

Correction - TP 7

## Statistiques descriptives bivariées

### Exercice 1

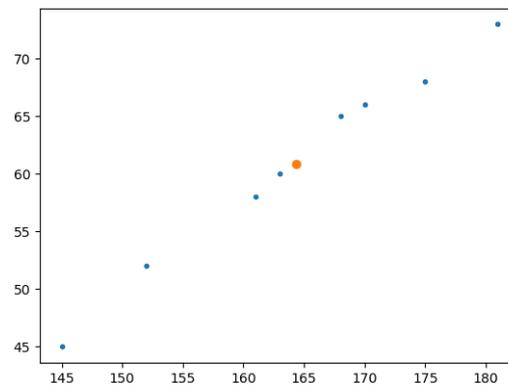
1. A priori, le poids dépend de la taille donc la variable explicative est  $X$  et la variable à expliquer est  $Y$ .
2. Voici les instructions pour tracer le nuage de points et le point moyen :

```

1 #Tracer du nuage de points
2 x = np.array([161, 170, 152, 181, 163, 145, 168, 175])
3 y = np.array([58, 66, 52, 73, 60, 45, 65, 68])
4 plt.plot(x, y, ".")
5
6 #Tracer du point moyen
7 plt.plot(np.mean(x), np.mean(y), "o")
8
9 plt.show()

```

On obtient le tracer suivant :



On remarque qu'il y a une forte corrélation linéaire entre la taille  $X$  et le poids  $Y$ , c'est-à-dire qu'il y a une relation du type  $Y \simeq aX + b$ .

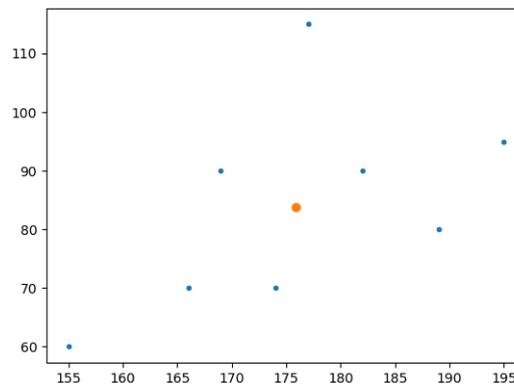
3. Voici les instructions pour tracer le nuage de points et le point moyen :

```

1 #Tracer du nuage de points
2 x = np.array([169, 195, 177, 182, 166, 155, 189, 174])
3 y = np.array([90, 95, 115, 90, 70, 60, 80, 70])
4 plt.plot(x, y, ".")
5
6 #Tracer du point moyen
7 plt.plot(np.mean(x), np.mean(y), "o")
8
9 plt.show()

```

On obtient le tracer suivant :



Il est difficile de voir apparaître une fonction qui relie  $X$  et  $Y$ . On remarque qu'il y a tout de même une variation selon une certaine direction oblique, on a donc  $Y \simeq aX + b$  mais cette corrélation linéaire est faible.

4. Aux USA, il y a un pourcentage important de personnes au comportement alimentaire déréglé... Ceci explique la faible corrélation linéaire entre la taille et le poids.

## Exercice 2

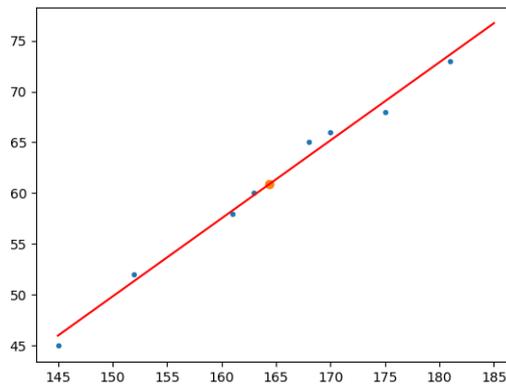
1. • Pour l'exemple du Japon :

```

1 #Tracer du nuage de points
2 x = np.array([161, 170, 152, 181, 163, 145, 168, 175])
3 y = np.array([58, 66, 52, 73, 60, 45, 65, 68])
4 plt.plot(x, y, ".")
5
6 #Tracer du point moyen
7 plt.plot(np.mean(x), np.mean(y), "o")
8
9 #Coefficients a et b de la droite de régression
10 s = np.mean(x.*y)-np.mean(x)*np.mean(y)
11 a = s/np.var(x)
12 b = np.mean(y)-a*np.mean(x)
13
14 #Coefficient de corrélation linéaire
15 r = s/(np.std(x)*np.std(y))
16 print(r)
17
18 #Tracer de la droite de régression linéaire
19 def f(t):
20     return(a*t+b)
21
22 t = np.linspace(145,185,100)
23 plt.plot(t, f(t), color="red")
24
25 plt.show()

```

On obtient le tracer suivant :



On utilise l'équation de la droite de régression linéaire :

```
>>> f(180)
72.873447
```

On peut ainsi estimer qu'un individu mesurant 1m80 pèse environ 72,9 kg. Étant donné que la corrélation linéaire entre  $X$  et  $Y$  est forte, cette analyse prédictive paraît pertinente.

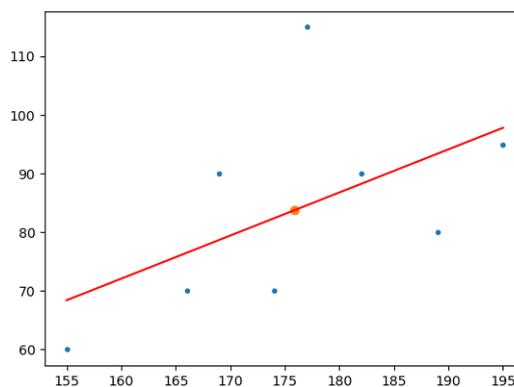
- Pour l'exemple des USA :

```

1 #Tracer du nuage de points
2 x = np.array([169, 195, 177, 182, 166, 155, 189, 174])
3 y = np.array([90, 95, 115, 90, 70, 60, 80, 70])
4 plt.plot(x, y, ".")
5
6 #Tracer du point moyen
7 plt.plot(np.mean(x), np.mean(y), "o")
8
9 # Coefficients a et b de la droite de régression
10 s = mean(x.*y)-mean(x)*mean(y)
11 a = s/np.var(x)
12 b = np.mean(y)-a*np.mean(x)
13
14 # Coefficient de corrélation linéaire
15 r = s/(np.std(x)*np.std(y))
16 print(r)
17
18 #Tracer de la droite de régression linéaire
19 def f(t):
20     return(a*t+b)
21
22 t = np.linspace(155,195,100)
23 plt.plot(t, f(t), color="red")
24
25 plt.show()

```

On obtient le tracer suivant :



On utilise l'équation de la droite de régression linéaire :

```
>>> f(180)
      86.783676
```

On peut ainsi estimer qu'un individu mesurant 1m80 pèse environ 86,8 kg. Étant donné que la corrélation linéaire entre  $X$  et  $Y$  est faible, cette analyse prédictive ne paraît pas pertinente.

2. (a) Dans les deux cas, on obtient un coefficient directeur  $a$  de signe positif ( $a = 0,76$  dans le cas du Japon et  $a = 0,73$  dans le cas des USA). Ceci s'interprète par le fait que les caractères  $X$  et  $Y$  ont tendance à varier dans le même sens. Cela paraît conforme à l'intuition, le poids et la taille aillant tendance à varier effectivement dans le même sens.
- (b) Dans le cas du Japon,  $r_{x,y} = 0.9953155$ . La corrélation linéaire entre les séries  $X$  et  $Y$  est donc forte. Cela correspond à ce qu'on s'attendait étant donné l'allure du nuage de points.  
Dans le cas des USA,  $r_{x,y} = 0.5419821$ . La corrélation linéaire entre les séries  $X$  et  $Y$  est donc faible. Cela correspond à ce qu'on s'attendait étant donné l'allure du nuage de points.

### Exercice 3

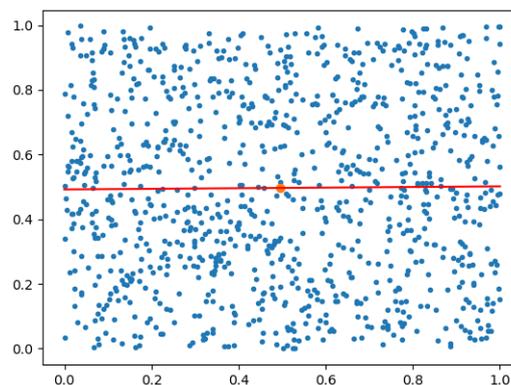
1. On utilise le code suivant :

```

1 #Tracer du nuage de points
2 x = rd.random(1000)
3 y = rd.random(1000)
4 plt.plot(x, y, ".")
5
6 #Tracer du point moyen
7 plt.plot(np.mean(x), np.mean(y), "o")
8
9 #Tracer de la droite de régression linéaire
10 s = mean(x.*y)-mean(x)*mean(y)
11 a = s/np.var(x)
12 b = np.mean(y)-a*np.mean(x)
13
14 def f(t):
15     return(a*t+b)
16
17 t = np.linspace(0,1,10)
18 plt.plot(t, f(t), color="red")
19
20 plt.show()
21
22 # Coefficient de corrélation linéaire
23 r = s/(np.std(x)*np.std(y))
24 print(r)

```

On obtient le graphe suivant :



La droite de régression linéaire est quasi horizontale. Elle ne semble pas du tout adaptée pour approximer notre nuage de points, il ne semble pas y avoir de corrélation linéaire entre les deux séries statistiques.

Le coefficient de corrélation linéaire obtenu est 0.0181605. Il confirme l'indépendance linéaire qu'on avait constaté auparavant. Il explique également pourquoi la droite de régression linéaire est horizontale. Rappelons en effet que le coefficient directeur de cette droite fait intervenir la covariance, qui dans notre cas est voisin de 0.

Tous ces résultats ne sont pas étonnants : les séries statistiques  $x$  et  $y$  ont été choisies de manière indépendante l'une de l'autre, la fonction `rd.random` générant des réalisations indépendantes de la loi  $\mathcal{U}([0,1])$ . Elles sont en particulier linéairement indépendantes, ce qui implique que leur coefficient de corrélation linéaire est proche de 0.

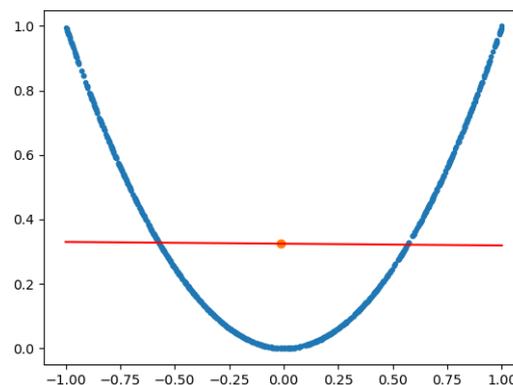
2. On adapte le code précédent :

```

1 #Tracer du nuage de points
2 x = rd.uniform(-1, 1, 1000)
3 y = x**2
4 plt.plot(x, y, ".")
5
6 #Tracer de la droite de régression linéaire
7 s = mean(x.*y)-mean(x)*mean(y)
8 a = s/np.var(x)
9 b = np.mean(y)-a*np.mean(x)
10
11 def f(t):
12     return(a*t+b)
13
14 t = np.linspace(-1,1,10)
15 plt.plot(t, f(t), color="red")
16
17 plt.show()
18
19 # Coefficient de corrélation linéaire
20 r = s/(np.std(x)*np.std(y))
21 print(r)

```

On obtient le graphe suivant :



La droite de régression linéaire est là aussi quasi horizontale (ce qui indique un coefficient de corrélation linéaire proche de 0) et ne semble également pas du tout adaptée pour approximer notre nuage de points. Le coefficient de corrélation linéaire obtenu le confirme : il est de  $-0.0056861$ .

On peut conclure que les séries statistiques  $x$  et  $y$  sont linéairement indépendantes. Mais attention à ne pas conclure à l'indépendance de ces séries : elles ne le sont clairement pas puisque  $y = x^2$ .

---

#### Exercice 4

1. et 2. On utilise le code suivant :

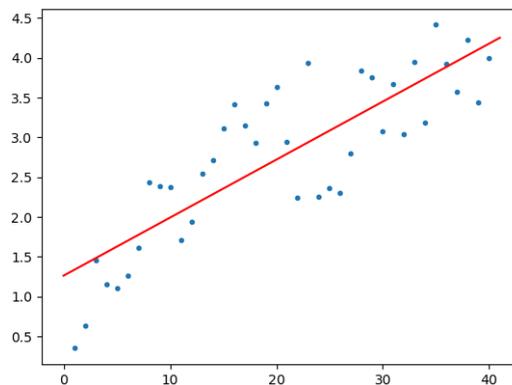
```

1 #Tracer du nuage de points
2 x = np.arange(1, 41)
3 y = np.log(x)+rd.uniform(-1, 1, 40)
4 plt.plot(x, y, ".")
5
6 #Tracer de la droite de régression linéaire
7 s = mean(x.*y)-mean(x)*mean(y)
8 a = s/np.var(x)
9 b = np.mean(y)-a*np.mean(x)
10
11 def f(t):
12     return(a*t+b)
13
14 t = np.linspace(0,41,10)
15 plt.plot(t, f(t), color="red")
16
17 plt.show()
18
19 # Coefficient de corrélation linéaire
20 r = s/(np.std(x)*np.std(y))
21 print(r)

```

On obtient 0.7537838 comme coefficient de corrélation linéaire, bien en deçà de la barre des 0.9. La corrélation linéaire entre les deux séries statistiques n'est donc pas bonne.

On obtient la représentation suivante du nuage de points et de la droite de régression linéaire associée :



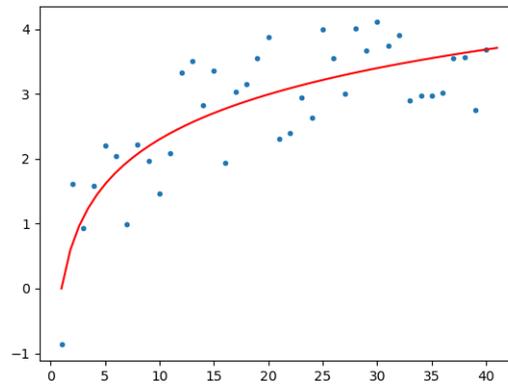
3. Avec le code :

```

1 #Tracer du nuage de points
2 x = np.arange(1, 41)
3 y = np.log(x)+rd.uniform(-1, 1, 40)
4 plt.plot(x, y, ".")
5
6 #Courbe du log
7 t = linspace(1,41,50)
8 plt.plot(t, np.log(t), color="red")

```

on obtient :



Le nuage de point se superpose effectivement bien avec la courbe représentative du logarithme.

4. En exécutant les commandes

```

1 //Coefficient de corrélation linéaire
2 z = np.log(x)
3 r = s/(np.std(z)*np.std(y))
4 print(r)

```

on obtient 0.8377479 comme coefficient de corrélation linéaire. C'est proche des 0.9, ce qui confirme une forte corrélation linéaire entre  $\log(x)$  et  $y$ .

### Exercice 5

1. Avec le code :

```

1 x = np.array([15.85, 15.7, 15.39, 15.32, 14.85])
2 y = np.array([230, 280, 360, 410, 525])
3 s = np.mean(x*y)-np.mean(x)*np.mean(y)
4 r = s/(np.std(x)*np.std(y))
5 print(r)

```

on obtient un coefficient de corrélation linéaire de -0.9951611.

2. Ces deux séries statistiques sont donc fortement corrélées linéairement. Pourtant, il semble bien difficile d'envisager un lien rationnel entre le nombre de morts sur les routes et les importations de citrons mexicains !