

TP6 : Analyse et prédiction des infections COVID-19

Importer les bibliothèques nécessaires

```
In [1]: import numpy as np
import pandas as pd
```

Importer le fichier de données après son nettoyage

```
In [2]: df = pd.read_csv("covid_data_cleaned.csv")
print("Les noms des colonnes: ",df.columns.values)
```

Les noms des colonnes:

```
['age', 'sex', 'chronic_disease_binary', 'asymptomatic', 'cough', 'fatigue', 'fever', 'headache', 'malaise', 'pneumonia', 'respiratory_problems', 'runny_nose', 'sore_throat', 'diabetes', 'hypertension', 'outcome']
```

```
In [3]: df.shape
Out[3]: (1577, 16)
```

```
In [4]: df
Out[6]:
```

	age	sex	chronic_disease_binary	asymptomatic	cough	fatigue	fever	headache	malaise	pneumonia	respiratory_problems	runny_nose	sore_throat
0	60	1	0	0	0	0	1	0	0	1	0	0	0
1	45	1	0	0	0	0	1	0	0	0	0	0	0
2	25	1	0	0	1	0	1	0	0	0	0	0	1
3	45	1	0	0	1	1	0	1	0	0	0	0	0
4	85	1	0	0	0	0	1	0	0	0	0	0	0
...
1572	65	1	0	0	1	0	0	0	0	0	0	0	0
1573	45	1	0	0	0	0	0	0	0	0	0	0	0
1574	35	1	1	0	1	0	1	0	0	0	0	0	0
1575	65	1	0	0	1	0	0	0	0	0	0	0	0
1576	39	0	0	0	0	0	0	0	0	0	0	0	0

1577 rows x 16 columns

1. Calculez les corrélations entre les variables. Quelles sont les variables les plus corrélées avec la cible (result)? Expliquez les résultats.

```
In [5]: df.corr()
Out[5]:
```

	age	sex	chronic_disease_binary	asymptomatic	cough	fatigue	fever	headache	malaise	pneu
age	1.000000	0.074035	0.026971	-0.092924	0.024976	0.075835	-0.036621	-0.075039	0.042810	0.0
sex	0.074035	1.000000	0.070941	-0.015488	-0.024624	-0.028209	0.017434	-0.020835	0.049048	0.0
chronic_disease_binary	0.026971	0.070941	1.000000	-0.053872	-0.052701	-0.030147	-0.172809	-0.050477	-0.052212	0.0
asymptomatic	-0.092924	-0.015488	-0.053872	1.000000	-0.130809	-0.099694	-0.194426	-0.038509	-0.032291	-0.0
cough	0.024976	-0.024624	-0.052701	-0.130809	1.000000	-0.019657	0.168018	-0.015364	-0.001978	-0.0
fatigue	0.075835	-0.028209	-0.030147	-0.099694	-0.019657	1.000000	0.020323	-0.000250	-0.038471	-0.0
fever	-0.036621	0.017434	-0.172809	-0.194426	0.168018	0.020323	1.000000	-0.003531	-0.099677	-0.0
headache	-0.075039	-0.020835	-0.050477	-0.38509	0.015364	-0.000250	-0.003531	1.000000	0.182293	-0.0
malaise	0.042810	0.049048	-0.052212	-0.032291	-0.001978	-0.038471	-0.099677	0.182293	1.000000	0.0
pneumonia	0.182247	0.063021	0.518432	-0.048326	-0.156608	-0.057397	-0.188569	-0.030234	0.012922	1.0
respiratory_problems	0.015409	0.002867	0.482077	-0.043363	-0.127954	-0.051502	-0.201340	-0.049964	-0.042027	0.0
runny_nose	-0.088150	-0.047703	-0.047338	-0.029492	0.055013	-0.035028	0.030465	0.065998	-0.028463	-0.0
sore_throat	-0.071339	-0.054563	-0.016424	-0.050060	0.046038	0.036587	-0.026851	0.127453	-0.099677	-0.0
diabetes	0.191144	0.060799	0.654556	-0.036671	-0.067095	-0.011836	-0.127674	-0.025950	-0.035541	-0.0
hypertension	0.246387	0.046194	0.573860	-0.041982	-0.102893	-0.021824	-0.162685	-0.033961	-0.064089	0.0
outcome	-0.296449	-0.082325	-0.728110	0.056613	0.093553	-0.008516	0.171968	0.054106	0.054869	-0.0

```
In [6]: import seaborn as sns
import matplotlib inline

# calculate the correlation matrix
corr = df.corr()

# plot the heatmap
sns.heatmap(corr,
            yticklabels=corr.columns,
            xticklabels=corr.columns)
```

```
Out[6]: <matplotlib.axes._subplots.AxesSubplot at 0x2702504b808>
```

On peut remarquer que les variables 'chronic_disease_binary' et 'hypertension' sont les plus corrélées avec la cible (result) ou (outcome), cela signifie que nous avons un lien et une dépendance entre ces variables et la variable cible de notre jeu de données

2. Visualisez les données en deux dimensions en passant par l'ACP (analyse en composantes principales). Pouvez-vous utiliser une autre méthode ?

Nous allons commencer par ACP

```
In [7]: X = df[['age', 'sex', 'chronic_disease_binary', 'asymptomatic', 'cough', 'fatigue', 'fever', 'headache', 'malaise', 'pneumonia', 'respiratory_problems', 'runny_nose', 'sore_throat', 'diabetes', 'hypertension', 'outcome']]
Y = df[['outcome']]

#cols = df.columns.tolist()
#X = df[cols[0:15]]
#Y = df[cols[15:16]]

from sklearn.decomposition import PCA
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis

pca = PCA(n_components=2)
pca.fit(X)
lda = LinearDiscriminantAnalysis(n_components=1)
lda.fit(X)
```

Pour la méthode LDA, elle utilise au maximum (n-1) colonnes. Il s'agit du nombre maximum de colonnes que nous voulons obtenir. Dans notre cas, nous avons 2 classes ('Dead', 'Discharged'), donc nous obtenons un maximum de 1 colonne.

```
In [9]: principal_df_pca = pd.DataFrame(data = principal_pca,
                                     columns = ['principal component 1', 'principal component 2'])

principal_df_lda = pd.DataFrame(data = principal_lda,
                               columns = ['principal component 1'])

In [10]: principal_df_pca.head()
```

```
Out[10]:
```

	principal component 1	principal component 2
0	-9.741707	-0.040233
1	5.260392	-0.179588
2	25.260267	-0.684727
3	5.259751	-0.725725
4	-34.738667	-0.279784

```
In [11]: principal_df_lda.head()
Out[11]:
```

	principal component 1
0	-0.486583
1	0.479503
2	0.957996
3	0.537840
4	0.098562

```
In [12]: import matplotlib.pyplot as plt

plt.figure()
plt.scatter(principal_pca[:, 0], principal_pca[:, 1], color='g', alpha=0.5, lw=1)
plt.scatter(principal_lda[:, 0], principal_lda[:, 1], color='r', alpha=0.5, lw=1)

plt.legend(loc='best', shadow=False, scatterpoints=1)
plt.title('PCA of dataset')

plt.show()
```

PCA of dataset

Dans la suite, nous utilisons une méthode d'apprentissage automatique afin de prédire la classe des patients soit «décédés» ('died') soit «sortis» ('discharged') de l'hôpital. Vous pouvez utiliser la classification par K-Nearest Neighbours (K-NN), l'arbre de décision ou le classificateur Bayes

3. Les résultats obtenus doivent être validés en utilisant certains indices externes comme l'erreur de prédiction (matrice de confusion et précision) ou d'autres comme Rappel, F-Mesure, ...

```
In [13]: X = df[['age', 'sex', 'chronic_disease_binary', 'asymptomatic', 'cough', 'fatigue', 'fever', 'headache', 'malaise', 'pneumonia', 'respiratory_problems', 'runny_nose', 'sore_throat', 'diabetes', 'hypertension', 'outcome']]
Y = df[['outcome']]
```

K-Nearest Neighbours (K-NN)

Créer une fonction qui cherche le meilleur k

```
In [26]: from sklearn.model_selection import train_test_split
import numpy as np
import random
from sklearn.neighbors import KNeighborsClassifier

def best_k_KNN(X_train, y_train, X_test, y_test):
    ks = 5
    mean_acc = np.zeros((Ks-1))
    for k in range(1,Ks):
        neigh = KNeighborsClassifier(n_neighbors = k).fit(X_train,y_train)
        mean_acc[k-1] = neigh.score(X_test, y_test)

    plt.plot(range(1,Ks),mean_acc,'g')
    plt.legend(['Accuracy'],loc='best')
    plt.ylabel('Accuracy')
    plt.xlabel('Number of neighbors (K)')
    plt.tight_layout()
    plt.show()
    print("The best accuracy was ", mean_acc.max(), " with k=", mean_acc.argmax()+1)

    return mean_acc.argmax()+1
```

diviser l'ensemble de données en données de test et d'apprentissage

```
In [27]: X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.2, random_state=random.seed())
```

Trouver le meilleur k

```
In [28]: best_k = best_k_KNN(X_train, y_train, X_test, y_test)

100% 49/49 [00:01<00:00, 40.66it/s]
```

The best accuracy was 0.9620253164556962 with k= 5

```
In [29]: knn = KNeighborsClassifier(n_neighbors = best_k).fit(X_train,y_train)
y_pred = knn.predict(X_test)

# Compute the mean squared error of our predictions
mse = ((y_pred - y_test) ** 2).sum() / len(y_pred)
print("Le pourcentage d'erreur pour le modèle KNN est : ", mse*100, "%")
```

Le pourcentage d'erreur pour le modèle KNN est : 3.79746835443038 %

```
In [30]: res_knn = knn.score(X_test,y_test)
print("Accuracy pour le modèle KNN est : ", res_knn*100, "%")

Accuracy pour le modèle KNN est : 96.20253164556962 %
```

```
In [31]: import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.metrics import plot_confusion_matrix

disp = plot_confusion_matrix(knn, X_test, y_test,cmap=plt.cm.Blues)

disp.ax_.set_title("Confusion Matrix")
plt.show()
#matplotlib inline
```

Confusion Matrix

```
In [32]: from sklearn.metrics import accuracy_score
print("Accuracy score: ",accuracy_score(y_test, y_pred))
from sklearn.metrics import f1_score
print("F1-Measure: ",f1_score(y_test, y_pred,average='macro'))
from sklearn.metrics import precision_score
print("Precision score: ",precision_score(y_test, y_pred, average='macro'))
from sklearn.metrics import recall_score
print("Recall score: ",recall_score(y_test, y_pred, average='macro'))
```

Accuracy score: 0.9620253164556962
F-Measure: 0.9670152287636067
Precision score: 0.968018918918919
Recall score: 0.96149184149184149

```
In [33]: # Classification report
from sklearn.metrics import classification_report
print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0	0.95	0.63	0.76	30
1	0.96	1.00	0.98	286
accuracy			0.96	316
macro avg	0.96	0.81	0.87	316
weighted avg	0.96	0.96	0.96	316

Arbre de décision

```
In [34]: from sklearn.tree import DecisionTreeClassifier
clf = tree.DecisionTreeClassifier()
clf.fit(X_train,y_train)

y_pred2 = clf.predict(X_test)
```

```
In [35]: import matplotlib.pyplot as plt
from sklearn.metrics import plot_confusion_matrix
disp = plot_confusion_matrix(clf, X_test, y_test,cmap=plt.cm.Blues)

disp.ax_.set_title("Confusion Matrix")
plt.show()
#matplotlib inline
```

Confusion Matrix

```
In [36]: from sklearn.metrics import accuracy_score
print("Accuracy score: ",accuracy_score(y_test, y_pred2))
from sklearn.metrics import f1_score
print("F1-Measure: ",f1_score(y_test, y_pred2,average='macro'))
from sklearn.metrics import precision_score
print("Precision score: ",precision_score(y_test, pred, average='macro'))
from sklearn.metrics import recall_score
print("Recall score: ",recall_score(y_test, pred, average='macro'))
```

Accuracy score: 0.956962025316456
F-Measure: 0.9670152287636067
Precision score: 0.968018918918919
Recall score: 0.96149184149184149

```
In [37]: # Classification report
from sklearn.metrics import classification_report
print(classification_report(y_test, y_pred2))
```

	precision	recall	f1-score	support
0	0.79	0.73	0.76	30
1	0.97	0.98	0.98	286
accuracy			0.96	316
macro avg	0.88	0.86	0.87	316
weighted avg	0.95	0.96	0.96	316

Classificateur Bayes

```
In [38]: from sklearn.naive_bayes import MultinomialNB
clf_CB = MultinomialNB()
clf_CB.fit(X_train, y_train)
y_pred3 = clf_CB.predict(X_test)
```

```
In [39]: import matplotlib.pyplot as plt
from sklearn.metrics import plot_confusion_matrix
disp = plot_confusion_matrix(clf_CB, X_test, y_test,cmap=plt.cm.Blues)

disp.ax_.set_title("Confusion Matrix")
plt.show()
```

Confusion Matrix

```
In [40]: from sklearn.metrics import accuracy_score
print("Accuracy score: ",accuracy_score(y_test, y_pred3))
from sklearn.metrics import f1_score
print("F1-Measure: ",f1_score(y_test, y_pred3,average='macro'))
from sklearn.metrics import precision_score
print("Precision score: ",precision_score(y_test, pred, average='macro'))
from sklearn.metrics import recall_score
print("Recall score: ",recall_score(y_test, pred, average='macro'))
```

Accuracy score: 0.9778481012658228
F-Measure: 0.9278135914230793
Precision score: 0.968018918918919
Recall score: 0.96149184149184149

```
In [41]: # Classification report
from sklearn.metrics import classification_report
print(classification_report(y_test, y_pred3))
```

	precision	recall	f1-score	support
0	1.00	0.77	0.87	30
1	0.98	1.00	0.99	286
accuracy			0.99	316
macro avg	0.99	0.88	0.93	316
weighted avg	0.98	0.98	0.98	316

4. Utilisez la régression pour prédire l'âge (age) des personnes en fonction d'autres variables. Vous avez le choix sur ces variables explicatives ? Comment choisir-vous ces variables ? Calculez la qualité de la prédiction à l'aide de l'erreur MSE (Mean Squared Error)

Tester sur tous les variables

```
In [42]: X2 = df[['sex', 'chronic_disease_binary', 'asymptomatic', 'cough', 'fatigue', 'fever', 'headache', 'malaise', 'pneumonia', 'respiratory_problems', 'runny_nose', 'sore_throat', 'diabetes', 'hypertension', 'outcome']]
Y2 = df[['age']]

In [43]: X_train2, X_test2, y_train2, y_test2 = train_test_split(X2, Y2, test_size=0.2, random_state=random.seed())

In [44]: from sklearn.linear_model import LinearRegression
reg = LinearRegression()
reg.fit(X_train2, y_train2)
```

```
Out[44]: LinearRegression()

In [45]: y2_pred = reg.predict(X_test2)
```

```
In [46]: from sklearn.metrics import mean_squared_error
MSE = mean_squared_error(y_test2, y2_pred)
print("l'erreur MSE : ", MSE)
```

l'erreur MSE : 321.601764906366

Utiliser RFE pour sélectionner les variables explicatives :

```
In [47]: from sklearn.feature_selection import RFE
selecteur = RFE(estimator = reg)

sel = selecteur.fit(X_test2, y_test2)

print(sel.n_features_)
print(sel.support_)

7
[False False False False True True False False False True True True
 False True True]
```

```
In [48]: # get the feature support boolean list
feature_mask = sel.support_
# reduce the dataframe to just those features
X_scaled_rfe_reduced = X.iloc[:,feature_mask]
```

```
Out[48]: X_scaled_rfe_reduced
```

	cough	fatigue	pneumonia	respiratory_problems	runny_nose	diabetes	hypertension
0	0	0	1	0	0	0	0
1	0	0	0	0	0	0	0
2	1	0	0	0	0	0	0
3	1	0	0	0	0	0	0
4	0	0	0	0	0	0	0
...
1572	1	0	0	0	0	0	0
1573	0	0	0	0	0	0	0
1574	1	0	0	0	0	0	0
1575	1	0	0	0	0	0	0
1576	0	0	0	0	0	0	0

1577 rows x 7 columns

```
In [50]: X2_RFE = X_scaled_rfe_reduced
Y2_RFE = df[['age']]

In [51]: X_train_RFE, X_test_RFE, y_train_RFE, y_test_RFE = train_test_split(X2_RFE, Y2_RFE, test_size=0.2, random_state=random.seed())

In [52]: reg_RFE = LinearRegression()
reg_RFE.fit(X_train_RFE, y_train_RFE)
y2_pred_RFE = reg_RFE.predict(X_test_RFE)
```

```
In [53]: from sklearn.metrics import mean_squared_error
MSE = mean_squared_error(y_test_RFE, y2_pred_RFE)
print("l'erreur MSE : ", MSE)
```

l'erreur MSE : 292.0825627294294

5. Appliquez trois méthodes de clustering (K-means, NMF et CAH) de l'ensemble de données pour segmenter les personnes en différents groupes. Utilisez l'index de Silhouette pour connaître le meilleur nombre de clusters.

K-means

```
In [57]: from sklearn.metrics import silhouette_score
from sklearn.cluster import KMeans

best_clusters = 0
silh_score = 0.0
silhouette_coefficients = []

# On commence à 2 clusters pour le coefficient de silhouette
for k in range(2, 21):
    kmeans = KMeans(n_clusters=k, init='random', max_iter=600)
    score = silhouette_score(principal_pca, kmeans.labels_)
    silhouette_coefficients.append(score)
    if score > silh_score:
        silh_score = score
        best_clusters = k

plt.style.use('fivethirtyeight')
plt.plot(range(2, 21), silhouette_coefficients)
plt.xticks(range(2, 21))
plt.xlabel('Number of Clusters')
plt.ylabel('Silhouette Coefficient')
plt.show()
```

D'après l'index de Silhouette, le meilleur nombre de clusters est: 'best_clusters,' avec un score de

```
Out[57]: 0.6250901392392571
```

D'après l'index de Silhouette, le meilleur nombre de clusters est: 18 avec un score de : 0.6335462472764803

```
In [58]: from sklearn import metrics

kmeans = KMeans(n_clusters=best_clusters, random_state=0).fit(principal_pca)
y_kmeans = kmeans.predict(principal_pca)
labels = kmeans.labels_
metrics.silhouette_score(principal_pca, labels, metric='euclidean')
```

```
Out[58]: 0.6250901392392571
```

NMF

```
In [60]: from sklearn.decomposition import NMF
import sklearn.preprocessing as preprocess

best_clusters_nmf = 0
silh_score_nmf = 0.0
silhouette_coefficients = []

# assurer que toutes les données ne sont pas négatives
X_scaled = preprocess.minmax_scale(principal_pca)

# On commence à 2 clusters pour le coefficient de silhouette
for k in range(2, 21):
    nmf = NMF(n_components=k, init='random', max_iter=600)
    D = nmf.fit_transform(X_scaled)
    W = nmf.components_
    nmf_labels = np.argmax(D, axis=1)
    score = silhouette_score(X_scaled, nmf_labels)
    silhouette_coefficients.append(score)
    if score > silh_score_nmf:
        silh_score_nmf = score
        best_clusters_nmf = k

plt.style.use('fivethirtyeight')
plt.plot(range(2, 21), silhouette_coefficients)
plt.xticks(range(2, 21))
plt.xlabel('Number of Clusters')
plt.ylabel('Silhouette Coefficient')
plt.show()
```

D'après l'index de Silhouette, le meilleur nombre de clusters est: 'best_clusters_nmf,' avec un score de

```
Out[60]: 0.2938627163070187
```

D'après l'index de Silhouette, le meilleur nombre de clusters est: 2 avec un score de : 0.2938627163070187

```
In [61]: nmf = NMF(n_components=best_clusters_nmf, init='random', random_state=0, max_iter=200)
W = nmf.fit_transform(X)
H = nmf.components_

In [62]: W = np.around(W, decimals = 2)
print("W :\n",W)
```

W :

```
[[[4.02 0.02]
 [2.46 0.02]
 [1.46 0.04]
 [1.17 0.04]
 [4.42 0.02]
 [2.74 0. ]]]
```

```
In [63]: H = np.around(H, decimals = 2)
print("H :\n",H)
```

H :

```
[[[1.4230e+01 7.0000e-02 3.0000e-02 1.0000e-02 0.0000e+00 1.0000e-02
 [4.445e+01 6.3000e-01 9.0000e-02 ... 8.0000e-02 3.0000e-02 6.0000e-02]
 [5.534e+01 9.9000e-01 7.0000e-02 ... 8.0000e-02 2.0000e-02 4.0000e-02]
 [1.1647e+02 2.1000e+01 0.0000e+00 0.0000e+00 2.1930e+01 9.0000e-02
 [4.4230e+01 1.2300e+01 0.0000e+00 0.0000e+00 2.0000e+00 1.1700e+00]
 [1.5700e+00 0.0000e+00 
```


