

TP7 Python

TESTS STATISTIQUES

1 Introduction : principe d'un test statistique

On considère une population Ω dont les individus possèdent un caractère X (modélisé par une variable aléatoire) dont la loi (ou un certains nombres de paramètres) sont inconnus.

Exemple : on dispose d'une pièce qu'on peut modéliser par une variable aléatoire suivant une loi de Bernoulli de paramètre inconnue p (la probabilité de faire pile par exemple).

On souhaite formuler une hypothèse sur le(s) paramètre(s) inconnu(s) (par exemple, une hypothèse sur sa valeur) et porter un jugement sur cette hypothèse (est-elle raisonnable?), en se basant sur l'observation d'un échantillon prélevé.

Exemple (suite) : on souhaite déterminer si notre pièce est équilibrée ou non c'est-à-dire si $p = \frac{1}{2}$.

Pour cela, on peut par exemple lancer la pièce un grand nombre de fois et observer la proportion de piles et de faces obtenus.

Principe des tests :

- **Un test d'hypothèse** est une démarche qui a pour but de fournir une règle de décision permettant, sur la base de résultats d'échantillon, de faire un choix entre deux hypothèses.
- **Hypothèse nulle, H_0** : c'est l'hypothèse que l'on cherche à tester. Par exemple, on fait l'hypothèse que la probabilité de faire pile vaut $\frac{1}{2}$.
- **Seuil de signification d'un test d'hypothèse** : c'est le risque α consenti à l'avance de rejeter à tort l'hypothèse H_0 , alors qu'elle est vraie.
On utilisera en général $\alpha = 5\%$.
- **Statistique de test ou variable d'échantillonnage** : c'est une variable aléatoire T (en lien avec le problème initial) dont on connaît la loi sous l'hypothèse H_0 .
On sépare alors \mathbb{R} en deux zones : une zone de rejet notée R_{rejet} et une zone de non-rejet notée $R_{\text{non-rejet}}$ telles que :

$$\mathbb{P}(T \in R_{\text{rejet}}) = \alpha \quad ; \quad \mathbb{P}(T \in R_{\text{non-rejet}}) = 1 - \alpha.$$

-Utilisation du test : à partir d'un échantillon, on calcule une valeur observée de T , notée t_{obs} .

- Si $t_{\text{obs}} \notin R_{\text{non-rejet}}$ c'est-à-dire $t_{\text{obs}} \in R_{\text{rejet}}$ alors l'hypothèse H_0 est rejetée.
- Sinon, l'hypothèse H_0 n'est pas rejetée.

Exemple (suite) : on reprend l'exemple précédent. On lance n fois la pièce et on note X_i la variable valant 1 si on obtient pile et 0 sinon.

La variable $\bar{X}_n = \frac{X_1 + \dots + X_n}{n}$ donne la proportion de piles parmi l'échantillon de lancers et sous l'hypothèse $H_0 : p = \frac{1}{2}$, $n\bar{X}_n$ suit une loi $\mathcal{B}\left(n, \frac{1}{2}\right)$.

Travail demandé

1. Établir avec l'inégalité de Bienaymée-Tchebychev que sous l'hypothèse $H_0 : p = \frac{1}{2}$ on a :

$$\mathbb{P}\left(\bar{X}_n \in \left[\frac{1}{2} - \varepsilon, \frac{1}{2} + \varepsilon\right]\right) \geq 1 - \frac{1}{4n\varepsilon^2}.$$

2. En déduire que l'intervalle $\left[\frac{1}{2} - \varepsilon, \frac{1}{2} + \varepsilon\right]$ est un intervalle de non-rejet avec un risque consenti inférieur à $\alpha = \frac{1}{4n\varepsilon^2}$.
3. On fixe $\varepsilon = 0.01$. Combien de tirages doit-on faire pour avoir un seuil de 5% ? Commenter.

2 Le théorème central limite

La zone de rejet obtenu dans le paragraphe précédent avec l'inégalité de Bienaymée-Tchebychev est assez mauvaise car pour avoir une précision acceptable (ε petit) et un seuil α de 5% il faut faire un (trop) grand nombre de tirages (si on pense à des sondages, on ne peut pas se permettre d'interroger 50 000 personnes).

Le but de cette partie est de visualiser expérimentalement ce qu'on appelle le théorème central limite (qu'on verra dans un chapitre ultérieur) puis de voir comment il permet d'améliorer le test précédent.

2.1 Approche expérimentale

Soit $(X_n)_{n \in \mathbb{N}}$ une suite de variables aléatoires indépendantes et de même loi dont on note μ l'espérance et σ l'écart-type, que l'on suppose non nul. On note pour tout $n \in \mathbb{N}^*$:

$$S_n = \frac{\sum_{k=1}^n X_k - n\mu}{\sigma\sqrt{n}} = \frac{\bar{X}_n - \mu}{\sqrt{\frac{\sigma^2}{n}}}.$$

On note f une densité d'une variable de loi $\mathcal{N}(0, 1)$ et Φ sa fonction de répartition (voir l'annexe).

Travail demandé

1. Écrire une fonction `f(x)` et qui renvoie $f(x)$.
2. Écrire une fonction `S_poisson(n, lambda_)` qui simule une réalisation de S_n lorsque les variables aléatoires $(X_n)_n$ suivent la loi $\mathcal{P}(\lambda)$.

Indication : voir annexe.

3. Simuler 10 000 réalisations de S_{1000} (pour différentes valeurs de `lambda_`) et afficher l'histogramme des valeurs obtenues.
Sur le même graphique, afficher la courbe de f .
4. Reprendre les deux questions précédentes lorsque les variables aléatoires $(X_n)_n$ suivent la loi :
 - géométrique ;
 - exponentielle.
5. Que constate-t-on ?

2.2 Le théorème central limite

Le phénomène observé sur les exemples précédents peut se traduire de la façon suivante :

$$\forall -\infty \leq a \leq b \leq +\infty, \quad \lim_{n \rightarrow +\infty} \mathbb{P}(a \leq S_n \leq b) = \int_a^b \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dt = \Phi(b) - \Phi(a).$$

Cela signifie que la suite $(S_n)_n$ « tend¹ » vers une variable aléatoire suivant une loi normale centrée réduite.

Pour n assez grand (en pratique $n > 30$), on pourra approximer $\mathbb{P}(a \leq S_n \leq b)$ par $\Phi(b) - \Phi(a)$.

Exemple : on reprend le même exemple. Remarquons que pour tout n on a :

$$S_n = \sqrt{n} \frac{\bar{X}_n - p}{\sqrt{p(1-p)}}.$$

On en déduit que sous l'hypothèse $H_0 : p = \frac{1}{2}$:

$$\mathbb{P}\left(\frac{1}{2} - \frac{t}{2\sqrt{n}} \leq \bar{X}_n \leq \frac{1}{2} + \frac{t}{2\sqrt{n}}\right) = \mathbb{P}(-t \leq S_n \leq t) \\ \xrightarrow[n \rightarrow +\infty]{} 2\Phi(t) - 1.$$

En prenant $t = t_\alpha$ tel que $\Phi(t_\alpha) = 1 - \frac{\alpha}{2}$ on obtient :

$$\mathbb{P}\left(\frac{1}{2} - \frac{t_\alpha}{2\sqrt{n}} \leq \bar{X}_n \leq \frac{1}{2} + \frac{t_\alpha}{2\sqrt{n}}\right) = \mathbb{P}(-t_\alpha \leq S_n \leq t_\alpha) \\ \simeq 1 - \alpha.$$

L'intervalle $\left[p_0 - \frac{t_\alpha}{2\sqrt{n}}, p_0 + \frac{t_\alpha}{2\sqrt{n}}\right]$ est donc une zone de non rejet avec un risque consenti (environ) égal à α .

1. Le sens précis derrière cette notion sera vu dans un chapitre ultérieur.

Travail demandé

1. Pour $\alpha = 5\%$, on peut déterminer que $t_\alpha \simeq 1.96$. Déterminer le nombre de tirages nécessaires pour avoir un intervalle de non-rejet de longueur inférieure à 0.02. Comparer avec le nombre de tirages obtenus pour l'intervalle établi avec l'inégalité de Bienaymée-Tchebychev.
2. Sur cahier de prépa, dans le dossier du TP7, télécharger le fichier `piece.txt` puis copier-le dans le dossier où votre TP est enregistré **sans l'ouvrir**. Le fichier contient la probabilité p (en %) que la pièce donne Pile.

- (a) Recopier le code suivant qui permet de lire et stocker la probabilité p dans la variable `p`.

```
1 fichier = open('../piece.txt')
2 p = int(fichier.read())/100
3 fichier.close()
```

Pour jouer le jeu et que la probabilité p demeure inconnue, il est important de ne pas ouvrir le fichier et de ne pas faire afficher la valeur de `p` dans Python.

- (b) Écrire une fonction qui détermine un intervalle de non-rejet de longueur 0.02, de seuil $\alpha = 5\%$ et qui décide si on rejette ou non l'hypothèse $H_0 : p = \frac{1}{2}$.
- (c) Modifier la fonction précédente pour qu'elle affiche également une valeur approchée de p .

Vous pouvez maintenant ouvrir le fichier `piece.txt` et comparer avec les résultats trouvés expérimentalement.

3 Test de conformité à la moyenne sous hypothèse gaussienne

Les méthodes précédentes (inégalités de Markov et théorème central limite) comportent des inconvénients importants :

- la taille des échantillons nécessaires doit être importante ;
- la variance était connue (en fonction du paramètre à estimer).

On va donc chercher à mettre en place un test qui fonctionne sur les petits échantillons même lorsque la variance est inconnue. Pour simplifier, on se place dans le cas où le caractère à étudier X suit une loi normale $\mathcal{N}(\mu, \sigma^2)$ de paramètres μ et σ inconnus.

Contexte : on dispose d'une série statistique $x = (x_k)_{k \in \llbracket 1, n \rrbracket}$ que l'on suppose être une réalisation d'un n -échantillon (X_1, \dots, X_n) d'une certaine variable aléatoire réelle X de loi $\mathcal{N}(\mu, \sigma^2)$.

Un nombre réel m étant donné, on souhaite tester la conformité de la série x à l'hypothèse :

$$H_0 : \mu = m.$$

Si on reprend la démarche du paragraphe 2, il est tentant de considérer les variables :

$$S_n = \frac{\bar{X}_n - \mu}{\sqrt{\frac{\sigma^2}{n}}}$$

pour construire une zone de rejet.

Cependant, le paramètre σ est ici inconnu donc on ne peut pas calculer la réalisation de S_n .

3.1 Moyenne et variance empiriques

Pour contourner ce problème, on va d'abord estimer la valeur de σ grâce à notre échantillon (X_1, \dots, X_n) . On note pour tout $n \in \mathbb{N}^*$:

$$\bar{X}_n = \frac{1}{n} \sum_{k=1}^n X_k \text{ (moyenne empirique)} ; \Sigma_n^2 = \frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X}_n)^2 \text{ (variance empirique)}.$$

Travail demandé

1. Écrire une fonction `moyenne` qui prend en entrée une liste et renvoie la moyenne empirique.
2. Écrire une fonction `variance` qui prend en entrée une liste et renvoie la variance empirique.
3. Pour de grandes valeurs de n ($n = 100, 1000$) calculer la moyenne et la variance empiriques de n -échantillons de loi normales de paramètre $\mathcal{N}(\mu, \sigma^2)$ pour différentes valeurs de μ et σ^2 .

Que constate-t-on ?

3.2 Loi du χ^2 , loi de Student

Les observations ci-dessus poussent à considérer les variables suivantes $T_n = \frac{\bar{X}_n - \mu}{\sqrt{\frac{\Sigma_n^2}{n}}}$

plutôt que S_n . Le principe du test que nous allons établir repose alors sur le fait que la loi de T_n est connue :

Théorème .0

Soit Z, X_1, \dots, X_n des variables aléatoires indépendantes de loi $\mathcal{N}(0, 1)$. Alors :

1. $U = \sum_{k=1}^n X_k^2$ est à densité de densité : $f_U : x \mapsto \begin{cases} c_k x^{\frac{k}{2}-1} e^{-\frac{x}{2}} & \text{si } x > 0 \\ 0 & \text{sinon} \end{cases}$

où c_k est une constante positive (faite pour que $\int_{\mathbb{R}} f_U = 1$).

On dit que U suit **la loi du χ^2 à k degrés de liberté, aussi notée χ_k^2** .

2. $\frac{Z}{\sqrt{\frac{U}{k}}}$ est à densité de densité $f : x \mapsto \frac{d_k}{\left(1 + \frac{x^2}{k}\right)^{\frac{k+1}{2}}}$ où d_k est une constante

positive (faite pour que $\int_{\mathbb{R}} f = 1$).

On dit que cette variable aléatoire **suit la loi de Student à k degrés de liberté**.

On peut alors prouver² le résultat suivant.

Théorème .0

La variable $T_n = \frac{\bar{X}_n - \mu}{\sqrt{\frac{\Sigma_n^2}{n}}}$ suit la loi de Student à $n - 1$ degrés de liberté.

3.3 Test de conformité à la moyenne

Sous l'hypothèse $H_0 : \mu = m$, la loi :

$$T_n = \frac{\bar{X}_n - m}{\sqrt{\frac{\Sigma_n^2}{n}}}$$

est donc connue ; une densité est donnée par :

$$f : x \mapsto \frac{d_{n-1}}{\left(1 + \frac{x^2}{n-1}\right)^{\frac{n}{2}}}.$$

Pour un seuil α donnée, on peut donc évaluer la quantité t_α pour que :

$$\mathbb{P}(-t_\alpha \leq T_n \leq t_\alpha) = 1 - \alpha.$$

Comme :

$$\mathbb{P}(-t_\alpha \leq T_n \leq t_\alpha) = \mathbb{P}\left(m \in \left[\bar{X}_n - t_\alpha \sqrt{\frac{\Sigma_n^2}{n}}; \bar{X}_n + t_\alpha \sqrt{\frac{\Sigma_n^2}{n}}\right]\right)$$

on peut à partir de notre échantillon déterminer l'intervalle $\left[\bar{X}_n - t_\alpha \sqrt{\frac{\Sigma_n^2}{n}}; \bar{X}_n + t_\alpha \sqrt{\frac{\Sigma_n^2}{n}}\right]$
et

- si $m \notin \left[\bar{X}_n - t_\alpha \sqrt{\frac{\Sigma_n^2}{n}}; \bar{X}_n + t_\alpha \sqrt{\frac{\Sigma_n^2}{n}}\right]$, l'hypothèse $H_0 : \mu = m$ est rejetée ;
- sinon l'hypothèse est acceptée.

Travail demandé

1. Écrire une fonction `Student(L, m)` qui prend en entrée un échantillon sous forme de liste `L` et une valeur `m` et qui renvoie la réalisation de la variable T_n sous l'hypothèse $H_0 : \mu = m$.
2. Écrire une fonction `t(alpha, n)` qui renvoie la valeur t_α pour laquelle $\mathbb{P}(-t_\alpha \leq T \leq t_\alpha) = 1 - \alpha$ pour T suivant la loi de Student à $n - 1$ degrés de liberté.

Indications : on pourra se référer à l'annexe.

2. mais c'est très dur!

3. Recopier et compléter le programme suivant :

```

1 sigma2 = # nombre strictement positif tiré au hasard
2
3 def EchantillonGaussien(n=5, mu=2):
4     ''' renvoie une réalisation d'un n-échantillon
5         d'une va de loi N(mu, sigma2)
6     '''
7     # lignes à compléter
8     return # à compléter

```

4. Écrire une fonction `Test` avec :

- en entrées : un échantillon sous forme de liste `L`, une valeur seuil `alpha` et la valeur `m` qu'on teste ;
- en sortie : la valeur de T_n , l'intervalle de non-rejet de seuil `alpha` et si oui ou non l'hypothèse $H_0 : \mu = m$ est acceptée.

Tester votre fonction à l'aide du script précédent.

5. On fixe α (par exemple $\alpha = 5\%$) et on répète $N=1000$ fois l'enchaînement suivant :
- on tire au sort un échantillon `x=EchantillonGaussien(5,2)`
 - on test l'hypothèse $H_0 : \mu = 2$.

Dans combien de cas (approximativement) l'hypothèse H_0 est-elle rejetée ?

Écrire une fonction `verification(alpha=0.05, N=1000, n=5, mu=2)` permettant le vérifier.

Annexe

Lois usuelles à densité

Définition .1

- Soit μ et $\sigma > 0$ deux réels. On dit qu'une variable aléatoire X suit la **loi normale de paramètres μ et σ^2** et on note $X \leftrightarrow \mathcal{N}(\mu, \sigma^2)$ si X a pour densité la fonction f définie par

$$\forall x \in \mathbb{R}, \quad f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}.$$

- Soit $\lambda > 0$. On dit qu'une variable aléatoire X suit la loi **exponentielle de paramètre $\lambda > 0$** et on note $X \leftrightarrow \mathcal{E}(\lambda)$ si X a pour densité la fonction f définie par

$$\forall x \in \mathbb{R}, \quad f(x) = \begin{cases} \lambda e^{-\lambda x} & \text{si } x \geq 0 \\ 0 & \text{sinon} \end{cases}.$$

3.4 Simulations de lois usuelles

En important le module `numpy.random` avec la commande :

```
import numpy.random as rd
```

on peut simuler les lois usuelles :

- `rd.exponential(lambda,n)` : simule un n échantillon de loi $\mathcal{E}(\lambda^{-1})$,
- `rd.normal(mu,sigma,n)` : simule un n échantillon de loi $\mathcal{N}(\mu,\sigma^2)$,
- `rd.poisson(lambda,n)` : simule un n échantillon de loi $\mathcal{P}(\lambda)$,
- ...

En important le module `matplotlib.pyplot` avec la commande :

```
import matplotlib.pyplot as plt
```

la commande

```
hist(L, density = True)
```

permet de tracer l'histogramme des fréquences de la liste `L`.

Loi de Student

En important le module `scipy.stats` avec la commande :

```
import scipy.stats as stats
```

la commande `stats.t.ppf(a,df =n)` donne la valeur t_a pour lequel on a

$$\mathbb{P}(T \leq t_a) = a$$

pour T suivant une loi de Student à n degrés de liberté.

Comme la densité donnée ci-dessus est paire, on a :

$$\mathbb{P}(T \leq -t) = \mathbb{P}(T \geq t)$$

donc

$$\mathbb{P}(-t \leq T \leq t) = \mathbb{P}(T \leq t) - \mathbb{P}(T \leq -t) = 2\mathbb{P}(T \leq t) - 1.$$

Pour $a = 1 - \frac{\alpha}{2}$, on a alors :

$$\mathbb{P}(-t_a \leq T \leq t_a) = 2a - 1 = 1 - \alpha.$$