quelle valeur de C choisiriez-vous ?
- (b) Quel noyau choisiriez-vous ?

3 Réseaux de neurones

18. Divisez la base en deux :

```
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=0.15, random_state=6)
```

3.1 Couches Cachées

19. Variez le nombre de couches cachées :

```
1 import matplotlib.pyplot as plt
2 from sklearn.neural_network import MLPClassifier
3
4 n_couche = [1,10,20,50]
5 train_results = []
6 test_results = []
7 for i in n_couche:
8     rf = MLPClassifier(hidden_layer_sizes=(i), solver='lbfgs',
9                         random_state=21,)
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(n_couche, train_results, label='Erreur_d_apprentissage')
18 line2, = plt.plot(n_couche, 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_de_couches_cachees')
22 fig1 = plt.gcf()
23 fig1.savefig('/Users/maitre/Dropbox/Assas/cours/machine_learning/slides/lecture11/hidden_layer_size.png')
24 plt.show()

```

3.2 Fonction d'activation

20. Comparez la performance des modèles selon la fonction d'activation :

```

1 from sklearn.neural_network import MLPClassifier
2
3 f_activation = ['identity', 'logistic', 'tanh', 'relu']
4 train_results = []
5 test_results = []
6 for i in f_activation:
7     rf = MLPClassifier(hidden_layer_sizes=(1), activation = i, solver = 'lbfgs',
8                        random_state=2)
9     rf.fit(X_train, y_train)
10    print(f'Test : {rf.score(X_test, y_test):.3f}')
11    print(f'Train : {rf.score(X_train, y_train):.3f}')

```

3.3 Questions

21. Vrai ou faux :

- Dans le cadre des réseaux de neurones, la performance du modèle augmente nécessairement avec le nombre de couches cachées
- Dans le cadre des réseaux de neurones, la fonction d'activation logistique produit nécessairement des meilleures performances que la fonction linéaire
- Dans le cadre des réseaux de neurones, on initialise les poids aléatoirement
- Dans le cadre des réseaux de neurones, on parle en général de deep learning quand le nombre de couches cachées est inférieur à 3
- Dans le cadre des machines à vecteur de support, on parle en général de deep learning quand on utilise un noyau gaussien
- L'algorithme de rétro-propagation a trois composantes fondamentales : reverse pass, forward pass et step-wise pass