

# Probab-Stats avec Python

## Remarque

• Avec les données en exemple qui suivent et les applications numériques (#AN), vous de-

vriez voir quelles sont les commandes correspondantes sur vos calculatrices.

## I IMPORTATIONS

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from scipy.stats import linregress
```

## II DONNÉES EN EXEMPLE

x =	-1.61	-1.20	-0.97	-0.51	-0.42
y =	2.22	2.27	2.38	2.60	2.65

```
1 X = [-1.61, -1.20, -0.97, -0.51, -0.42]
2 Y = [2.22, 2.27, 2.38, 2.60, 2.65]
```

## III STATISTIQUES UNIVARIÉES

```
1 # Moyenne :
2 mx = np.mean(X)           #AN : -0.94146716215
3
4 # Variance des données :
5 s2x = np.var(X, ddof=0)   #AN : 0.195570714089
6
7 # Écart-type des données :
8 sx = np.std(X, ddof=0)    #AN : 0.442233777644
9
10 # Variance non biaisée :
11 np.var(X, ddof=1)        #AN : 0.244463392612
12
13 # Écart-type non biaisé :
14 np.std(X, ddof=1)       #AN : 0.494432394379
```

### Rem.1 |

a. Python calcule de la façon suivante :

$$\text{np. std}(X, \text{ddof}=D) = \frac{1}{n-D} \sum_{i=1}^n (X[i] - \bar{x})^2,$$

qui vaut bien  $s_x^2$  quand  $\text{ddof}=0$ , et la variance non biaisée si  $\text{ddof}=1$ .

b.  $\text{ddof}$  signifie : delta degrees of freedom.

c. On voit encore sur cet exemple le fait que la variance (resp. l'écart-type) empirique sous-estime systématiquement la variance (resp. l'écart-type) théorique.

## IV STATISTIQUES BIVARIÉES

```
1 # Covariance de la série :
2 M = np.cov(X, Y, ddof=0)
3
4 #AN :
5 array([[ 0.19557071,  0.07506966],
6        [ 0.07506966,  0.03029122]])
7
8 sxy = M[0,1]           #AN : 0.07506966
9
10 # Coefficient de corrélation : voir V. A) !
```

**Rem.2 |** On obtient pour la covariance non pas un flottant, mais une matrice (symétrique réelle). C'est la matrice :

$$M = \begin{pmatrix} s_x & s_{xy} \\ s_{xy} & s_y \end{pmatrix},$$

dite *matrice de covariance*. Par lecture du tableau M :

```
1 # sx = 0.19557071
2 # sy = 0.03029122
3 # sxy = 0.07506966
```

## V RÉGRESSION LINÉAIRE

### A) COEFFICIENTS DE LA DROITE DE RÉGRESSION . . . . .

Les trois premiers arguments de sortie de la fonction `linregress` sont les coefficients  $a, b$  de la droite de régression, ainsi que le coefficient de corrélation  $r = r_{x,y}$ .

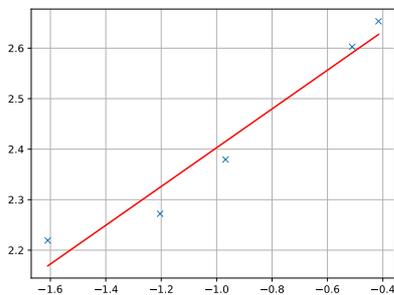
```
1 a, b, r = linregress(X, Y)[:3]
2
3 #AN :
4 # a = 0.38384918507505805
5 # b = 2.7867428337153797
6 # r = 0.9753357022963296
```

## B) TRACÉS DU NUAGE DE POINTS ET DE LA DROITE DE RÉGRESSION .....

```

1 # Nuage de points
2 plt.plot(X,Y, 'x')
3 # Droite de régression
4 x = np.array(X)
5 y = a*x+b # avantage de np.array sur list !
6 plt.plot(x,y, 'r')
7 plt.grid('on')
8 plt.show()

```

VI AUTOUR DE LA LOI  $\mathcal{N}(m, \sigma^2)$ 

## A) EXEMPLE .....

**Exple.1** | Dans une population, il y a 20% de malades. On prélève indépendamment un échantillon de 500 individus.

- Calculer la probabilité  $p_1$  que l'échantillon contienne 80 malades
- Calculer la probabilité  $p_2$  que l'échantillon contienne entre 50 et 150 malades.

## B) SOLUTION .....

La variable aléatoire  $S$  qui compte le nombre de malades dans l'échantillon de 500 individus suit une loi  $\mathcal{B}(N, p)$  où  $N = 500$ , et  $p = 0,2$ . Les hypothèses du théorème de Moivre-Laplace sont vérifiées ici, puisque :

- $N = 500 \geq 30$
- $Np = 100 \geq 5$
- $Nq = 400 \geq 5$

On peut donc remplacer la loi de  $S$  par la loi d'une variable aléatoire  $Z$  de loi  $\mathcal{N}(E(S), V(S))$ . On obtient donc les réponses suivantes :

- À l'aide d'une correction de continuité pour commencer :

$$\begin{aligned}
 p_1 = P(S = 80) &= P(S \in [79,5; 80,5]) \\
 &\simeq P(Z \in [79,5; 80,5]) \\
 &= F_Z(80,5) - F_Z(79,5),
 \end{aligned}$$

où  $F_Z$  est la fonction de répartition de la variable  $Z$ .

- De même :

$$\begin{aligned}
 p_2 = P(S \in [100, 150]) &\simeq P(Z \in [100; 150]) \\
 &= F_Z(150) - F_Z(100),
 \end{aligned}$$

## C) CALCULS SOUS PYTHON .....

Ces deux dernières expressions approchées des probabilités  $p_1, p_2$  se calculent facilement sous Python puisque la fonction  $F_Z$  y est implémentée :

- La commande `norm.cdf(t,m,sigma)` renvoie  $F_Z(t)$  pour une variable aléatoire  $Z$  de loi  $\mathcal{N}(m, \sigma^2)$  :

$$F_Z(t) = \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^t \exp\left(-\frac{(x-m)^2}{2\sigma^2}\right) dx.$$

- Attention, c'est bien  $\sigma$ , et pas  $\sigma^2$  qui est pris en argument de la fonction `cdf` !
- Par défaut,  $(m, \sigma) = (0, 1)$  : ainsi `norm.cdf(t)` renvoie la valeur de  $\Phi(t)$ .
- `cdf` signifie : cumulative distribution function.

```

1 N,p = 500, 0.2 # Param. de la loi binomiale
2 q = 1-p
3 def Fz(t):
4     m,s = N*p,np.sqrt(N*p*q) # évite d'avoir à
5     return norm.cdf(t,m,s) # retaper m,s
6
7 # Réponse à la question a.
8 p1 = F_z(80.5) - Fz(79.5)
9     #AN : 0.003668872638773692
10
11 # Réponse à la question b.
12 p2 = F_z(150) - Fz(50)
13     #AN : 0.99999997731525148

```

**Rem.3** | Sur cet exemple, l'écart type de  $S$  est  $\sigma_s \simeq 8.94$ , et  $m = E(S) = Np = 100$ . Or  $[100; 150] = [m - 50; m + 50]$ , et  $50 \simeq 5.6\sigma_s$ . Il est donc tout à fait normal (!) que la valeur de  $p_2$  soit aussi proche de 1 (en effet, d'après les propriétés de la loi normale, ce serait déjà le cas pour un intervalle de rayon  $3\sigma_s$  centré sur  $m$ ).