

ECHANTILLONNAGE

▶ Vidéo <https://youtu.be/EXEcSJE31QY>

I. Notion d'échantillon

1) Définition

Exemples :

1) Sur l'ensemble des cartes à puce produites par une entreprise en une semaine, on en prélève 200. On dit que cet ensemble de 200 cartes à puce constitue un **échantillon de taille 200** de la population de toutes les cartes à puce produites en une semaine.

2) On s'intéresse aux intentions de vote lors d'une élection. On sonde 1000 personnes en leur demandant leur intention de vote. L'ensemble de ces 1000 personnes constitue un **échantillon de taille 1000** de la population totale des électeurs.

3) On lance une pièce de monnaie 50 fois de suite et on note les résultats obtenus. L'ensemble de ces 50 lancers constitue un **échantillon de taille 50**.

Définition :

Un **échantillon de taille n** est constitué des résultats de n répétitions indépendantes de la même expérience sur l'ensemble des personnes ou objets sur lesquels porte l'étude statistique (la population).

Un échantillon issu d'une population est donc l'ensemble de quelques éléments de cette population.

2) Simulation d'une expérience aléatoire

On considère l'expérience aléatoire qui consiste à lancer un dé à 6 faces. Le programme Python suivant permet de simuler cette expérience.

Le fonction **randint** renvoie un nombre aléatoire entier de 1 à 6.

```
from random import *

def dé():
    r=randint(1,6)
    return(r)
```

```
>>> dé()
1
```

On exécute le programme et on obtient l'affichage ci-contre. Cela signifie que le logiciel a simulé un lancer de dé et on a obtenu un « 1 ».



La règle du jeu veut que si le résultat est « 1 » ou « 6 », on gagne. Dans le cas contraire, on perd. On répète n fois de suite cette expérience à deux issues (gagner ou perdre) consistant à lancer le dé.

On modifie et complète le programme Python afin de simuler n lancers de dé. Le programme affiche le nombre de fois que l'on gagne.

La variable n désigne le nombre de lancers. La variable s permet de compter le nombre de fois que l'on gagne : le dé s'arrête sur « 1 » ou sur « 6 ».

```
from random import*

def dé(n):
    s=0
    for k in range(n):
        r=randint(1,6)
        if r==1 or r==6:
            s=s+1
    return(s)
```

```
>>> dé(10)
3
```

On exécute le programme et on obtient l'affichage ci-contre. Cela signifie que sur 10 lancers, on a gagné 3 fois.

II. Loi des grands nombres

Modifions le programme afin d'afficher en sortie la fréquence de jeux gagnés sur un échantillon de n lancers de dé.

Il suffit de remplacer dans la dernière ligne **return(s)** (l'effectif) par **return(s/n)** (la fréquence).

```
from random import*

def dé(n):
    s=0
    for k in range(n):
        r=randint(1,6)
        if r==1 or r==6:
            s=s+1
    return(s/n)
```

```
>>> dé(10)
0.2
>>> dé(100)
0.32
>>> dé(1000)
0.328
>>> dé(5000)
0.3372
>>> dé(100000)
0.33353
```

On exécute le programme pour des valeurs de n de plus en plus grandes. Ci-contre les résultats obtenus à l'aide du logiciel.

On constate que, plus n devient grand, plus les fréquences observées semblent se rapprocher d'une valeur théorique égale à $\frac{1}{3}$.

En effet, la probabilité de gagner (obtenir un « 1 » ou un « 6 ») est égale à $\frac{2}{6} = \frac{1}{3}$.

Loi des grands nombres : Lorsque n devient grand, sauf exception, la fréquence observée est proche de la probabilité.

III. Estimation d'une probabilité

On se propose maintenant de répéter N fois la simulation de l'expérience aléatoire précédente. Dans chaque cas, pour n suffisamment grand, la fréquence observée f devrait être proche de la probabilité théorique $p = \frac{1}{3}$.

On veut calculer la proportion des cas pour lesquels l'écart entre f et p est inférieur ou égale à $\frac{1}{\sqrt{n}}$.

Après avoir importé le module **math**, nécessaire pour utiliser la fonction **abs** (valeur absolue), on complète le programme précédent avec la fonction **estim**.

abs(f-1/3) est l'écart entre f et $1/3$.
 \sqrt{n} se note **sqrt(n)**.

On teste **N** fois si **abs(f-1/3) <= 1/sqrt(n)**.
La variable **c** compte le nombre de fois où ce test est vérifié.

```
from math import*

def estim(N,n):
    c=0
    for k in range(N):
        f=dé(n)
        if abs(f-1/3)<=1/sqrt(n):
            c=c+1
    return(c/N)
```

Le programme complet au format texte se trouve sur la dernière page de ce document.

On exécute le programme pour différentes valeurs de N en choisissant n suffisamment grand, soit $n = 10000$.

On trouve des valeurs proches de 0,95 ce qui signifie que dans 95% des cas, l'écart entre la fréquence observée f et la probabilité p est inférieur ou égale à 0,01.

En effet : $\frac{1}{\sqrt{n}} = \frac{1}{\sqrt{10000}} = 0,01$.

```
>>> estim(10,10000)
0.9
>>> estim(50,10000)
0.96
>>> estim(100,10000)
0.96
>>> estim(100,10000)
0.94
```

Principe de l'estimation : Pour n assez grand, f donne une bonne estimation de p dans environ 95 % des cas.

Le programme complet :

```

from random import*
from math import*

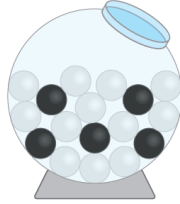
def dé(n):
    s=0
    for k in range(n):
        r=randint(1,6)
        if r==1 or r==6:
            s=s+1
    return(s/n)

def estim(N,n):
    c=0
    for k in range(N):
        f=dé(n)
        if abs(f-1/3)<=1/sqrt(n):
            c=c+1
    return(c/N)

```

IV. Estimation et échantillonnage

Dans ce paragraphe, on va étudier deux domaines des statistiques qu'il faut savoir distinguer :

Echantillonnage – Prise de décision	Estimation
<p>- Une urne contient un très grand nombre de boules blanches et de boules noires dont on connaît la proportion p de boules blanches. On tire avec remise n boules (échantillon) et on observe la fréquence d'apparition des boules blanches. Cette fréquence observée appartient à un intervalle, appelé intervalle de fluctuation de centre p.</p> <p>- Dans le cas où on ne connaît pas la proportion p mais on est capable de faire une hypothèse sur sa valeur, on parle de prise de décision. On veut par exemple savoir si un dé est bien équilibré. On peut faire l'hypothèse que l'apparition de chaque face est égale à $1/6$ et on va tester cette hypothèse à l'aide d'une expérience. Le résultat de l'expérience va nous permettre d'accepter ou rejeter l'hypothèse de départ.</p>	<p>Une urne contient un très grand nombre de boules blanches et de boules noires dont on ignore la proportion p de boules blanches. On tire avec remise n boules dans le but d'estimer la proportion p de boules blanches. On obtient ainsi une fréquence d'apparition qui va nous permettre d'estimer la proportion p à l'aide d'un intervalle de confiance.</p> <div data-bbox="1002 1630 1182 1832" style="text-align: center;">  </div>

Conditions sur les paramètres : Dans tout le chapitre, sauf mention contraire, la taille de l'échantillon n et la proportion p du caractère étudié dans la population vérifient :

$$n \geq 30, n \times p \geq 5 \text{ et } n \times (1 - p) \geq 5.$$

1) Intervalle de fluctuation asymptotique

Dans ce paragraphe, on suppose que la proportion p du caractère étudié est connue.

Exemple :

On dispose d'une urne contenant un grand nombre de boules blanches et noires. La proportion de boules blanches contenues dans l'urne est $p = 0,3$.

On tire successivement avec remise $n = 50$ boules.

Soit X_{50} la variable aléatoire dénombrant le nombre de boules blanches tirées.

X_{50} suit la loi binomiale $B(50 ; 0,3)$.

En effectuant 50 tirages dans cette urne, on va prouver dans ce chapitre que la fréquence d'apparition d'une boule blanche est comprise dans l'intervalle $[0,173 ; 0,427]$ avec une probabilité de 0,95.

Cet intervalle s'appelle l'intervalle de fluctuation asymptotique au seuil 0,95 (ou 95%).

Définition : X_n est une variable aléatoire qui suit une loi binomiale $B(n ; p)$.

La variable aléatoire $F_n = \frac{X_n}{n}$ s'appelle la **variable aléatoire fréquence de succès** pour un schéma de Bernoulli de paramètres n et p .

Propriété : Soit $\alpha \in]0 ; 1[$ et X_n une variable aléatoire qui suit une loi binomiale $B(n ; p)$.

La probabilité que la fréquence F_n prenne ses valeurs dans l'intervalle

$I_n = \left[p - u_\alpha \frac{\sqrt{p(1-p)}}{\sqrt{n}} ; p + u_\alpha \frac{\sqrt{p(1-p)}}{\sqrt{n}} \right]$ se rapproche de $1 - \alpha$ quand la taille de l'échantillon n devient grande. On note : $\lim_{n \rightarrow +\infty} P(F_n \in I_n) = 1 - \alpha$.

Définition : I_n est appelé **intervalle de fluctuation asymptotique** de la fréquence F_n au seuil $1 - \alpha$.

Démonstration :

X_n suit la loi binomiale $B(n ; p)$ donc la suite de variables aléatoires $Z_n = \frac{X_n - E(X_n)}{\sigma(X_n)}$ suit une loi normale centrée réduite $N(0 ; 1)$ et d'après le théorème de Moivre-Laplace, on a :

$$\lim_{n \rightarrow +\infty} P(a \leq Z_n \leq b) = \int_a^b \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx, \text{ pour tous réels } a \text{ et } b \text{ avec } a < b.$$

Or,

$$Z_n = \frac{X_n - E(X_n)}{\sigma(X_n)} = \frac{X_n - np}{\sqrt{np(1-p)}} = \frac{n \left(\frac{X_n}{n} - p \right)}{n \frac{\sqrt{p(1-p)}}{\sqrt{n}}} = \frac{F_n - p}{\frac{\sqrt{p(1-p)}}{\sqrt{n}}}$$

Donc,

$$\lim_{n \rightarrow +\infty} P \left(p + a \frac{\sqrt{p(1-p)}}{\sqrt{n}} \leq F_n \leq p + b \frac{\sqrt{p(1-p)}}{\sqrt{n}} \right) = \int_a^b \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx$$

Comme, pour tout réel $\alpha \in]0 ; 1[$, il existe un unique réel positif u_α tel que $P(-u_\alpha \leq X \leq u_\alpha) = 1 - \alpha$ où X suit une loi normale centrée réduite $N(0 ; 1)$, on a :

$$\int_{-u_\alpha}^{u_\alpha} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = 1 - \alpha$$

En prenant $a = -u_\alpha$ et $b = u_\alpha$, on a :

$$\lim_{n \rightarrow +\infty} P\left(p - u_\alpha \frac{\sqrt{p(1-p)}}{\sqrt{n}} \leq F_n \leq p + u_\alpha \frac{\sqrt{p(1-p)}}{\sqrt{n}}\right) = 1 - \alpha$$

Remarque :

La probabilité définie dans la propriété se rapproche de $1 - \alpha$ sans être nécessairement égale d'où l'emploi du terme "asymptotique".

Exemple :

 Vidéo https://youtu.be/k_Q2FN07jQ0

Démontrons le résultat donné dans l'exemple en début de paragraphe :

$$I_{50} = \left[0,3 - 1,96 \times \frac{\sqrt{0,3 \times 0,7}}{\sqrt{50}} ; 0,3 + 1,96 \times \frac{\sqrt{0,3 \times 0,7}}{\sqrt{50}}\right]$$

car $u_{0,05} = 1,96$.

Soit $I_{50} = [0,173 ; 0,427]$

Pour 500 tirages, on obtient :

$$I_{500} = \left[0,3 - 1,96 \times \frac{\sqrt{0,3 \times 0,7}}{\sqrt{500}} ; 0,3 + 1,96 \times \frac{\sqrt{0,3 \times 0,7}}{\sqrt{500}}\right] = [0,26 ; 0,34]$$

On constate que l'intervalle, pour un même seuil, se resserre fortement lorsqu'on augmente le nombre de tirages.

Définition : On appelle intervalle de fluctuation au seuil 0,95 de la variable aléatoire fréquence l'intervalle :

$$\left[p - 1,96 \times \frac{\sqrt{p(1-p)}}{\sqrt{n}} ; p + 1,96 \times \frac{\sqrt{p(1-p)}}{\sqrt{n}}\right]$$

2) Prise de décision

Dans ce paragraphe, la proportion du caractère étudié n'est pas connue mais est supposée être égale à p .

La prise de décision consiste à valider ou invalider l'hypothèse faite sur la proportion p .

Propriété (Règle de décision) : Soit f la fréquence du caractère étudié d'un échantillon de taille n .

Soit l'hypothèse : "La proportion de ce caractère dans la population est p ."

Soit I l'intervalle de fluctuation asymptotique au seuil 0,95.

- Si $f \in I$, alors on accepte l'hypothèse faite sur la proportion p .

- Si $f \notin I$, alors on rejette l'hypothèse faite sur la proportion p .

Remarque :

On peut interpréter cette propriété par le fait que la probabilité qu'on rejette à tort l'hypothèse sur p sachant qu'elle est vraie est approximativement égale à 5%.

Méthode : Prendre une décision à l'aide d'un intervalle de fluctuation

 Vidéo <https://youtu.be/QZ0YFthGI0Y>

Un fabricant d'alarme commande auprès de son fournisseur deux types de composants électroniques : RS017 et P412. Il demande 900 composants de chaque sorte.

Au moment de la livraison, le service de contrôle retire 50 composants et constate que 19 sont des modèles RS017.

Peut-on affirmer que la commande est respectée par le fournisseur ?

- Le fabricant a commandé autant de composants de chaque sorte. On peut donc supposer que la proportion de composants RS017 est égale à 0,5.

La taille de l'échantillon est $n = 50$.

La fréquence observée est donc $f = \frac{19}{50} = 0,38$.

- Vérifions si les paramètres n et p répondent aux conditions imposées :

$n = 50 \geq 30$, $n \times p = 50 \times 0,5 = 25 \geq 5$ et $n \times (1 - p) = 50 \times 0,5 = 25 \geq 5$

- L'intervalle de fluctuation asymptotique au seuil 0,95 est :

$$\left[0,5 - 1,96 \times \frac{\sqrt{0,5 \times 0,5}}{\sqrt{50}} ; 0,5 + 1,96 \times \frac{\sqrt{0,5 \times 0,5}}{\sqrt{50}} \right] = [0,361 ; 0,639]$$

La fréquence observée $f = 0,38$ appartient à l'intervalle de fluctuation asymptotique au seuil 0,95, d'après la règle de décision, l'hypothèse faite est acceptable.

3) Estimation

Dans ce paragraphe, on suppose que la proportion p du caractère étudié est inconnue.

C'est le problème inverse de celui de l'échantillonnage. À partir de la fréquence observée sur un échantillon, on va estimer la proportion p d'un caractère dans la population tout entière.

Propriété : X_n est une variable aléatoire qui suit une loi binomiale $B(n ; p)$.

$F_n = \frac{X_n}{n}$ est la fréquence associée à X_n .

Pour n suffisamment grand, p appartient à l'intervalle $J_n = \left[F_n - \frac{1}{\sqrt{n}} ; F_n + \frac{1}{\sqrt{n}} \right]$, avec une probabilité supérieure ou égale à 0,95.

- Admis -

Définition : Soit f une fréquence observée du caractère étudié sur un échantillon de taille n .

On appelle **intervalle de confiance** de la proportion p au niveau de confiance 0,95, l'intervalle $\left[f - \frac{1}{\sqrt{n}} ; f + \frac{1}{\sqrt{n}} \right]$.

Remarques :

- Un niveau de confiance 0,95 signifie que dans 95 cas sur 100, on affirme à juste titre que p appartient à l'intervalle de confiance.
- Il n'est pas vrai d'affirmer que p est égal au centre de l'intervalle de confiance. Il n'est pas possible d'évaluer la position de p dans l'intervalle de confiance.
- p étant inconnu, il n'est pas possible de vérifier si les conditions énoncées sur n et p en introduction de chapitre sont vérifiées.

Cependant, il faudra les vérifier sur la fréquence observée f :

$$n \geq 30, n \times f \geq 5 \text{ et } n \times (1 - f) \geq 5$$

Exemple :

On dispose d'une urne contenant un grand nombre de boules blanches et noires. La proportion de boules blanches contenues dans l'urne n'est pas connue.

On réalise un tirage de 100 boules et on obtient 54 boules blanches.

La fréquence observée est donc $f = 0,54$.

L'intervalle de confiance de la proportion de boule blanche dans l'urne au niveau de confiance 95% est $\left[0,54 - \frac{1}{\sqrt{100}} ; 0,54 + \frac{1}{\sqrt{100}} \right] = [0,44 ; 0,64]$.

Méthode : Estimer une proportion inconnue par un intervalle de confiance

 **Vidéo** <https://youtu.be/cU5cJICVAM8>

Un institut de sondage interroge 1052 personnes entre les deux tours de l'élection présidentielle sur leur intention de vote.

614 déclarent avoir l'intention de voter pour Martine Phinon.

En supposant que les votes seront conformes aux intentions, la candidate a-t-elle raison de croire qu'elle sera élue ?

- La proportion p des électeurs de Martine Phinon est inconnue.

La taille de l'échantillon est $n = 1052$.

La fréquence observée est $f = \frac{614}{1052} \approx 0,5887$.

- Vérifions si les paramètres n et f répondent aux conditions imposées :

$$n = 1052 \geq 30, n \times f = 1052 \times 0,5887 \approx 614 \geq 5$$

$$\text{et } n \times (1 - f) = 1052 \times (1 - 0,5887) \approx 438 \geq 5$$

- L'intervalle de confiance de la proportion p au niveau de confiance 0,95 est :

$$\left[0,5887 - \frac{1}{\sqrt{1052}} ; 0,5887 + \frac{1}{\sqrt{1052}} \right] \approx [0,553 ; 0,615].$$

Pour être élue, la proportion p doit être strictement supérieure à 0,5. Selon ce sondage, il est envisageable que Martine Phinon soit élue.

Intervalle de FLUCTUATION V.S. Intervalle de CONFIANCE :

 **Vidéo** <https://youtu.be/97vzxWsyie8>

Méthode : Déterminer une taille d'échantillon suffisante pour obtenir une estimation d'une proportion

 Vidéo <https://youtu.be/ogmMVpkBVgs>

Un constructeur automobile fait appel à un institut de sondage afin de mesurer le degré de satisfaction du service après-vente.

L'institut souhaite estimer la proportion de clients satisfaits au niveau de confiance 0,95 avec une amplitude d'au plus 5 centièmes.

Combien de personnes au minimum faut-il interroger ?

On appelle p la proportion de clients satisfaits. Cette proportion est inconnue.

Une estimation de cette proportion peut être obtenue à l'aide de l'intervalle de confiance au niveau de confiance 0,95 : $\left[f - \frac{1}{\sqrt{n}} ; f + \frac{1}{\sqrt{n}} \right]$, où f est la fréquence observée.

Cet intervalle a pour longueur $\frac{2}{\sqrt{n}}$.

Donc $\frac{2}{\sqrt{n}} \leq 0,05$ soit $n \geq \frac{4}{0,05^2} = 1600$.

L'institut de sondage devra donc interroger au moins 1600 personnes.



Hors du cadre de la classe, aucune reproduction, même partielle, autres que celles prévues à l'article L 122-5 du code de la propriété intellectuelle, ne peut être faite de ce site sans l'autorisation expresse de l'auteur.

www.maths-et-tiques.fr/index.php/mentions-legales