

# TP2: Statistiques descriptives bivariées

Tanguy ROUDAUT — Tadios QUINIO

FIPASE 24

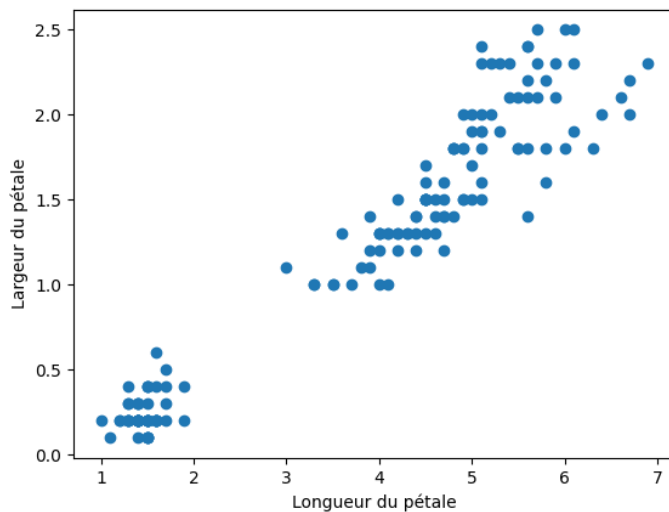
13 Septembre 2022

## 1 Statistiques descriptives bivariées sur des données d'Iris

### 1.1 Étude de la largeur du pétale en fonction de la longueur du pétale

#### Représentation graphique

**Question 1 :** Tracer le nuage de points de la longueur du pétale en fonction de la largeur du pétale pour les 150 iris contenus dans les données ((numpy.)plot, (numpy.)scatter). Ne pas oublier de mettre des titres sur les axes. Décrire le nuage de points.



Dans un premier temps, on remarque deux concentrations principale.

La première a une longueur de pétale qui varie entre 1 cm et 2 cm pour une largeur d'environ 0,2 cm et 0,6 cm. La seconde, une longueur variant de 3 cm à 7 cm avec une largeur de 1 cm à 2,5 cm.

On constate également que la longueur de pétale varie proportionnellement à la largeur.

FIGURE 1 – Nuage de points de la longueur en fonction de la largeur du pétale des 150 Iris

```
1 plt.scatter(petallength, petalwidth)
2 plt.ylabel("Largeur du pétale")
3 plt.xlabel("Longueur du pétale")
4 plt.title("Nuage de points de la longueur en fonction de la largeur du pétale")
5 plt.show()
```

Listing 1 – Code Python pour calculer le coefficient de corrélation

**Question 2 :** Rappeler la définition du coefficient de corrélation et le calculer par la fonction ((numpy.)corrcoef)

Le coefficient de corrélation est une valeur qui permet de déterminer s'il existe une relation linéaire entre deux variables, ce qui veut dire que si ces variables sont corrélées alors elles sont liées.

```
1 corr_coef_matrix = np.corrcoef(petallength, petalwidth)
2 print("Matrice de corrélation:\n",corr_coef_matrix)
```

Listing 2 – Code Python pour calculer le coefficient de corrélation

```
1 Matrice de corrélation:
2 [[1.          0.9627571]
3  [0.9627571  1.          ]]
```

Listing 3 – Résultat du code

**Question 3 :** Donner l'équation de la droite de régression linéaire, créer une fonction permettant de la calculer à partir de deux variables X et Y et tracer la sur le même graphique

### 1. Résultat et analyse :

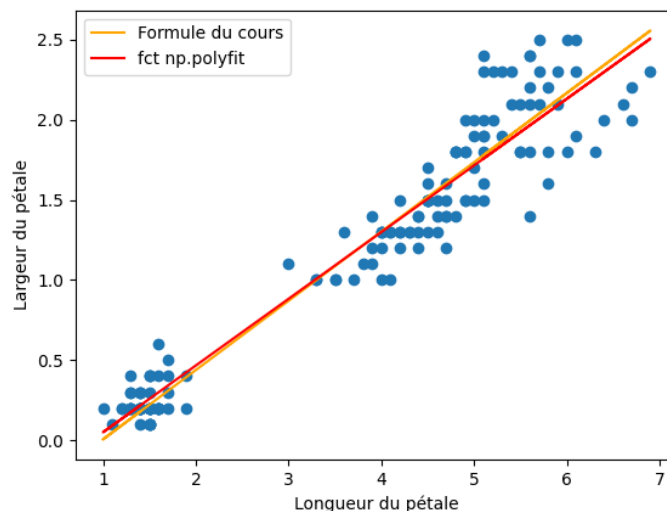


FIGURE 2 – Équation de la droite de régression linéaire sur le nuage de points

### 2. Formules utilisées :

$$\hat{a} = \rho_{X,Y} \cdot \frac{S_{n-1,Y}}{S_{n-1,X}} \quad \text{et} \quad \hat{b} = \bar{Y}_n \cdot \rho_{X,Y} \cdot \frac{S_{n-1,Y}}{S_{n-1,X}} \cdot \frac{1}{\bar{X}_n} \quad (1)$$

### 3. Code python :

```
1 # Méthode par la formule du cours
2 def regression_lineaire(corr_coef, x, y):
3     mean_x = np.mean(x)
4     var_x = np.var(x)
5     e_type_x = np.sqrt(var_x)
6
7     mean_y = np.mean(y)
8     var_y = np.var(y)
9     e_type_y = np.sqrt(var_y)
10
11     a = corr_coef * (e_type_y/e_type_x)
```

```

12     b = mean_y - corr_coef * ((e_type_y/e_type_x)) * mean_x
13
14     return a, b
15
16
17 a, b = regression_lineaire(corr_coef, petallength, petalwidth)
18
19 # Méthode avec polyfit de numpy
20 fit = np.polyfit(petallength, petalwidth, 1)
21 poly = petallength*fit[0]+fit[1]
22
23 plt.plot(petallength, a*petallength+b, color='orange', label='Formule du cours'
24 )
25 plt.plot(petallength, poly, color='red', label='fct np.polyfit')
26 plt.scatter(petallength, petalwidth)
27 plt.ylabel("Largeur du pétale")
28 plt.xlabel("Longueur du pétale")
29 plt.title("Equation de la droite de régression linéaire")
30 plt.legend()
31 plt.show()

```

Listing 4 – Code Python pour tracer la droite de régression linéaire

#### Question 4 : Analyser le lien entre les deux variables

On constate que les variables *petalLength* et *petalWidth* sont liées, elles évoluent proportionnellement l'une avec l'autre.

On peut donc soumettre l'hypothèse que ces variables sont corrélées puisqu'il existe deux réels tels que  $Y = aX + b$ .

## 1.2 Étude de la longueur de pétale selon les différentes espèces

### Représentations graphiques

**Question 5 :** Représenter sur une même figure, les trois histogrammes de la longueur de pétale : un histogramme pour chaque espèce avec une couleur. Commenter cette figure.

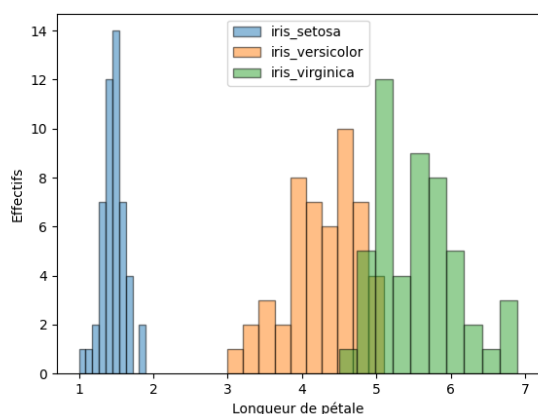


FIGURE 3 – Histogrammes de la longueur de pétale pour chaque espèce

Lors du *TP1 question 7*, nous avons pu remarquer que l'une des espèces a une taille de pétale plus faible, mais qui varie également moins que les deux autres. Le problème était que c'était seulement une hypothèse, puisque nous n'avions pas cette répartition entre les différentes espèces. L'histogramme de la figure 3, reprend le même schéma que celui de la *question 7 du TP1*, mais cette fois si en prenant en compte les différentes espèces.

On peut donc conclure que notre hypothèse est vraie, l'iris setosa est plus petite que les deux autres espèces.

```

1 iris_setosa_petallength = []
2 iris_versicolor_petallength = []
3 iris_virginica_petallength = []
4

```

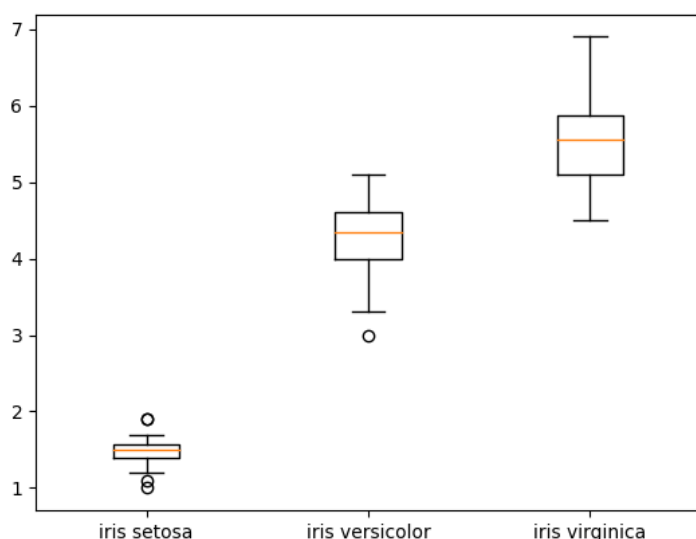
```

5 for i in range(len(species)):
6     if species[i] == 'Iris-setosa':
7         iris_setosa_petallength.append(petallength[i])
8     elif species[i] == 'Iris-versicolor':
9         iris_versicolor_petallength.append(petallength[i])
10    elif species[i] == 'Iris-virginica':
11        iris_virginica_petallength.append(petallength[i])
12
13 plt.hist(iris_setosa_petallength, edgecolor='black', alpha=0.5, label='iris_setosa')
14 plt.hist(iris_versicolor_petallength, edgecolor='black', alpha=0.5, label='
iris_versicolor')
15 plt.hist(iris_virginica_petallength, edgecolor='black', alpha=0.5, label='
iris_virginica')
16 plt.xlabel('Longueur de pétale')
17 plt.ylabel('Effectifs')
18 plt.title("Histogrammes de la longueur de pétale pour chaque espèce")
19 plt.legend()
20 plt.show()

```

Listing 5 – Code Python pour calculer le coefficient de corrélation

**Question 6 :** Représenter sur une même figure. une boîte à moustache par espèce. Commenter cette figure.



On remarque encore une fois que l'iris setosa à un intervalle de longueur de pétales plus faible que les deux autres espèces.

FIGURE 4 – Boîte à moustache de la longueur de pétale pour chaque espèce

```

1 boxplot_petallength = [iris_setosa_petallength, iris_versicolor_petallength,
2   iris_virginica_petallength]
3 plt.boxplot(boxplot_petallength, labels=['iris setosa', 'iris versicolor', 'iris
4   virginica'])
5 plt.title("Boîte à moustache de la longueur de pétale pour chaque espèce")
6 plt.show()

```

Listing 6 – Code Python pour calculer le coefficient de corrélation

## Représentations graphiques

**Question 7 :** Calculer le rapport de corrélation lié à la décomposition de la variance en variance intraclasse et interclasse. Qu'en concluez-vous ?

### 1. Formules utilisées :

– Variance interclasse :

$$S_B^2 = \frac{1}{n} \sum_{i=1}^p n_i (\bar{Y}_i - \bar{Y})^2 \quad (2)$$

– Variance intraclasse :

$$S_W^2 = \frac{1}{n} \sum_{i=1}^p n_i \cdot S_{n-1, Y_i}^2 \quad (3)$$

– Variance totale :

$$S_{n-1, Y}^2 = S_B^2 + S_W^2 \quad (4)$$

– Rapport de corrélation :

$$S_{\frac{Y}{X}} = \sqrt{\frac{S_B^2}{S_{n-1, Y}^2}} \in [0, 1] \quad (5)$$

### 2. Valeurs obtenues :

Variance interclasse	2.910
Variance intraclasse	0.181
Variance total	3.092
Rapport de corrélation	0.970

### 3. Conclusion :

La variance totale trouvée grâce à la formule 4, correspond avec celle trouvée avec *numpy*. Cette constatation nous permet de nous assurer que nos résultats sont corrects et de trouver par la suite le rapport de corrélation qui est de 0,970.

Le rapport de corrélation étant proche de 1 nous permet de confirmer notre hypothèse établie à la *question 4* concernant la corrélation des variables *petalLength* et *petalWidth*.

### 4. Code python :

```

1 n_cumul = len(species)
2
3 mean_species_petallength = [np.mean(iris_setosa_petallength), np.mean(
4     iris_versicolor_petallength), np.mean(iris_virginica_petallength)]
5
6 mean_petallength = np.mean(petallength)
7
8 var_species_petallength = [np.var(iris_setosa_petallength), np.var(
9     iris_versicolor_petallength), np.var(iris_virginica_petallength)]
10
11
12
13 def var_inter_classe_petallength(n_cumul, n_species, mean_species, mean):
14     ret = 0
15     for i in range(len(speciesname)):
16         ret += n_species[i] * (mean_species[i] - mean)**2
17
18     return ret/n_cumul
19
20

```

```
21
22 def var_intra_classe_petallength(n_cumul, n_species, var_species):
23     ret = 0
24     for i in range(len(speciesname)):
25         ret += n_species[i] * var_species[i]
26
27     return ret / n_cumul
28
29
30
31 var_inter_classe = var_inter_classe_petallength(n_cumul, n_species,
32     mean_species_petallength, mean_petallength)
33
34 var_intra_classe = var_intra_classe_petallength(n_cumul, n_species,
35     var_species_petallength)
36
37 var_totale_formule = var_inter_classe + var_intra_classe
38
39 var_totale = np.var(petallength)
40
41 rapport_corr = np.sqrt(var_inter_classe/var_totale_formule)
42
43 print("Variance interclasse: ", var_inter_classe)
44 print("Variance intraclasse: ", var_intra_classe)
45 print("Variance total trouvé avec la formule: ", var_totale_formule)
46 print("Variance total trouvé avec numpy: ", var_totale)
47 print("Rapport de corrélation: ", rapport_corr)
```

Listing 7 – Code Python pour calculer le coefficient de corrélation

```
1 Variance interclasse:  2.910958222222223
2 Variance intraclasse:  0.18146666666666667
3 Variance total trouvé avec la formule:  3.0924248888888894
4 Variance total trouvé avec numpy:  3.0924248888888889
5 Rapport de corrélation:  0.9702159417163924
6
```

Listing 8 – Résultat du code

## 1.3 Code complet

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4
5 df = pd.read_csv('iris.csv', sep=',')
6
7 sepallength = df["sepallength"].values
8 sepalwidth = df["sepalwidth"].values
9 petallength = df["petallength"].values
10 petalwidth = df["petalwidth"].values
11 species = df["class"].values
12 speciesname = np.unique(species)
13 variablename = df.keys().values[0:-1]
14
15
16 # question 1
17 plt.scatter(petallength, petalwidth)
18 plt.ylabel("Largeur du pétale")
19 plt.xlabel("Longueur du pétale")
20 plt.title("Nuage de points de la longueur en fonction de la largeur du pétale")
21 plt.show()
22
23
24 # question 2
25 corr_coef_matrix = np.corrcoef(petallength, petalwidth)
26 print(corr_coef_matrix)
27
28 corr_coef = corr_coef_matrix[1][1]
29
30
31 # question 3
32 def regression_lineaire(corr_coef, x, y):
33     mean_x = np.mean(x)
34     var_x = np.var(x)
35     e_type_x = np.sqrt(var_x)
36
37     mean_y = np.mean(y)
38     var_y = np.var(y)
39     e_type_y = np.sqrt(var_y)
40
41     a = corr_coef * (e_type_y/e_type_x)
42     b = mean_y - corr_coef * ((e_type_y/e_type_x)) * mean_x
43
44     return a, b
45
46
47 a, b = regression_lineaire(corr_coef, petallength, petalwidth)
48
49 fit = np.polyfit(petallength, petalwidth, 1)
50 poly = petallength*fit[0]+fit[1]
51
52 plt.plot(petallength, a*petallength+b, color='orange', label='Formule du cours')
53 plt.plot(petallength, poly, color='red', label='fct np.polyfit')
54 plt.scatter(petallength, petalwidth)
55 plt.ylabel("Largeur du pétale")
56 plt.xlabel("Longueur du pétale")
57 plt.title("Equation de la droite de régression linéaire")
58 plt.legend()
59 plt.show()
60
61
62 #question 5
63 iris_setosa_petallength = []
64 iris_versicolor_petallength = []

```

```

65 iris_virginica_petallength = []
66
67 for i in range(len(species)):
68     if species[i] == 'Iris-setosa':
69         iris_setosa_petallength.append(petallength[i])
70     elif species[i] == 'Iris-versicolor':
71         iris_versicolor_petallength.append(petallength[i])
72     elif species[i] == 'Iris-virginica':
73         iris_virginica_petallength.append(petallength[i])
74
75 plt.hist(iris_setosa_petallength, edgecolor='black', alpha=0.5, label='iris_setosa')
76 plt.hist(iris_versicolor_petallength, edgecolor='black', alpha=0.5, label='
iris_versicolor')
77 plt.hist(iris_virginica_petallength, edgecolor='black', alpha=0.5, label='
iris_virginica')
78 plt.xlabel('Longueur de pétale')
79 plt.ylabel('Effectifs')
80 plt.title("Histogrammes de la longueur de pétale pour chaque espèce")
81 plt.legend()
82 plt.show()
83
84
85 # question 6
86 boxplot_petallength = [iris_setosa_petallength, iris_versicolor_petallength,
iris_virginica_petallength]
87 plt.boxplot(boxplot_petallength, labels=['iris setosa', 'iris versicolor', 'iris
virginica'])
88 plt.title("Boite à moustache de la longueur de pétale pour chaque espèce")
89 plt.show()
90
91
92 #question 7
93 n_cumul = len(species)
94
95 mean_species_petallength = [np.mean(iris_setosa_petallength), np.mean(
iris_versicolor_petallength), np.mean(iris_virginica_petallength)]
96
97 mean_petallength = np.mean(petallength)
98
99 var_species_petallength = [np.var(iris_setosa_petallength), np.var(
iris_versicolor_petallength), np.var(iris_virginica_petallength)]
100
101 n_species = [len(iris_setosa_petallength), len(iris_versicolor_petallength), len(
iris_virginica_petallength)]
102
103
104
105 def var_inter_classe_petallength(n_cumul, n_species, mean_species, mean):
106     ret = 0
107     for i in range(len(speciesname)):
108         ret += n_species[i] * (mean_species[i] - mean)**2
109
110     return ret/n_cumul
111
112
113
114 def var_intra_classe_petallength(n_cumul, n_species, var_species):
115     ret = 0
116     for i in range(len(speciesname)):
117         ret += n_species[i] * var_species[i]
118
119     return ret / n_cumul
120
121
122

```



```
123 var_inter_classe = var_inter_classe_petallength(n_cumul, n_species,
    mean_species_petallength, mean_petallength)
124
125 var_intra_classe = var_intra_classe_petallength(n_cumul, n_species,
    var_species_petallength)
126
127 var_totale_formule = var_inter_classe + var_intra_classe
128
129 var_totale = np.var(petallength)
130
131 rapport_corr = np.sqrt(var_inter_classe/var_totale_formule)
132
133 print("Variance interclasse: ", var_inter_classe)
134 print("Variance intraclasses: ", var_intra_classe)
135 print("Variance total trouvé avec la formule: ", var_totale_formule)
136 print("Variance total trouvé avec numpy: ", var_totale)
137 print("Rapport de corrélation: ", rapport_corr)
```

Listing 9 – Code Python complet TP2