

Support Vector Machines (SVM)

Un Support Vector Machines (SVM) est un modèle de machine learning très puissant et polyvalent, capable d'effectuer une classification linéaire ou non linéaire, une régression et même une détection des outliers. C'est l'un des modèles les plus populaires de l'apprentissage automatique et toute personne intéressée par l'apprentissage automatique devrait l'avoir dans sa boîte à outils. Dans cet article, nous ferons une introduction aux SVM et implémenterons un SVM en python.

Support Vector Machines a.k.a. SVM, Kezako ?

Comme présenté en introduction, le SVM est un modèle d'apprentissage automatique supervisé qui est principalement utilisé pour les classifications (mais il peut aussi être utilisé pour la régression !). L'intuition derrière les Support Vector Machines est de simplement séparer des données en les délimitant (créer des frontières) afin de créer des groupes.

En d'autres termes, les SVM visent à résoudre les problèmes de classification en trouvant de bonnes frontières de décision (voir figure ci-dessous) entre deux ensembles de points appartenant à deux catégories différentes. Une frontière de décision peut être considérée comme une ligne ou une surface séparant vos données d'apprentissage en deux espaces correspondant à deux catégories. Pour classer de nouveaux points de données, il suffit de vérifier de quel côté de la frontière de décision ils se trouvent.

Les SVM procèdent à la recherche de ces frontières en deux étapes :

- 1- Les données sont mises en correspondance avec une nouvelle représentation à haute dimension où la frontière de décision peut être exprimée sous la forme d'un hyperplan (si les données étaient bidimensionnelles, comme dans la figure ci-dessous, un hyperplan serait une ligne droite).
- 2- Une bonne limite de décision (un hyperplan de séparation) est calculée en essayant de maximiser la distance entre l'hyperplan et les points de données les plus proches de chaque classe, une étape appelée maximisation de la marge. Cela permet à la frontière de bien s'adapter à de nouveaux échantillons en dehors de l'ensemble de données d'apprentissage.

Cette technique utilisée par les Support Vector Machines est appelée "kernel trick". Elle permet de transformer les données, puis, sur la base de ces transformations, il trouve une limite optimale entre les résultats possibles. En d'autres termes, il effectue des transformations de données extrêmement complexes, puis détermine comment séparer les données en fonction des labels que vous avez définis.

Les SVM en python

On importe le dataset que nous utiliserons dans ce tutoriel.

In [2]:

```
1 #Import du data set et des libs
2 %matplotlib inline
3 from sklearn import datasets
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import pandas as pd
7 import numpy as np
8 #chargement du jeu de donnees
9 wine_dataset = datasets.load_wine()
```

Exploration des données

Après avoir importé notre jeu de données, nous allons essayer de l'explorer et le visualiser afin d'avoir le maximum d'informations pour notre modélisation

In [3]:

```
1 #On a un dictionnaire python dont les différentes clés sont
2 # 'data', 'target', 'target_names', 'DESCR', 'feature_names'
3 print(wine_dataset.keys())
4
```

```
dict_keys(['data', 'target', 'frame', 'target_names', 'DESCR', 'feature_name
s'])
```

In [4]:

```
1 # Les features
2 print("Features: ", wine_dataset.feature_names)
3
```

```
Features: ['alcohol', 'malic_acid', 'ash', 'alcalinity_of_ash', 'magnesium', 'total_phenols', 'flavanoids', 'nonflavanoid_phenols', 'proanthocyanins', 'color_intensity', 'hue', 'od280/od315_of_diluted_wines', 'proline']
```

On voit donc qu'on a 3 labels ce qui pourrait dire qu'on a affaire à différentes catégories de vin. Allez, transformons tout ça en un DataFrame Pandas.

In [5]:

```
1 wine_dataframe = pd.DataFrame(data=wine_dataset['data'], columns=wine_dataset['feature_names'])
2 #On ajoute la colonne target ensuite
3 wine_dataframe['target'] = wine_dataset['target']
4 # Puis le nom des targets
5 wine_dataframe['class'] = wine_dataframe['target'].map(lambda target_name: wine_dataset['target_names'][target_name])
6
```

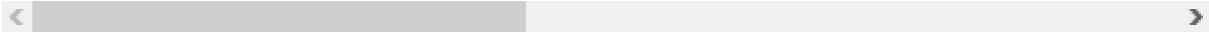
À ce niveau on peut faire un `head()` et un `describe()` pour mieux appréhender le DataFrame que nous venons de créer.

In [6]:

```
1 wine_dataframe.describe()
```

Out[6]:

	alcohol	malic_acid	ash	alcalinity_of_ash	magnesium	total_phenols	flavan
count	178.000000	178.000000	178.000000	178.000000	178.000000	178.000000	178.000000
mean	13.000618	2.336348	2.366517	19.494944	99.741573	2.295112	2.020613
std	0.811827	1.117146	0.274344	3.339564	14.282484	0.625851	0.990128
min	11.030000	0.740000	1.360000	10.600000	70.000000	0.980000	0.340000
25%	12.362500	1.602500	2.210000	17.200000	88.000000	1.742500	1.200000
50%	13.050000	1.865000	2.360000	19.500000	98.000000	2.355000	2.130000
75%	13.677500	3.082500	2.557500	21.500000	107.000000	2.800000	2.870000
max	14.830000	5.800000	3.230000	30.000000	162.000000	3.880000	5.080000



Un peu de visualisation...

Le describe() nous montre que la moyenne d'alcool contenu dans les vins est d'envrion 13%. Essayons de visualiser la distribution d'alcool pour chaque catégorie de vin.

In [7]:

```

1 for i in wine_dataframe.target.unique():
2     sns.distplot(wine_dataframe['alcohol'][wine_dataframe.target==i],
3                 kde=1,label='{}'.format(i))
4
5 plt.legend()

```

C:\Users\MSI\anaconda3\lib\site-packages\seaborn\distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

warnings.warn(msg, FutureWarning)

C:\Users\MSI\anaconda3\lib\site-packages\seaborn\distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

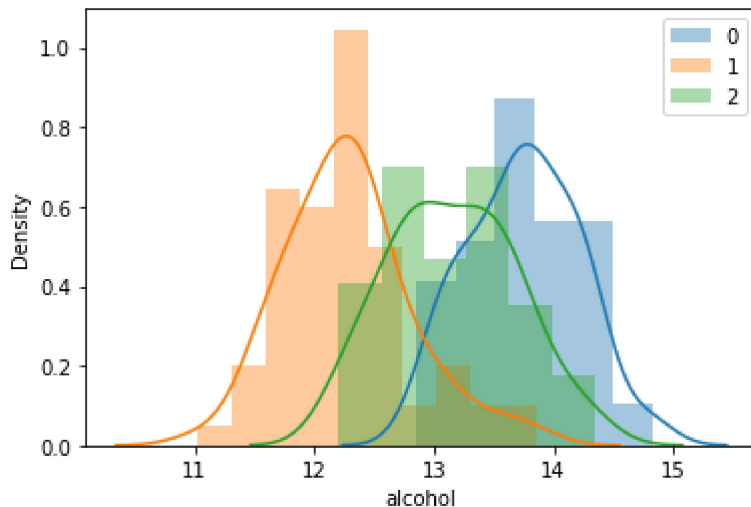
warnings.warn(msg, FutureWarning)

C:\Users\MSI\anaconda3\lib\site-packages\seaborn\distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

warnings.warn(msg, FutureWarning)

Out[7]:

<matplotlib.legend.Legend at 0x1836ba980a0>



On s'aperçoit que les distributions semblent plutôt normales et peuvent être classées en 3 groupes. Les vins faibles en teneur d'alcool, les vins à teneur d'alcool moyenne et les vins forts. On peut s'amuser à analyser d'autres caractéristiques, mais notre objectif dans ce tutoriel est de montrer l'utilisation du SVM.

Création de notre modèle SVM

In [8]:

```

1 from sklearn.model_selection import train_test_split
2
3 X_train, X_test, y_train, y_test = \
4     train_test_split(wine_dataset['data'], wine_dataset['target'],
5                     test_size=0.2)

```

Nous pouvons dès à présent créer notre SVM dans sklearn cela consiste à créer un objet SVC (support vector classifier). L'un des paramètres importants est le noyau 'kernel', comme présenté plus haut, c'est une fonction qui sert à transformer les données dans une représentation spécifique. Les SVM utilisent différents types de fonctions noyau. Ces fonctions sont de différents types, par exemple, linéaire, non linéaire, polynomiale, fonction de base radiale (RBF) et sigmoïde. Il faut donc avoir une attention particulière sur ce paramètre. Je vous laisse la documentation de scikit-learn sur les kernel pour mieux approfondir ce point. Dans notre cas nous choisirons un kernel linéaire (pourquoi 😊)

In [9]:

```

1 # Fit du Training set
2 from sklearn.svm import SVC
3 classifier = SVC(kernel = 'linear', random_state = 0)
4 classifier.fit(X_train, y_train)
5 #Prediction sur Le Test set
6 y_pred = classifier.predict(X_test)

```

In [10]:

```

1 from sklearn.metrics import classification_report
2 print(classification_report(y_test, y_pred))

```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	9
1	1.00	0.94	0.97	18
2	0.90	1.00	0.95	9
accuracy			0.97	36
macro avg	0.97	0.98	0.97	36
weighted avg	0.98	0.97	0.97	36

On regarde (ci-dessus) le rapport de classification pour voir ce que ça donne. Et ça semble pas mal.

Conclusion

Dans cet article nous avons présenté les Support Vector Machines qui sont des modèles très intéressants dont les principaux points forts sont les suivants :

- L'efficacité dans les espaces à dimension élevés
- L'efficacité dans les cas où le nombre de dimensions est supérieur au nombre d'échantillons.
- L'utilisation d'un sous-ensemble de points d'apprentissage dans la fonction de décision (appelés vecteurs de support), ce qui la rend également efficace en termes de mémoire.

- La polyvalence/ flexibilité : différentes fonctions de noyau peuvent être spécifiées pour la fonction de décision.

Cependant, il faut faire attention au cas particulier où le nombre de " features " est beaucoup plus grand que le nombre d'échantillons, car un mauvais le choix des fonctions de noyau peut entrainer l'over-fitting.

In []:

1	
---	--