



Master Informatique EID2

Traitement numérique des données

TP2 Prétraitement et visualisation de
données

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn import *
4 from sklearn.preprocessing import *
5 from sklearn.decomposition import PCA
6 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis

```

```

1 #####
2 # A. Normalisation de données #
3 #####
4
5 #-----#
6 ## 1- Créez la matrice X suivante ##
7 #-----#
8 X = np.array([[1,-1,2],[2,0,0],[0,1,-1]],float)

```

```

1 #-----#
2 ## 2- Visualisez X et calculez la moyenne et la variance de X ##
3 #-----#
4 print('X :\n',X)
5 print('Moyenne = ', X.mean())
6 print('Variance = ', X.var())

```

```

X :
[[ 1. -1.  2.]
 [ 2.  0.  0.]
 [ 0.  1. -1.]]
Moyenne =  0.4444444444444444
Variance =  1.1358024691358024

```

```

1 #-----#
2 ## 3- Utilisez la fonction scale pour normaliser la matrice X. Que constatez vous ? ##
3 #-----#
4 Xnorm = scale(X)
5 print('Xnorm : \n', Xnorm)

```

```

Xnorm :
[[ 0.         -1.22474487  1.33630621]
 [ 1.22474487  0.         -0.26726124]
 [-1.22474487  1.22474487 -1.06904497]]

```

```

1 #-----#
2 ## 4- Calculer la moyenne et la variance de la matrice X normalisé ##
3 #-----#
4
5 print('Moyenne = ', Xnorm.mean())
6 print('Variance = ', Xnorm.var())
7
8 #####
9 # ===== Expliquez le résultat obtenu ===== #
10 # Nous avons normalisé de façon à ce que la Variance soit égale à 1 pour #
11 # réduire l'écart entre nos valeurs Min et Max #
12 #####

```

```

Moyenne =  4.9343245538895844e-17
Variance =  1.0

```

```

1 #####
2 # B. Normalisation MinMax #
3 #####
4
5 #-----#
6 ## 1- Créez la matrice de données X2 suivante ##
7 #-----#
8 X2 = np.copy(X)

```

```

1 #-----#
2 ## 2- Visualisez la matrice et calculez la moyenne sur les variables ##
3 #-----#
4 print('X2 :\n',X2)
5 print('Moyenne = ', X2.mean())

```

```

X2 :
[[ 1. -1.  2.]
 [ 2.  0.  0.]
 [ 0.  1. -1.]]
Moyenne =  0.4444444444444444

```

```

1 #-----#
2 ##----- 3- Normalisez les données dans l'intervalle [0 1] -----##
3 ## Visualisez les données normalisées et calculez la moyenne sur les variables ##
4 #-----#
5
6 Xnorm2 = MinMaxScaler((0, 1)).fit_transform(X2)
7 print('Xnorm : \n', Xnorm2)
8 print('Moyenne = ', Xnorm2.mean())
9 print('Variance = ', Xnorm2.var())
10
11 #####
12 # ===== Que constatez-vous ? ===== #
13 # MinMaxScaler normalise de la façon suivante  $x = (X - \text{Min}) / (\text{Max} - \text{Min})$  #
14 # Dans notre cas nous avons des valeurs comprise entre 0 et 1 #
15 #####

```

```

Xnorm :
[[0.5      0.      1.      ]
 [1.      0.5     0.33333333]
 [0.      1.      0.      ]]
Moyenne =  0.48148148148148145
Variance =  0.16941015089163236

```

```

1 #####
2 # C. visualisation de données #
3 #####
4
5 #-----#
6 ## 1- Chargez les données Iris ##
7 #-----#
8 iris = datasets.load_iris()
9 data = iris.data
10 target = iris.target
11 var = iris.feature_names
12 target_n = iris.target_names
13 print(target_n)

```

```
['setosa' 'versicolor' 'virginica']
```

```

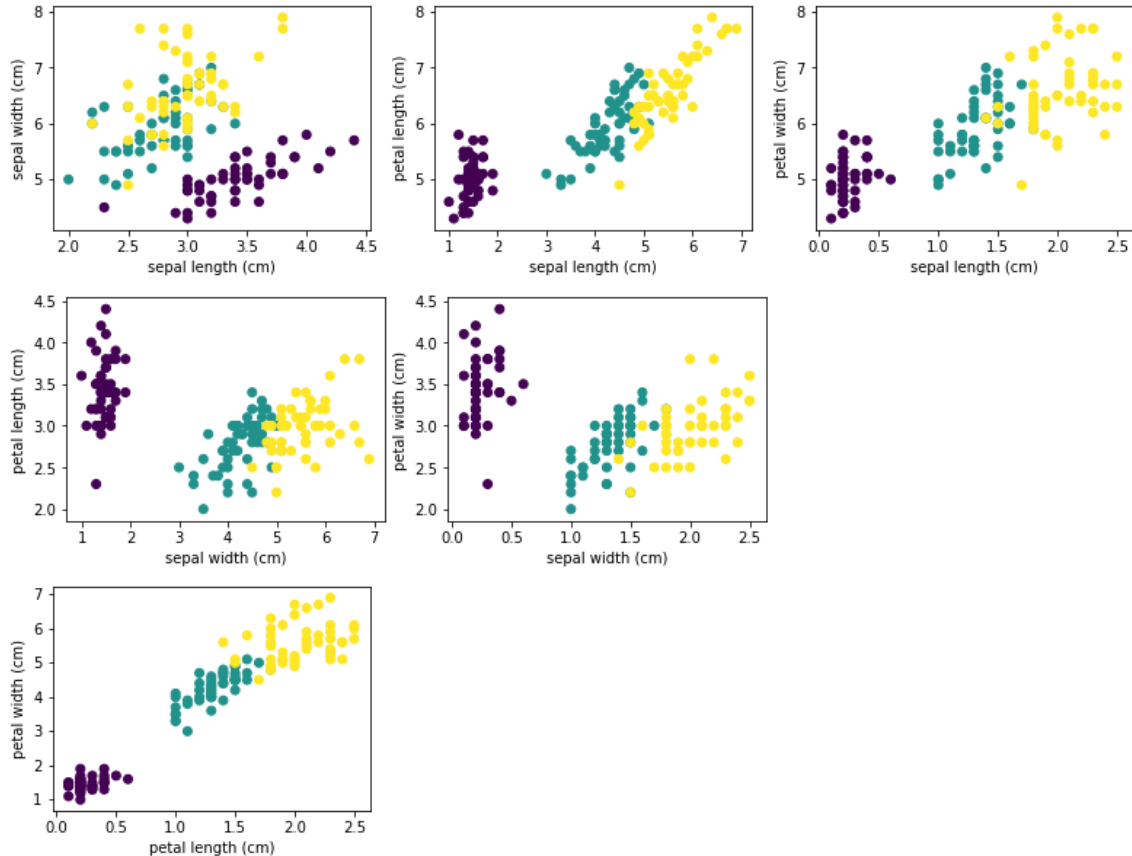
1 #-----#
2 ## 2- Visualisez le nuage de points en 2D avec des couleurs correspondant aux classes en ##
3 ##----- utilisant toutes les combinaisons de variables -----##
4 #-----#
5 for i in range(len(var)):
6     plt.figure(figsize=(20,15))
7     for j in range(i+1,len(var)):
8         plt.subplot(4,4,j+1)
9         plt.scatter(data[:,j],data[:,i], c=iris.target)
10        plt.xlabel(var[i])
11        plt.ylabel(var[j])
12    plt.show()
13
14
15

```

```

16 #####
17 # ===== Quelle est la meilleure visualisation? =====#
18 # ===== Justifiez votre réponse ===== #
19 # La meilleure visualisation est PetatWidth vs SepalWidth car on voit #
20 # une bonne séparation des 3 classes sans corrélation #
21 #####
22

```



```

1 #####
2 ## D. Réduction de dimensions et visualisation de données ##
3 #####
4
5 #-----#
6 ## 2- Les méthodes PCA et LDA peuvent être importées à partir des packages suivants ##
7 #-----#
8
9 pca = PCA(n_components=2)
10 irisPCA = pca.fit(data).transform(data)
11 print('irisPCA (10 premières lignes) :\n',irisPCA[:10])

```

```

irisPCA (10 premières lignes) :
[[-2.68412563  0.31939725]
 [-2.71414169 -0.17700123]
 [-2.88899057 -0.14494943]
 [-2.74534286 -0.31829898]
 [-2.72871654  0.32675451]
 [-2.28085963  0.74133045]
 [-2.82053775 -0.08946138]
 [-2.62614497  0.16338496]
 [-2.88638273 -0.57831175]
 [-2.6727558  -0.11377425]]

```

```

1 lda = LinearDiscriminantAnalysis(n_components=2)
2 irisLDA = lda.fit(data, target).transform(data)
3 print('irisLDA (10 premières lignes):\n',irisLDA[:10])

```

```

irisLDA (10 premières lignes):
[[ 8.06179978  0.30042062]
 [ 7.12868772 -0.78666043]
 [ 7.48982797 -0.26538449]
 [ 6.81320057 -0.67063107]
 [ 8.13230933  0.51446253]
 [ 7.70194674  1.46172097]
 [ 7.21261762  0.35583621]
 [ 7.60529355 -0.01163384]
 [ 6.56055159 -1.01516362]
 [ 7.34305989 -0.94731921]]

```

```

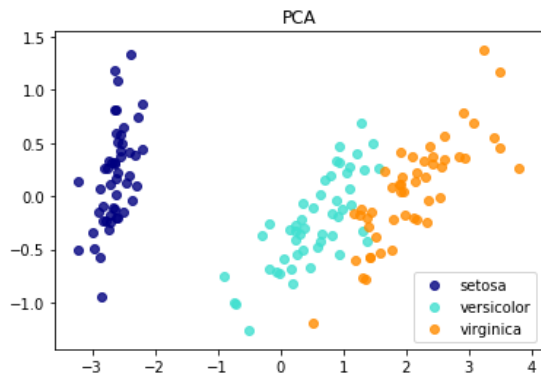
1 #-----#
2 ## 3- Visualisez les nuages de points avec les nouvelles axes obtenus : une image pour l'ACP ##
3 ##----- et une autre pour l'ADL et utiliser la classe de Iris comme couleur de points -----##
4 #-----#
5 def affichescatter(x,y,z,t):
6     color = ['navy', 'turquoise', 'darkorange']
7     plt.figure()
8     plt.title(t)
9     for color, i, target_name in zip(color, [0, 1, 2], y):
10         plt.scatter(x[y == i, 0], x[y == i, 1], color=color, alpha=.8, label=z[i])
11     plt.legend(loc='best', shadow=False, scatterpoints=1)
12     plt.show()

```

```

1 #####
2 #----- ACP -----#
3 #####
4 affichescatter(irisPCA,target,target_n,'PCA')

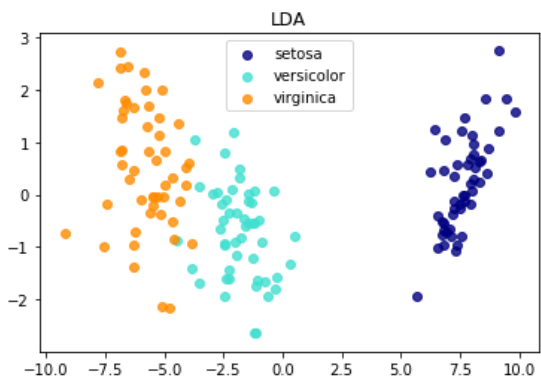
```



```

1 #####
2 #----- LDA -----#
3 #####
4 affichescatter(irisLDA,target,target_n,'LDA')

```



```
1 #####
2 ## Quelle différence constatez-vous entre les deux visualisations ? Expliquer votre raisonnement ##
3 ##-----##
4 # Pour IrisPCA, nous avons des données beaucoup plus proches les unes des autres (on le remarque à #
5 # l'échelle) tandis que pour IrisLDA nos données sont plus éloignées. Cela est dû au fait que PCA #
6 # cherche à réduire la dimension en trouvant des corrélations entre les colonnes tandis que LDA #
7 # cherche des corrélations sur les lignes #
8 #####
```