

En attendant 1, 2, 3 sigma ...

Compétences visées

- **chercher**, expérimenter – en particulier à l'aide d'outils logiciels ;
- **calculer**, appliquer des techniques et mettre en œuvre des algorithmes ;
- **communiquer** un résultat par oral ou par écrit, expliquer une démarche.

Ces compétences sont mises en œuvre au regard de l'extrait du programme de 2^{nde} GT ci-dessous :

« Pour des données réelles ou issues d'une simulation, lire et comprendre une fonction écrite en Python renvoyant la moyenne m , l'écart type s , et la proportion d'éléments appartenant à $[m - 2s, m + 2s]$. »

Situation déclenchante

Pour bien comprendre les formules statistiques et prendre un peu de recul ...

Comment distinguer dans une série numérique de grande taille représentant des données (a priori réelles) les valeurs aberrantes des valeurs extrêmes ?

Problématique

Créer des fonctions en python permettant de :

- calculer la moyenne d'une série de nombres ;
- calculer l'écart type d'une série de nombres ;
- calculer la proportion de valeurs dont la distance à la moyenne est inférieure à deux écarts types.

Fiche méthode

Proposition de résolution

Pour éviter d'avoir à saisir des listes de nombres à la main, on crée au préalable une fonction permettant de générer une liste de n nombres aléatoires entiers à partir de la fonction **randint**. Cependant, une liste obtenue directement à partir **randint** donnerait un résultat peu intéressant comme cela est montré à la fin de la fiche. C'est pourquoi, pour rendre le travail pertinent, on est amené à créer des listes de nombres de manière plus fine, comme décrit dans les étapes de résolution qui suivent.

- Sur la copie d'écran ci-contre, l'instruction **liste(0,5,20)** renvoie une liste de 20 nombres entiers compris entre 3×0 et 3×5 ; cette liste est stockée dans la variable **li**.

li est ensuite affichée (c'est à déconseiller par la suite lorsqu'on aura beaucoup de valeurs).

liste2(0,5,20,4) génère quant à elle une liste de 20 nombres entiers avec un mode plus « élaboré » décrit plus tard.

- Écrire la fonction **moy** et la fonction **ecty** qui prennent toutes deux comme paramètre une liste et qui renvoient respectivement la valeur moyenne et l'écart type de des valeurs de la liste
- Écrire une fonction **prop** qui prend en paramètre une liste et qui renvoie la proportion de valeurs de la liste qui sont à une distance inférieure à deux écarts types de la moyenne.
- Pour effectuer plusieurs essais, il est commode d'utiliser la flèche directionnelle vers le haut : un appui réécrit la dernière instruction utilisée (deux appuis donne l'instruction écrite deux lignes plus haut et ainsi de suite) ; en appuyant sur **entrer**, on exécute à nouveau l'instruction.

```
PYTHON SHELL
>>> li=liste(0,5,20)
>>> li
[4, 7, 10, 12, 8, 13, 4, 10, 4,
9, 10, 6, 1, 5, 4, 10, 12, 5, 3,
9]
>>> li2=liste2(0,5,20,4)
>>> li2
[4, 13, 12, 9, 12, 16, 6, 9, 11,
11, 12, 10, 16, 11, 16, 7, 10,
9, 7, 10]
>>> |
Fns... a A # Outils Éditer Script
```

```
PYTHON SHELL
>>> li=liste(1,5,500)
>>> moy(li)
8.891999999999999
>>> ecty(li)
2.50605985562995
>>> li=liste(1,5,500)
>>> moy(li)
8.9
>>> li=liste(1,5,500)
>>> |
Fns... a A # Outils Éditer Script
```

```
PYTHON SHELL
>>> prop(li)
0.9619999999999999
>>> li=liste2(0,10,500,4)
>>> prop(li)
0.966
>>> li=liste2(0,10,500,4)
>>> prop(li)
0.958
>>> li=liste2(0,10,500,4)
>>> prop(li)
0.96
Fns... a A # Outils Éditer Script
```

Remarque

Importation en préambule du code de la bibliothèque « random » par « **from random import *** » pour pouvoir utiliser la fonction **randint**

Importation de la bibliothèque « math » par « **from math import *** » pour pouvoir utiliser les fonctions **sqrt** (racine carrée) et **fabs** (valeur absolue)

```
ÉDITEUR : STATS
LIGNE DU SCRIPT 0011
from random import *
from math import *
```

Pour profiter de tutoriels vidéos, Flasher le QRCode ou cliquer dessus



Fiche méthode

Etapes de résolution

L'instruction **liste(a,b,n)** génère une liste de n valeurs entières comprises entre $3a$ et $3b$.

li = [] définit une variable de type list et l'initialise en une liste vide.

li.append() permet d'ajouter à la liste **li** la valeur passée en paramètre: ici, on ajoute $\text{randint}(a,b) + \text{randint}(a,b)$ à la liste **li**.

La fonction **moy** a comme paramètre une liste et renvoie la moyenne des valeurs contenues dans la liste.

sum est une fonction native qui prend une liste en paramètre et renvoie la somme des valeurs contenues dans la liste.

len est une fonction native qui prend une liste en paramètre et renvoie le nombre de valeurs de la liste.

La fonction **ecty** a pour paramètre une liste et renvoie l'écart type de la série de valeurs contenues dans la liste.

La fonction **ecty** utilise la fonction **moy** précédemment définie, ce qui rend son écriture proche de la formule du cours.

for x in liste permet à **x** de parcourir les valeurs contenue dans la liste nommée **liste**.

Il faut comprendre le rôle joué par la variable **v** (v pour variance) qui ajoute au fur et à mesure que l'on parcourt la liste les valeurs notées usuellement $(x_i - m)^2$.

La fonction **prop** a comme paramètre une liste et renvoie la proportion des valeurs de la liste situées à une distance supérieure à deux écarts type de la moyenne.

Pour des raisons de lisibilité, on a choisi de noter: **n=len(liste)**, **m=moy(liste)**, **s=ecty(liste)**, ce qui permet d'avoir un test écrit de manière proche de l'écriture naturelle :

« si $| \text{valeur} - m | \leq 2s$, alors ajouter 1 à c »

```
ÉDITEUR : STATS
LIGNE DU SCRIPT 0010
from random import *
from math import *

def liste(a,b,n):
    li=[]
    for i in range(n):
        li.append(randint(a,b)+randint(a,b))
    return li
```

```
ÉDITEUR : STATS
LIGNE DU SCRIPT 0010
def moy(liste):
    return sum(liste)/len(liste)

def ecty(liste):
    v=0
    n=len(liste)
    m=moy(liste)
    for x in liste:
        v=v+(x-m)**2
    v=v/n
    return sqrt(v)
```

```
ÉDITEUR : STATS
LIGNE DU SCRIPT 0030
def prop(liste):
    c=0
    n=len(liste)
    m=moy(liste)
    s=ecty(liste)
    for x in liste:
        if fabs(x-m)<=2*s:
            c=c+1
    return c/n
```

Pour profiter de tutoriels vidéos, Flasher le QRCode ou cliquer dessus



Fiche méthode

Retour à la situation déclenchante

Lorsqu'une liste est issue d'une situation réelle (valeurs données par un capteur, valeurs issues d'un sondage, etc.), il se peut que certaines soient erronées.

S'il est possible de remonter à la source, on pourra éventuellement corriger la valeur (erreur de saisie par exemple), sinon, il sera pertinent de l'éliminer.

Comment examiner ces valeurs ?

Il faut les traiter au cas par cas, et le travail précédent va permettre de le faire.

En effet, on examinera les valeurs inférieures à $m - 2s$ et celles supérieures à $m + 2s$ et on fera fonctionner son esprit critique ! Pour distinguer des valeurs acceptables qui seraient extrêmes des valeurs aberrantes.

Si par exemple un sondage traite de la taille des individus (en cm), la valeur 1.78 aura sans doute été saisie par quelqu'un répondant en m. On pourra la modifier.

liste(50,70,100) génère 100 nombres entiers compris entre 150 et 210 : ils peuvent simuler un échantillon de tailles d'adultes. On observe les valeurs extrêmes, hors de l'intervalle $[m - 2s; m + 2s]$.

```
ÉDITEUR : STATS
LIGNE DU SCRIPT 0031

def test(liste):
    examine=[]
    n=len(liste)
    m=moy(liste)
    s=ecty(liste)
    for x in liste:
        if fabs(x-m)>2*s:
            examine.append(x)
    return examine

Fns... a A # Outils Exéc Script
```

```
PYTHON SHELL

>>>
>>> li=liste(50,70,100)
>>> li.append(1.78)
>>> test(li)
[1.78]
>>> li=liste(50,70,100)
>>> test(li)
[158, 157, 202, 203, 158]
>>> |

Fns... a A # Outils Éditer Script
```

Remarque

On peut améliorer le processus permettant de générer une liste à partir de la fonction **randint** par la fonction ci-contre qui va ajouter r fois des nombres aléatoires entiers compris entre les entiers a et b .

La fonction **liste2(a,b,n,r)** génère une liste de n valeurs entières comprises entre $a \times r$ et $b \times r$ en ajoutant r fois un nombre issu de **randint(a,b)**.

```
ÉDITEUR : STATS
LIGNE DU SCRIPT 0019

def liste2(a,b,n,r):
    li=[]
    for i in range(n):
        val=0
        for i in range(r):
            val=val+randint(a,b)
        li.append(val)
    return li

Fns... a A # Outils Exéc Script
```

Pour profiter de tutoriels vidéos, Flasher le QRCode ou cliquer dessus



Fiche méthode

Un peu de théorie

Pourquoi avoir créé des fonctions pour générer des listes de valeurs ?

On aurait pu écrire le code ci-contre, générant une liste directement à partir de `randint` : on obtient n nombres entiers répartis 'uniformément' entre a et b .

Mais alors, toutes les valeurs appartiennent à l'intervalle $[m - 2s; m + 2s]$. On va le montrer ici.

Si X est une variable aléatoire qui suit la loi de probabilité suivante :

a_i	1	2	...	n
$P(X = a_i)$	$\frac{1}{n}$	$\frac{1}{n}$		$\frac{1}{n}$

Son espérance mathématique est : $E(X) = \frac{1}{n} \sum_{k=1}^n k = \frac{1}{n} \cdot \frac{n(n+1)}{2} = \frac{n+1}{2}$

La variance se calcule par :

$$V(X) = \frac{1}{n} \sum_{k=1}^n \left(k - \frac{n+1}{2}\right)^2 = \frac{1}{n} \left(\sum_{k=1}^n k^2 - (n+1) \sum_{k=1}^n k + \sum_{k=1}^n \frac{(n+1)^2}{4} \right)$$

$$V(X) = \frac{1}{n} \left(\frac{n(n+1)(2n+1)}{6} - 2(n+1) \frac{n(n+1)}{4} + n \frac{(n+1)^2}{4} \right) = \frac{n^2-1}{12}$$

Au final, l'écart type est égal à : $\sigma(X) = \sqrt{V(X)} = \sqrt{\frac{n^2-1}{12}}$

On peut confronter ces résultats à la moyenne et l'écart type obtenu à partir de `liste3(1,10)`.

Théoriquement, la moyenne vaut 5,5 et l'écart type $\sqrt{\frac{99}{12}} \approx 2,87$

Ainsi, une telle liste ne présente pas d'intérêt quant à la proportion d'éléments distants de moins de deux écarts-types de la moyenne.

En effet, pour $n \in \mathbb{N}^*$: $E(X) - 2\sigma(X) = \frac{n+1}{2} - 2\sqrt{\frac{n^2-1}{12}} \leq 1$

$$E(X) + 2\sigma(X) = \frac{n+1}{2} + 2\sqrt{\frac{n^2-1}{12}} \geq n$$

Ce qui montre que toutes les valeurs comprises entre 1 et n sont dans l'intervalle $[E(X) - 2\sigma(X); E(X) + 2\sigma(X)]$

C'est pourquoi on a créé les fonctions `liste` et `liste2` qui génèrent des listes qui n'auront pas systématiquement toutes leurs valeurs comprises entre $m - 2s$ et $m + 2s$.

```

EDITEUR : STATS
LIGNE DU SCRIPT 0029

def liste3(a,b,n):
    li=[]
    for i in range(n):
        li.append(randint(a,b))
    return li
    
```

```

PYTHON SHELL

>>> prop(li)
0.0
>>> li=liste3(1,10,5)
>>> prop(li)
1.0
>>> li=liste3(1,10,500)
>>> prop(li)
1.0
>>> li=liste3(1,10,500)
>>> prop(li)
1.0
    
```

```

PYTHON SHELL

>>> li=liste3(1,10,500)
>>> moy(li)
5.328
>>> ecty(li)
2.8327400163093
>>> li=liste3(1,10,500)
>>> moy(li)
5.546
>>> ecty(li)
2.874697201445743
>>> |
    
```

```

PYTHON SHELL

>>> l1=liste(1,10,500)
>>> prop(l1)
0.9719999999999999
>>> l1=liste(1,10,500)
>>> prop(l1)
0.952
>>> l1=liste(1,10,500)
>>> prop(l1)
0.968
>>> |
    
```

Pour profiter de tutoriels vidéos, Flasher le QRCode ou cliquer dessus

