

# ACP – Python

Yéro Diamanka

2024-02-12

Nous allons explorer le jeu de données `trees.csv` qui donne des mesures de diamètre (girth, en pouces), hauteur (height, en pieds) et volume (volume, en pieds cubes) de cerisiers noirs.

## 0.1 Importation de bibliothèques !

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.decomposition import PCA
6 from sklearn.preprocessing import StandardScaler
```

```
1 trees=pd.read_csv('trees.csv')
```

1 - Centrer et réduire les données

```
1 scaler = StandardScaler()
2 scaler.fit(trees)
3 Z = scaler.transform(trees)
4 Z
```

2 - Lancer l'ACP

```
1 pca = PCA()
2 pca.fit(Z)
```

3 - Examiner les ratios cumulés de variance. Combien de variabilité est expliquée par les deux premiers axes ?

```
1 np.cumsum(pca.explained_variance_ratio_)
```

99% de la variabilité est expliquée par les deux premiers axes, c'est très élevé.

4 - Afficher les valeurs propres

```
1 l = 3*pca.explained_variance_ratio_
2 l
```

5 - Extraire les facteurs

```
1 Gn=pca.components_
2 G=np.empty(shape=Gn.shape)
3 for i in range(0, Gn.shape[0]):
4     G[i,:]=Gn[i,:]*np.sqrt(l[i])
5 G
```

6 - Tracer le cercle des corrélations dans le premier plan factoriel pour les variables

```

1 fig,ax=plt.subplots(figsize=(5,5))
2 for i in range(0, Gn.shape[1]):
3     ax.arrow(0,0, # la flèche part de l'origine
4             G[0, i], G[1, i], # et arrive en (G_1i,G_2i)
5             head_width=0.05,head_length=0.07,length_includes_head=True)
6     ax.text(G[0, i] + 0.01,G[1, i]-0.02, trees.columns[i],fontsize=8)
7 # affichage des lignes horizontales et verticales
8 ax.plot([-1, 1], [0, 0], color='grey', ls='--')
9 ax.plot([0, 0], [-1, 1], color='grey', ls='--')
10 # nom des axes, avec le pourcentage d'inertie expliqué
11 ax.set_xlabel('G{} ({}%)'.format(1,
12 round(100*pca.explained_variance_ratio_[0],1)))
13 ax.set_ylabel('G{} ({}%)'.format(2,
14 round(100*pca.explained_variance_ratio_[1],1)))
15 ax.set_title("Cercle des corrélations (G{} et G{}).format(1, 2))
16 an = np.linspace(0, 2 * np.pi, 100)
17 ax.plot(np.cos(an), np.sin(an))

```

Examiner le cercle des corrélation dans le premier plan factoriel pour les variables.

7 - Quelles variables sont bien représentées dans ce plan ?

Toutes : les pointes des flèches sont très proches du cercle unité. On le confirme en calculant les cos carrés.

```

1 Gsq = G**2
2 print(Gsq[0, :]+Gsq[1, :])

```

8 - Que pensez vous des positions relatives des variables Girth et Volume dans ce plan ?

Elles sont proches, donc très corrélées positivement.

9 - Quelles sont les variables qui contribuent le plus au premier axe factoriel ?

```

1 Contrib=(Gn**2)/np.sum(Gn**2,axis=0)
2 print(Contrib[0, :])

```

Ce sont les variables de diamètre et de volume.

10 - Quelles sont les variables qui contribuent le plus au deuxième axe factoriel ?

```

1 print(Contrib[1, :])

```

C'est la variable de hauteur.

11 - Extraire les composantes principales

```

1 F=pca.fit_transform(Z)
2 F

```

12 - Tracer le nuage de point des individus projetés dans le premier plan factoriel.

```

1 fig,ax=plt.subplots(figsize=(5,5))
2 # individus

```

```

3 ax.scatter(F[:,0],F[:,1])
4 for i in range(trees.shape[0]):
5     ax.text(F[i,0]+0.1,F[i,1],'{}'.format(i),fontsize=8)
6 ax.set_xlabel('F{} ({}%)'.format(1,
round(100*pca.explained_variance_ratio_[0],1)))
7 ax.set_ylabel('F{} ({}%)'.format(2,
round(100*pca.explained_variance_ratio_[1],1)))
8 ax.set_title("Individus projetés (F{} et F{}).format(1, 2))
9 ax.grid()
10 ax.plot([min(F[:,0]), max(F[:,0])],[0,0], linestyle="--", color='C7')
11 ax.plot([0, 0],[min(F[:,1]), max(F[:,1])], linestyle="--", color='C7')

```

13 - Quels sont les individus qui sont bien représentés dans ce plan ?

```

1 cos2ind = pd.DataFrame(
2     columns=['axe 1', 'axe 2', 'somme'],
3     index=np.arange(trees.shape[0]))
4 for i in np.arange(trees.shape[0]):
5     for k in np.arange(2):
6         cos2ind.iloc[i,k] = F[i,k]**2/(sum(Z[i,:]**2))
7     cos2ind.iloc[i,2]=cos2ind.iloc[i,0]+cos2ind.iloc[i,1]
8 cos2ind

```

Tous les individus sont très bien représentés sur ce plan. Seul l'arbre 15 est un peu moins bien représenté que les autres.

14 - Que peut-on dire des arbres qui sont le plus à droite sur premier plan factoriel ?

Ce sont ceux de plus grand diamètre/volume

15 - Que peut-on dire des arbres qui sont le plus à gauche sur premier plan factoriel ?

Ce sont ceux de plus petit diamètre/volume

16 - Que peut-on dire des arbres qui sont le plus en haut (resp. le plus en bas) sur premier plan factoriel ?

Ce sont ceux de plus petite (resp. grande) hauteur.