

TP Informatique n° 14/15

Estimation ponctuelle et intervalle de confiance

I. Aspects théoriques

I.1. Contexte de l'estimation

Dans la suite, on considère l'exemple suivant : une urne contient des boules vertes et des boules rouges dont on ignore le nombre et la proportion. On effectue un tirage dans cette urne et on note X la v.a.r. égale à 1 si la boule obtenue est verte et à 0 sinon. Le but est d'estimer la proportion p de boules vertes.

Cet exemple permet de comprendre le contexte de l'estimation.

- On considère un phénomène aléatoire et une v.a.r. X qui lui est liée.
- La loi de X est globalement connue (X suit une loi de Bernoulli dans l'exemple précédent) mais pas entièrement spécifiée (le paramètre p est à déterminer).
- Le problème de l'estimation consiste alors à estimer la valeur du paramètre θ (non spécifié) de cette loi (p dans notre exemple).

I.2. Estimateur

Afin d'estimer le paramètre θ , on observe n fois le phénomène (on tire n fois dans l'urne).

- On obtient alors un n -uplet de données appelé **échantillon observé**.
- Cette observation est une **réalisation** de n v.a.r. X_1, \dots, X_n mutuellement indépendantes et de même loi que X et définies sur le même espace probabilisable (X_i est la v.a.r. égale à 1 si le $i^{\text{ème}}$ tirage donne une boule verte et à 0 sinon).

Un estimateur T_n du paramètre θ est une v.a.r. $\varphi(X_1, \dots, X_n)$ fonction de l'échantillon (X_1, \dots, X_n) .

Toute réalisation $\varphi(x_1, \dots, x_n)$ est appelé estimation de θ .

I.3. Critères de qualité de l'estimateur

Pour tenter d'évaluer la qualité d'un estimateur on dispose des critères suivants.

- **Estimateur sans biais** : $\mathbb{E}_\theta(T_n) = \theta$ (défini si T_n admet une espérance)

On note $b_\theta(T_n) = \mathbb{E}_\theta(T_n) - \theta$ le biais de l'estimateur T_n .

- **Estimateur asymptotiquement sans biais** : $\lim_{n \rightarrow +\infty} \mathbb{E}_\theta(T_n) = \theta$

- **Estimateur convergent** : $\forall \varepsilon > 0, \lim_{n \rightarrow +\infty} \mathbb{P}_\theta(|T_n - \theta| \geq \varepsilon) = 0$

La notion d'estimateur convergent peut être caractérisée à l'aide du **risque quadratique** :

$$\begin{aligned} r_\theta(T_n) &= \mathbb{E}_\theta((T_n - \theta)^2) \\ &= \mathbb{V}_\theta(T_n) + (b_\theta(T_n))^2 \end{aligned} \quad (\text{défini si, pour tout } \theta, T_n^2 \text{ admet une espérance})$$

On a alors : $\lim_{n \rightarrow +\infty} r_\theta(T_n) = 0 \Rightarrow T_n$ est un estimateur convergent

II. Estimation du paramètre p d'une loi de Bernoulli

On prend (une nouvelle fois) l'exemple du cours : une urne contient des boules vertes et des boules rouges dont on ignore le nombre et la proportion. On effectue un tirage dans cette urne et on note X la v.a.r. égale à 1 si la boule obtenue est verte et à 0 sinon.

Autrement dit, $X \hookrightarrow \mathcal{B}(p)$ où p est un paramètre à déterminer.

II.1. Aspect théorique

- Quel estimateur T_n peut-on proposer ?

Il s'agit ici d'estimer une espérance. On choisit (X_1, \dots, X_n) mutuellement indépendantes et de même loi que X . Autrement dit, $X_i \hookrightarrow \mathcal{B}(p)$.

On choisit comme estimateur la **moyenne empirique** : $\overline{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$.

- Calculer le biais de cet estimateur.

Par linéarité de l'espérance, on a :

$$\mathbb{E}(T_n) = \mathbb{E}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(X_i) = \frac{1}{n} \sum_{i=1}^n p = \frac{1}{n} n \times p = p$$

Ainsi, $b(T_n) = \mathbb{E}(T_n) - p = 0$ et T_n est sans biais.

- Calculer la variance de cet estimateur.

Les v.a.r. (X_1, \dots, X_n) étant mutuellement indépendantes, on a :

$$\begin{aligned} \mathbb{V}(T_n) &= \mathbb{V}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \mathbb{V}\left(\sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n \mathbb{V}(X_i) = \frac{1}{n^2} \sum_{i=1}^n p(1-p) \\ &= \frac{1}{n^2} n \times p(1-p) = \frac{p(1-p)}{n} \end{aligned}$$

- En déduire que cet estimateur est convergent.

$$r(T_n) = \mathbb{V}(T_n) + (b(T_n))^2 = \frac{p(1-p)}{n} \xrightarrow{n \rightarrow +\infty} 0.$$

L'estimateur T_n est donc convergent.

- Quel résultat permet, sans calcul de $\mathbb{V}(T_n)$, de conclure que $T_n = \overline{X}_n$ est convergent ?

- Cet estimateur d'espérance \overline{X}_n joue un rôle particulier.

En effet, énoncer sa convergence revient à citer la loi faible des grands nombres.

Soit $(X_n)_{n \in \mathbb{N}^*}$ une suite de variables aléatoires indépendantes de même espérance m et de même variance.

Alors :
$$\forall \varepsilon > 0, \lim_{n \rightarrow +\infty} \mathbb{P}(|\overline{X}_n - m| \geq \varepsilon) = 0$$

- Ce résultat signifie la convergence en probabilité de la v.a.r. \overline{X}_n vers la v.a.r. constante m . Cependant, ce résultat ne dit rien sur la vitesse à laquelle a lieu cette convergence. On y reviendra plus loin.

II.2. Modélisation en Python

- Dans la suite, on considère que l'urne contient 300 boules dont 192 vertes.
- Pour estimer la proportion p de boules vertes dans cette urne :
 - × on effectue n tirages successifs et avec remise,
 - × on calcule la moyenne de ces observations.
(c'est une réalisation de la moyenne empirique)
- Pour les simulations, on pourra utiliser la fonction `binomial`, accessible depuis l'espace de nommage `numpy.random` et qui permet de simuler une v.a.r. suivant une loi binomiale. Plus précisément, `numpy.random(n, p, (nbL, nbC))` renvoie un tableau de taille $\text{nbL} \times \text{nbC}$ dont chaque case est la simulation d'une v.a.r. X telle que $X \hookrightarrow \mathcal{B}(n, p)$.
- On commencera par l'importation habituelle.

```

1 import numpy as np
2 import matplotlib.pyplot as plt

```

- Quelle est la valeur exacte de p ?

Il y a 192 boules vertes et 300 boules au total donc $p = \frac{192}{300} = \frac{64}{100} = 0,64$.

- Quel appel permet de simuler une v.a.r. qui suit une loi de Bernoulli de paramètre p ?
Comment simule-t-on n v.a.r. indépendantes suivant une loi de Bernoulli de paramètre p ?

Pour 1 v.a.r. : `np.random.binomial(1,p,(1,1))`, et pour n v.a.r. :
`np.random.binomial(1,p,(1,n))`.

- Comment peut-on calculer la moyenne d'un tableau `Obs` contenant n éléments.

Au choix : `np.mean(Obs)` ou `np.sum(Obs)/n` ou `np.sum(Obs)/np.size(Obs)`.

- Écrire une fonction `estimME` qui prend en paramètre un réel p et un entier n et qui renvoie une simulation de l'estimateur T_n .

```

1 def estimME(p, n) :
2     Obs = np.random.binomial(1,p,(1,n))
3     return np.mean(Obs)

```

- En théorie, avec cette méthode, quelle est la proportion minimale que l'on peut observer?
Et la maximale?

- On peut ne tirer que des boules rouges (n'obtenir que des échecs).
D'où une proportion minimale de 0.
- On peut ne tirer que des boules vertes (n'obtenir que des succès).
D'où une proportion maximale de 1.

- Exécuter 8 fois de suite la fonction `estimME` avec les paramètres $p=0.64$ et $n=15$.
Noter ci-dessous les résultats.

0.67 0.8 0.47 0.6 0.8 0.8 0.67 0.67

- ▶ Exécuter 8 fois de suite la fonction `estimME` avec les paramètres $p=0.64$ et $n=50$.
Noter ci-dessous les résultats.

0.54	0.68	0.72	0.54	0.76	0.66	0.6	0.64
------	------	------	------	------	------	-----	------

- ▶ Commenter les résultats obtenus.

On observe une certaine variabilité des résultats qui semble plus grande dans le premier cas que dans le deuxième.
--

II.3. Garantie offerte par cette méthode

Les résultats précédents sont difficilement exploitables tel quels. En effet, ils dépendent fortement de la manière dont les tirages ont été effectués. A-t-on eu plutôt la main chaude (tendance à obtenir plus de succès que prévu) ou plutôt la main froide (tendance à obtenir plus d'échecs)? Il est compliqué de s'entraîner à avoir la main tiède. Il nous faut donc un procédé mathématique permettant d'estimer la qualité des résultats obtenus.

II.3.a) Loi approchée de T_n par simulation

Dans le chapitre précédent, on a simulé T_n . Essayons de déterminer la loi de T_n par simulation. Pour ce faire :

- × on se fixe un n ($n1 = 15$ puis $n2 = 50$ et enfin $n3 = 150$),
- × on simule N réalisations de la v.a.r. T_n (prendre N grand, $N = 10000$ par exemple),
- × on trace le diagramme en bâtons de ces N valeurs.
- ▶ À n fixé (prenons $n = 15$ par exemple), quel est le support de T_n ?

T_{15} compte la moyenne du nombre de succès obtenu lors d'une répétition de 15 épreuves de Bernoulli indépendantes (et de même paramètre p). Ainsi :
--

$$T_{15}(\Omega) = \left\{ 0, \frac{1}{15}, \frac{2}{15}, \dots, \frac{14}{15}, 1 \right\}. \text{ De manière générale, on a :}$$

$$T_n(\Omega) = \left\{ 0, \frac{1}{n}, \frac{2}{n}, \dots, \frac{n-1}{n}, 1 \right\} \text{ et } T_n \text{ est une v.a.r. discrète.}$$

- ▶ Si n est fourni, comment obtient-on le support de T_n en **Python**?

Il suffit d'écrire <code>np.linspace(0,1,n+1)</code> .
--

- ▶ Écrire un programme **Python** qui :

- × réalise les affectations initiales $p = 0.64$, $n1 = 15$ et $N = 10000$,
- × crée une liste `ObsT1` de N simulations de T_n ,
- × affiche le diagramme en bâtons associé aux valeurs de ce vecteur.
(on utilisera la fonction `hist` et on prendra soin de prendre en paramètre un nombre de classes cohérent)

<pre> 1 p = 0.64 2 n1 = 15 3 N = 10000 4 5 ObsT1 = [estimME(p,n1) for k in range(N)] 6 plt.hist(ObsT1, n1, normed = True, color = 'r') 7 plt.show() </pre>
--

- Ajouter au programme précédent les variables $n2 = 50$ et $n3 = 150$ et tracer les diagrammes des vecteurs `ObsT2` et `ObsT3` correspondants.
Les 3 diagrammes doivent être affichés côte à côte. (*on utilisera la fonction `subplot`*)

On ajoute les lignes suivantes.

```

1  n2 = 50
2  n3 = 150
3
4  ObsT2 = [estimME(p,n2) for k in range(N)]
5  plt.hist(ObsT2, n2, normed = True, color = 'g')
6
7  ObsT3 = [estimME(p,n3) for k in range(N)]
8  plt.hist(ObsT3, n3, normed = True, color = 'b')
```

- Commenter les graphiques obtenus :

- 1) quel type de courbe obtient-on ?
- 2) quel théorème permettait de prévoir ce résultat à l'avance ?

1) On obtient des courbes en cloche (des gaussiennes) centrées en 0,64.

2) Ce résultat illustre le TCL.

- Soit $(X_n)_{n \in \mathbb{N}^*}$ suite de v.a.r. indépendantes, de même loi, d'espérance m et d'écart-type σ .

$$\text{Notons } \overline{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \text{ et } \overline{X}_n^* = \frac{\overline{X}_n - m}{\frac{\sigma}{\sqrt{n}}}.$$

(\overline{X}_n^* est la v.a.r. centrée réduite associée à \overline{X}_n)

Alors la suite (\overline{X}_n^*) converge en loi vers une v.a.r. N telle que $N \leftrightarrow \mathcal{N}(0, 1)$.

- L'idée est donc que \overline{X}_n^* va se comporter comme N pour n grand (très rapidement en pratique). Mieux : \overline{X}_n va se comporter comme M_n avec $M_n \leftrightarrow \mathcal{N}\left(m, \frac{\sigma}{\sqrt{n}}\right)$ pour n grand.

- Dans notre exemple, T_n a donc une loi proche de $\mathcal{N}\left(p, \frac{\sqrt{p(1-p)}}{\sqrt{n}}\right)$.

- Où lit-on les éléments caractéristiques (espérance, écart-type) d'une loi normale ? Déterminer les éléments caractéristiques pour les 3 courbes précédentes (lecture graphique et calcul de la valeur attendue).

• La cloche est centrée en l'espérance m ; le point au sommet de cloche a pour ordonnée $\frac{1}{\sigma \sqrt{2\pi}}$ où σ est l'écart-type de la loi normale considérée.

• Dans les 3 exemples, on a toujours $m = p = 0,64$ et :

$$\times \sigma_1 \simeq \frac{\sqrt{p(1-p)}}{\sqrt{15}} \simeq 0,124 \quad \text{d'où} \quad \frac{1}{\sigma_1 \sqrt{2\pi}} \simeq 3,22$$

$$\times \sigma_2 \simeq \frac{\sqrt{p(1-p)}}{\sqrt{50}} \simeq 0,068 \quad \text{d'où} \quad \frac{1}{\sigma_2 \sqrt{2\pi}} \simeq 5,88$$

$$\times \sigma_3 \simeq \frac{\sqrt{p(1-p)}}{\sqrt{150}} \simeq 0,039 \quad \text{d'où} \quad \frac{1}{\sigma_3 \sqrt{2\pi}} \simeq 10,18$$

- En quoi la 3^{ème} expérience fournit-elle des résultats meilleurs que les 2 premières ?

Plus n est grand plus la largeur de la cloche est faible. La cloche est centrée en p , valeur que l'on souhaite obtenir. Une cloche de largeur faible assure que la méthode va fournir, avec une probabilité forte, une valeur proche de la valeur p réelle.

II.3.b) Intervalle de confiance asymptotique : aspect théorique

On cherche maintenant à estimer quelle garantie on peut accorder au résultat précédent. A priori, on peut obtenir toutes les valeurs de $[0, 1]$ par nos tirages. Cependant, certaines valeurs (celles très éloignées de p) n'ont qu'une probabilité très faible d'être obtenues alors que d'autres (qui sont proches de p) ont une probabilité forte d'être obtenues.

Ceci invite à poser la notion de garantie sous la forme :

« le résultat obtenu est ε -proche du résultat réel avec une probabilité d'au moins $1 - \alpha$ »

(si on prend $\alpha = 0,05$, on obtient un niveau de confiance de 95%)

- Comment exprime-t-on, à l'aide de l'opérateur $|\cdot|$, le fait que deux réels t et p soient écartés au plus de ε ? À quel encadrement de p cela correspond-il ?

$$|t - p| \leq \varepsilon \Leftrightarrow -\varepsilon \leq t - p \leq \varepsilon \Leftrightarrow -t - \varepsilon \leq -p \leq -t + \varepsilon \Leftrightarrow t - \varepsilon \leq p \leq t + \varepsilon$$

- En déduire une expression de la garantie qu'il semble raisonnable d'exiger dans notre exemple.

• Il semble raisonnable de demander : $\mathbb{P}(T_n - \varepsilon \leq p \leq T_n + \varepsilon) \geq 1 - \alpha$.
 • En fait, le résultat obtenu par le TCL est un peu plus faible et s'exprime sous la forme : $\lim_{n \rightarrow +\infty} \mathbb{P}(T_n - \varepsilon \leq p \leq T_n + \varepsilon) \geq 1 - \alpha$.

Cette garantie n'est rien d'autre que la notion d'intervalle de confiance.

Intervalles de confiance et de confiance asymptotique

Soit $(U_n)_{n \in \mathbb{N}^*}$ et $(V_n)_{n \in \mathbb{N}^*}$ deux suites d'estimateurs de θ .

- 1) Pour $n \in \mathbb{N}^*$, on dit que $[U_n, V_n]$ est un intervalle de confiance de θ au niveau de confiance $1 - \alpha$ (où le risque α est un élément de $]0, 1[$) si :

$$\forall \theta \in \Theta, \quad \mathbb{P}_\theta(U_n \leq \theta \leq V_n) \geq 1 - \alpha$$

(classiquement obtenu par l'inégalité de Bienaymé-Chebychev)

- 2) On dit que $([U_n, V_n])_{n \in \mathbb{N}^*}$ est un intervalle de confiance asymptotique de θ au niveau de confiance $1 - \alpha$ (où le risque α est un élément de $]0, 1[$) si :

$$\forall \theta \in \Theta, \quad \lim_{n \rightarrow +\infty} \mathbb{P}_\theta(U_n \leq \theta \leq V_n) \geq 1 - \alpha$$

(classiquement obtenu par le TCL)

Deux notions entrent ici en compétition :

- × la précision de l'intervalle (*i.e.* l'écart $V_n - U_n$),
- × et le niveau de confiance $(1 - \alpha)$ que l'on peut accorder à cet intervalle.



- Plus la précision exigée augmente, plus le niveau de confiance diminue.
(une précision importante signifie que l'écart $V_n - U_n$ est faible)
- Plus la précision exigée diminue, plus le niveau de confiance augmente.
(une précision faible signifie que l'écart $V_n - U_n$ est grand)

II.3.c) Intervalle de confiance : retour à notre exemple

- Le TCL permet d'affirmer que la suite (\overline{X}_n^*) converge en loi vers une v.a.r. N telle que $N \hookrightarrow \mathcal{N}(0, 1)$. Que cela signifie-t-il sur la suite $(\mathbb{P}(-x_0 \leq \overline{X}_n^* \leq x_0))_n$? (avec $x_0 \in \mathbb{R}$)
(on exprimera le résultat à l'aide de Φ)

Cela signifie que cette suite est convergente et que :

$$\begin{aligned} \lim_{n \rightarrow +\infty} \mathbb{P}(-x_0 \leq \overline{X}_n^* \leq x_0) &= \mathbb{P}(-x_0 \leq N \leq x_0) \\ &= \Phi(x_0) - \Phi(-x_0) = \Phi(x_0) - (1 - \Phi(x_0)) \\ &= 2\Phi(x_0) - 1 \end{aligned}$$

- À quel encadrement de p l'inégalité $-x_0 \leq \overline{X}_n^* \leq x_0$ est-elle équivalente ?

$$\begin{aligned}
 -x_0 \leq \overline{X}_n^* \leq x_0 &\Leftrightarrow -x_0 \leq \frac{\overline{X}_n - p}{\frac{\sqrt{p(1-p)}}{\sqrt{n}}} \leq x_0 \\
 &\Leftrightarrow -x_0 \frac{\sqrt{p(1-p)}}{\sqrt{n}} \leq \overline{X}_n - p \leq x_0 \frac{\sqrt{p(1-p)}}{\sqrt{n}} \\
 &\Leftrightarrow -\overline{X}_n - x_0 \frac{\sqrt{p(1-p)}}{\sqrt{n}} \leq -p \leq -\overline{X}_n + x_0 \frac{\sqrt{p(1-p)}}{\sqrt{n}} \\
 &\Leftrightarrow \overline{X}_n - x_0 \frac{\sqrt{p(1-p)}}{\sqrt{n}} \leq p \leq \overline{X}_n + x_0 \frac{\sqrt{p(1-p)}}{\sqrt{n}} \\
 &\Leftrightarrow p \in \left[\overline{X}_n - x_0 \frac{\sqrt{p(1-p)}}{\sqrt{n}}, \overline{X}_n + x_0 \frac{\sqrt{p(1-p)}}{\sqrt{n}} \right]
 \end{aligned}$$

- En déduire l'intervalle de confiance associé à notre problème.

- On ne peut prendre $U_n = \overline{X}_n - x_0 \frac{\sqrt{p(1-p)}}{\sqrt{n}}$ et $V_n = \overline{X}_n + x_0 \frac{\sqrt{p(1-p)}}{\sqrt{n}}$ puisque ces estimateurs dépendent de p , quantité que l'on cherche à estimer.
- La méthode consiste à majorer $p(1-p)$ par $\frac{1}{4}$.
(l'intervalle de confiance obtenu est alors plus grand)
- On choisit donc $U_n = \overline{X}_n - x_0 \frac{1}{2\sqrt{n}}$ et $V_n = \overline{X}_n + x_0 \frac{1}{2\sqrt{n}}$.

- Quelle valeur de x_0 (expression en fonction de Φ^{-1}) permet d'assurer un niveau de confiance de $1 - \alpha$? Quelle est la largeur de l'intervalle (précision) associée ?

- On doit choisir x_0 tel que $2 \Phi(x_0) - 1 \geq 1 - \alpha$. Or :

$$2 \Phi(x_0) - 1 \geq 1 - \alpha \Leftrightarrow \Phi(x_0) \geq 1 - \frac{\alpha}{2} \Leftrightarrow x_0 \geq \Phi^{-1} \left(1 - \frac{\alpha}{2} \right)$$
- La largeur de l'intervalle est alors :

$$V_n - U_n = \left(\overline{X}_n + x_0 \frac{1}{2\sqrt{n}} \right) - \left(\overline{X}_n - x_0 \frac{1}{2\sqrt{n}} \right) = \frac{x_0}{\sqrt{n}}$$

- Quel résultat retrouve-t-on alors ?

- Φ^{-1} est strictement croissante (car de même monotonie que Φ).
- On a ainsi que x_0 augmente lorsque $1 - \frac{\alpha}{2}$ augmente.
- Ce qui permet de conclure que lorsque le niveau de confiance ($1 - \alpha$) augmente alors x_0 augmente aussi, ce qui signifie que la largeur de l'intervalle augmente et donc que la précision diminue.

La fonction Φ^{-1} , inverse de la fonction de répartition Φ , est appelée fonction **quantile** (parfois appelée *percent point function* en anglais). En **Python**, cette fonction est accessible dans l'espace de nommage **norm** qu'il faut importer depuis la bibliothèque **scipy.stats**.

```
1 from scipy.stats import norm
```

L'appel $\Phi^{-1}(v)$ est réalisé en **Python** par la commande `norm.ppf(v)`

- ▶ Quelle valeur de x_0 permet d'assurer un niveau de confiance à 85% ?
- Quelle valeur de x_0 permet d'assurer un niveau de confiance à 90% ?
- Quelle valeur de x_0 permet d'assurer un niveau de confiance à 95% ?
- On déterminera les précisions (largeur de l'intervalle) associées dans le cas où $n = 150$.

On choisit $n = 150$ et on réalise l'appel `x0 = norm.ppf(1-alpha/2)`.

- Si `alpha = 0.15`, on obtient `x0 ≈ 1.44` et `x0 / sqrt(n) ≈ 0.12`.
- Si `alpha = 0.1`, on obtient `x0 ≈ 1.64` et `x0 / sqrt(n) ≈ 0.13`.
- Si `alpha = 0.05`, on obtient `x0 ≈ 1.96` et `x0 / sqrt(n) ≈ 0.16`.

Encore une fois, lorsque le niveau de confiance augmente, la précision diminue.

II.3.d) Intervalle de confiance asymptotique : visualisation

On fixe maintenant $n = 150$, $\alpha = 0.05$ (niveau de confiance de 95%).

Afin de visualiser la notion d'intervalle de confiance, on agit comme suit :

- × on simule $N = 100$ réalisations de T_{150} ,
- × pour chaque réalisation de T_{150} , on obtient une réalisation $[u, v]$ de $[U_{150}, V_{150}]$ associée,
 - (i) si p est dans l'intervalle de confiance $[u, v]$, on trace un trait vert vertical d'ordonnée minimale u et maximale v .
 - (ii) si p n'est pas dans l'intervalle de confiance $[u, v]$, on trace un trait rouge vertical d'ordonnée minimale u et maximale v .
- × on trace enfin un trait horizontal pour visualiser la valeur de p .
- ▶ Compléter le programme suivant afin qu'il réalise la stratégie précédente.

```

1  p = 0.64
2  n = 150
3  N = 100
4  alpha = 0.05
5
6  x0 = norm.ppf(1-alpha/2)
7
8  for i in range(N) :
9      T = estimME(p, n)
10     eps = x0 / (2*np.sqrt(n))
11     if abs(T-p) <= eps :
12         plt.plot([i+1, i+1], [T-eps, T+eps], 'g')
13     else :
14         plt.plot([i+1, i+1], [T-eps, T+eps], 'r')
15 plt.plot([1, N], [p, p])
16 # "Astuce" pour que les ordonnées soient entre 0 et 1
17 plt.plot([0, 0], [0, 1], "b")

```

► Commenter le résultat obtenu.

- On obtient 3 intervalles rouges et 97 intervalles verts.
- Vu que l'on calcule des IC à niveau de confiance 95%, sur 100 intervalles calculés on doit en obtenir au maximum 5 de faux *i.e.* 5 qui ne contiennent pas la valeur estimée p .
- La simulation obtenue est donc en parfait accord avec la théorie mathématique.

III. Test de conformité à la moyenne

III.1. Introduction

- On considère n réalisations (x_1, \dots, x_n) censées provenir d'un échantillon (X_1, \dots, X_n) de v.a.r. indépendantes, identiquement distribuées et d'espérance m . Le test de conformité à la moyenne consiste à confronter la moyenne de ces réalisations à la moyenne théorique m . L'idée est de pouvoir répondre à la question : les réalisations obtenues sont-elles conformes à nos hypothèses sur les v.a.r. X_i et leur moyenne.
- Dans la pratique, il s'agit souvent de tester si un échantillon est représentatif d'une population.
- Cela peut servir pour des tests de qualité. Considérons par exemple le test de qualité d'une chaîne de production de médicaments. Chaque médicament est censé contenir 12 mg de principe actif. On teste alors un échantillon pour déterminer s'il appartient à la famille des médicaments conformes de la chaîne.

Plus précisément, on calcule la moyenne de cet échantillon prélevé.

- 1) soit le poids moyen du produit actif ne s'écarte pas significativement des 12 mg. Dans ce cas, on rejette l'hypothèse de conformité.
- 2) dans le cas contraire, on ne peut rejeter l'hypothèse.

III.2. Rappels théoriques

III.2.a) Si l'on connaît la variance théorique

Dans la section précédente, nous avons montré comment obtenir des intervalles de confiance asymptotique. On a en fait détaillé le théorème suivant (c'est un corollaire immédiat du TCL).

Intervalles de confiance asymptotique via le TCL

- Soit $(X_n)_{n \in \mathbb{N}^*}$ une suite de v.a.r. **indépendantes**, identiquement distribuées, d'espérance m et de variance σ .
- Soit $\alpha \in]0, 1[$. On note $u_{1-\frac{\alpha}{2}}$ le quantile d'ordre $1 - \frac{\alpha}{2}$ de la loi $\mathcal{N}(0, 1)$. Par définition, $u_{1-\frac{\alpha}{2}}$ est l'unique réel tel que : $\Phi(u_{1-\frac{\alpha}{2}}) = 1 - \frac{\alpha}{2}$ *i.e.* $u_{1-\frac{\alpha}{2}} = \Phi^{-1}(1 - \frac{\alpha}{2})$.
- On note $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ (*moyenne empirique*).

$$\text{On a alors : } \lim_{n \rightarrow +\infty} \mathbb{P} \left(\bar{X}_n - u_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \leq m \leq \bar{X}_n + u_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \right) = 1 - \alpha$$

$$\text{Ou encore : } \lim_{n \rightarrow +\infty} \mathbb{P} \left(m \notin \left[\bar{X}_n - u_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}, \bar{X}_n + u_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \right] \right) = \alpha$$

$$\text{Ou enfin : } \lim_{n \rightarrow +\infty} \mathbb{P} \left(|\bar{X}_n - m| > u_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \right) = \alpha$$

On considère dans la suite un n -uplet censé provenir d'un échantillon (X_1, \dots, X_n) de v.a.r. telles que $X_i \hookrightarrow \mathcal{N}(\mu, \sigma^2)$ dont on connaît la variance $\sigma^2 = (0.2)^2$.

Le but de l'exercice est de savoir si l'on peut rejeter l'hypothèse H_0 selon laquelle $\mu = 0.64$.

En **Python**, la génération de n réalisations d'un échantillon (X_1, \dots, X_n) de v.a.r. suivant une même loi normale peut se faire à l'aide de la fonction `rvs`, accessible dans l'espace de nommage `norm`. Plus précisément, on doit réaliser l'appel :

```
1 norm.rvs(loc, scale, size)
```

où `loc` est la moyenne de la loi normale considérée, `scale` est son écart-type et `size` est la taille de l'échantillon considéré.

- Compléter le programme suivant. Il doit permettre de répondre à la question : « peut-on rejeter l'hypothèse $H_0 : \mu = 0.64$? ».

```
1 # Valeur des paramètres
2 p = 0.64
3 sig = 0.2
4 n = 50
5 alpha = 0.05
6
7 # Quantile d'ordre 1-alpha/2
8 x0 = norm.ppf(1-alpha/2)
9
10 Obs = norm.rvs(p, sig, n)
11 Xb = np.mean(Obs)
12 # Bornes de l'intervalle de confiance
13 U = Xb - x0 * sig / np.sqrt(n)
14 V = Xb + x0 * sig / np.sqrt(n)
15
16 # Test de conformité à la moyenne
17 print((p < U) | (p > V))
```

III.2.b) Si l'on ne connaît pas la variance théorique

Dans le cas où la variance théorique σ n'est pas connue, on utilise le théorème suivant (deuxième version du TCL).

Intervalles de confiance asymptotique via le TCL

- Soit $(X_n)_{n \in \mathbb{N}^*}$ une suite de v.a.r. **indépendantes**, identiquement distribuées, d'espérance m et de variance σ .
- Soit $\alpha \in]0, 1[$. On note $u_{1-\frac{\alpha}{2}}$ le quantile d'ordre $1 - \frac{\alpha}{2}$ de la loi $\mathcal{N}(0, 1)$. Par définition, $u_{1-\frac{\alpha}{2}}$ est l'unique réel tel que : $\Phi(u_{1-\frac{\alpha}{2}}) = 1 - \frac{\alpha}{2}$ i.e. $u_{1-\frac{\alpha}{2}} = \Phi^{-1}(1 - \frac{\alpha}{2})$.
- On note $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ (moyenne empirique).
- On note $S_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = \frac{1}{n} \left(\sum_{i=1}^n X_i^2 \right) - \bar{X}_n^2$ (variance empirique).

$$\text{On a alors : } \lim_{n \rightarrow +\infty} \mathbb{P} \left(\bar{X}_n - u_{1-\frac{\alpha}{2}} \frac{S_n}{\sqrt{n}} \leq m \leq \bar{X}_n + u_{1-\frac{\alpha}{2}} \frac{S_n}{\sqrt{n}} \right) = 1 - \alpha$$

$$\text{Ou encore : } \lim_{n \rightarrow +\infty} \mathbb{P} \left(m \notin \left[\bar{X}_n - u_{1-\frac{\alpha}{2}} \frac{S_n}{\sqrt{n}}, \bar{X}_n + u_{1-\frac{\alpha}{2}} \frac{S_n}{\sqrt{n}} \right] \right) = \alpha$$

$$\text{Ou enfin : } \lim_{n \rightarrow +\infty} \mathbb{P} \left(|\bar{X}_n - m| > u_{1-\frac{\alpha}{2}} \frac{S_n}{\sqrt{n}} \right) = \alpha$$

On reprend l'exemple précédent du n -uplet censé provenir d'un échantillon (X_1, \dots, X_n) de v.a.r. telles que $X_i \hookrightarrow \mathcal{N}(\mu, \sigma^2)$ dont on ne connaît pas ici la variance σ^2 . Le but de l'exercice est de savoir si l'on peut rejeter l'hypothèse H_0 selon laquelle $\mu = 0.64$.

- Compléter le programme suivant. Il doit permettre de répondre à la question : « peut-on rejeter l'hypothèse $H_0 : \mu = 0.64$? ».

```

1  # Valeur des paramètres
2  p = 0.64
3  sig = np.random.rand()
4  n = 50
5  alpha = 0.05
6
7  # Quantile d'ordre 1-alpha/2
8  x0 = norm.ppf(1-alpha/2)
9
10 Obs = norm.rvs(p, sig, n)
11 Xb = np.mean(Obs)
12 S2 = np.mean([x**2 for x in Obs]) - Xb**2
13 # Test de conformité à la moyenne
14 print(np.abs(Xb - p) > x0 * np.sqrt(S2 / n))

```