

ESSEC II 2020

Lorsque l'on effectue des sondages, de nombreux biais statistiques peuvent apparaître : on peut par exemple avoir considéré un échantillon non-représentatif de la population, il peut y avoir un biais dans les réponses des personnes sondées... On va s'intéresser dans ce problème à ce que l'on appelle le biais par la taille : il provient du fait que si l'on choisit une personne au hasard dans la population, celle-ci a plus de chances de faire partie d'une catégorie nombreuse de la population.

Le biais par la taille est la source de nombreux « paradoxes » probabilistes, comme le fait que les gagnants du loto vivent en moyenne plus longtemps (parce que les gagnants sont ceux qui ont pu jouer au loto plus longtemps) ou le fait que vos amis ont en moyenne plus d'amis que vous (car les gens qui ont un très grand nombre d'amis font sûrement partie de vos amis). On verra ici comment formaliser le biais par la taille, et l'utiliser dans différents contextes.

Toutes les variables aléatoires intervenant dans le problème sont définies sur un espace probabilisé $(\Omega, \mathcal{A}, \mathbb{P})$. Pour toute variable aléatoire X , on notera $\mathbb{E}(X)$ son espérance (resp. $\mathbb{V}(X)$ sa variance) lorsqu'elles existent.

Première partie : Biais par la taille, exemples discrets

1. On suppose que le nombre d'enfants dans une famille française est une variable aléatoire X . Pour connaître la loi de X , une idée serait d'interroger les élèves d'une école pour connaître le nombre d'enfants dans leur famille.

On va voir que cette approche introduit un biais, en considérant une situation particulière. Supposons que X suive la loi binomiale de paramètres $n = 10$ et $p = \frac{1}{5}$. On note $p_k = \mathbb{P}([X = k])$ pour $k \in \{0, 1, \dots, 10\}$.

a) (i) Rappeler l'expression de p_k pour $k \in \{0, 1, \dots, 10\}$.

(ii) Que vaut $\mathbb{E}(X)$?

(iii) Donner $\mathbb{V}(X)$, et en déduire $\mathbb{E}(X^2)$.

b) Soit M_k le nombre de familles à k enfants, $M = \sum_{k=0}^{10} M_k$ le nombre total de familles (donc

$p_k = \frac{M_k}{M}$). Soit N_k le nombre total d'enfants (c'est-à-dire dans toute la population) qui font

partie d'une famille à k enfants, et $N = \sum_{k=0}^{10} N_k$ le nombre total d'enfants de la population.

(i) Démontrer : $N_k = k p_k M$.

(ii) Démontrer : $\frac{N}{M} = 2$.

(iii) Montrer que la proportion des enfants provenant d'une famille à k enfants est : $p_k^* = \frac{k p_k}{2}$.

c) On choisit une personne au hasard dans la rue, à qui l'on demande combien d'enfants ses parents ont eu (lui ou elle inclus). On note Y ce nombre d'enfants.

(i) Pour tout entier k élément de $\{1, 2, \dots, 10\}$, démontrer : $\mathbb{P}([Y = k]) = \frac{k p_k}{2}$.

(ii) Démontrer : $\mathbb{E}(Y) = \frac{\mathbb{E}(X^2)}{\mathbb{E}(X)}$.

(iii) En déduire $\mathbb{E}(Y)$ et le comparer à $\mathbb{E}(X)$.

2. Soit X une variable aléatoire à valeurs dans \mathbb{N} , non identiquement nulle et admettant une espérance.

Pour tout entier $i > 0$, on pose : $q_i = \frac{i}{\mathbb{E}(X)} \mathbb{P}([X = i])$.

a) Calculer $\sum_{i=1}^{+\infty} q_i$.

La suite $(q_i)_{i>0}$ définie ci-dessus définit donc bien une loi de probabilité. On considère la variable aléatoire X^* dont la loi est donnée par les q_i , c'est-à-dire, pour tout i entier naturel non nul :

$$\mathbb{P}([X^* = i]) = \frac{i}{\mathbb{E}(X)} \mathbb{P}([X = i])$$

On dit que X^* suit la loi de X biaisée par la taille.

b) On suppose que X admet un moment d'ordre 2. Démontrer : $\mathbb{E}(X^*) = \frac{\mathbb{E}(X^2)}{\mathbb{E}(X)}$

c) En déduire que si $\mathbb{E}(X^2)$ existe, on a : $\mathbb{V}(X) = \mathbb{E}(X) (\mathbb{E}(X^*) - \mathbb{E}(X))$.

d) Conclure : $\mathbb{E}(X^*) \geq \mathbb{E}(X)$

3. a) Soit λ un réel strictement positif. On suppose que X est une variable aléatoire qui suit la loi de Poisson de paramètre λ . Soit X^* une variable aléatoire suivant la loi de X biaisée par la taille.

(i) Donner la loi de X^* .

(ii) Vérifier que X^* suit la même loi que $X + 1$.

b) Réciproquement, on suppose que X est une variable aléatoire à valeurs dans \mathbb{N} admettant une espérance non nulle, telle que X^* et $X + 1$ suivent la même loi.

(i) Montrer que pour tout $k \geq 1$: $\mathbb{P}([X = k]) = \frac{\mathbb{E}(X)}{k} \mathbb{P}([X = k - 1])$.

(ii) Montrer que pour tout k entier naturel : $\mathbb{P}([X = k]) = \frac{(\mathbb{E}(X))^k}{k!} \mathbb{P}([X = 0])$.

(iii) En déduire la loi de X .

4. *Le paradoxe du temps d'attente du bus.*

Soit $n \geq 1$ un entier naturel, et soit X une variable aléatoire à valeurs dans $\{1, \dots, n\}$ telle que pour tout $1 \leq k \leq n$: $\mathbb{P}([X = k]) > 0$. On suppose qu'à un arrêt de bus donné, les intervalles de temps entre deux bus consécutifs, exprimés en minutes, sont des variables aléatoires indépendantes, de même loi que X . Une personne arrive à cet arrêt à un instant aléatoire, et se demande combien de temps elle va attendre.

a) Une première idée est que la personne arrive à un instant uniforme entre deux arrivées de bus, séparées par un intervalle de X minutes. On note T la variable aléatoire qui représente le temps d'attente (à valeurs dans $\{1, \dots, n\}$) et on suppose donc que pour tout entier k élément de $\{1, \dots, n\}$:

$$\mathbb{P}_{[X=k]}([T = j]) = \begin{cases} \frac{1}{k} & \text{si } j \in \{1, \dots, k\} \\ 0 & \text{si } j > k \end{cases}$$

(i) Montrer que pour tout entier $k \in \{1, \dots, n\}$, on a : $\sum_{j=1}^n j \mathbb{P}_{[X=k]}([T = j]) = \frac{k+1}{2}$.

(ii) En déduire : $\sum_{k=1}^n \sum_{j=1}^n j \mathbb{P}([X = k]) \mathbb{P}_{[X=k]}([T = j]) = \frac{\mathbb{E}(X+1)}{2}$.

(iii) Démontrer : $\mathbb{E}(T) = \sum_{j=1}^n \sum_{k=1}^n j \mathbb{P}([X = k]) \mathbb{P}_{[X=k]}([T = j])$.

(iv) Démontrer : $\mathbb{E}(T) = \frac{\mathbb{E}(X + 1)}{2}$.

- b) En réalité, en arrivant à l'arrêt de bus, on « tombe » dans un intervalle entre deux bus de manière proportionnelle à sa taille (plus l'intervalle est long, plus on a de chances de « tomber » dedans) : l'intervalle de temps est X^* , suivant la loi de X biaisée par la taille. Le temps d'attente T^* vérifie donc en fait, pour tout $k \in \{1, \dots, n\}$:

$$\mathbb{P}_{[X^*=k]}([T^* = j]) = \begin{cases} \frac{1}{k} & \text{si } j \in \{1, \dots, k\} \\ 0 & \text{si } j > k \end{cases}$$

(i) Montrer que pour tout entier $k \in \{1, \dots, n\}$, on a : $\sum_{j=1}^n j \mathbb{P}_{[X^*=k]}([T^* = j]) = \frac{k+1}{2}$.

(ii) Démontrer : $\mathbb{E}(T^*) = \sum_{j=1}^n \sum_{k=1}^n j \mathbb{P}([X^* = k]) \mathbb{P}_{[X^*=k]}([T^* = j])$.

(iii) Démontrer : $\mathbb{E}(T^*) = \frac{\mathbb{E}(X^* + 1)}{2}$.

(iv) En déduire qu'on a : $\mathbb{E}(T^*) \geq \mathbb{E}(T)$.

Deuxième partie : Biais par la taille, propriétés

Dans cette partie, on démontre de nombreuses propriétés des variables aléatoires biaisées par la taille.

5. Biais par la taille : le cas de variables à densité.

Soit X une variable aléatoire **positive** de densité f et admettant une espérance $\mathbb{E}(X)$ strictement positive (donc $f(x) = 0$ pour tout x strictement négatif).

On définit la fonction g sur \mathbb{R} par : $g : x \mapsto \frac{x}{\mathbb{E}(X)} f(x)$.

- a) Montrer que g définit une densité d'une variable aléatoire positive.

Soit une variable aléatoire X^* dont une densité est g . On dit que X^* suit *la loi de X biaisée par la taille*.

- b) Soit a un réel strictement positif.

(i) Montrer que la variable aléatoire aX possède pour densité $x \mapsto \frac{1}{a} f\left(\frac{x}{a}\right)$.

(ii) En déduire que $(aX)^*$ et $a \cdot X^*$ possèdent la même loi.

- c) *Une propriété importante.*

Soit $h : [0, +\infty[\rightarrow \mathbb{R}$ une fonction bornée et continue sauf éventuellement en un nombre fini de points. Montrer que $\mathbb{E}(Xh(X))$ est bien défini et :

$$\mathbb{E}(h(X^*)) = \frac{1}{\mathbb{E}(X)} \mathbb{E}(Xh(X))$$

On pose alors la définition suivante (**que la variable X soit à densité ou non**) : si X est une variable aléatoire réelle positive d'espérance $\mathbb{E}(X)$ strictement positive, on dit que **la variable aléatoire positive Y suit la loi de X biaisée par la taille** si on a :

$$\mathbb{E}(h(Y)) = \frac{1}{\mathbb{E}(X)} \mathbb{E}(Xh(X))$$

pour toute fonction $h : [0, +\infty[\rightarrow \mathbb{R}$ bornée et continue sauf éventuellement en un nombre fini de points.

6. Dans cette question, on se fixe $f : \mathbb{R} \rightarrow \mathbb{R}$ et $g : \mathbb{R} \rightarrow \mathbb{R}$ deux fonctions **croissantes**. Soit X une variable aléatoire telle que les espérances $\mathbb{E}(f(X))$, $\mathbb{E}(g(X))$ et $\mathbb{E}(f(X)g(X))$ sont bien définies.

a) Montrer que quels que soient les réels x_1 et x_2 , on a : $(f(x_1) - f(x_2))(g(x_1) - g(x_2)) \geq 0$.

b) Soient X_1 et X_2 deux variables aléatoires indépendantes, de même loi que X . Démontrer :

$$\mathbb{E}\left((f(X_1) - f(X_2))(g(X_1) - g(X_2))\right) = 2\mathbb{E}(f(X)g(X)) - 2\mathbb{E}(f(X))\mathbb{E}(g(X))$$

c) En déduire : $\mathbb{E}(f(X)g(X)) \geq \mathbb{E}(f(X))\mathbb{E}(g(X))$.

7. Dans cette question, on suppose que X est une variable aléatoire positive d'espérance strictement positive, et telle que $\mathbb{E}(X^{m+1})$ existe pour un entier $m \geq 1$ donné.

a) Soit p un entier naturel tel que $1 \leq p \leq m$.

(i) Montrer que pour tout réel $x \geq 0$, on a : $0 \leq x^p \leq 1 + x^{m+1}$.

(ii) Montrer que $\mathbb{E}(X^p)$ existe.

b) Démontrer : $\mathbb{E}(X^{m+1}) \geq \mathbb{E}(X)\mathbb{E}(X^m)$.

c) En déduire : $\mathbb{E}((X^*)^m) \geq \mathbb{E}(X^m)$.

8. Pour A un événement, on note $\mathbb{1}_A$ la variable aléatoire définie par :

$$\mathbb{1}_A : \omega \mapsto \begin{cases} 1 & \text{si } \omega \in A \\ 0 & \text{sinon} \end{cases}$$

Pour tout t réel, on définit la fonction g_t sur \mathbb{R} par : $g_t : x \mapsto \mathbb{1}_{]t, +\infty[}(x)$.

a) Montrer que la fonction $x \mapsto g_t(x)$ est croissante sur \mathbb{R} .

b) Soit X une variable aléatoire positive, admettant une espérance. Montrer que pour tout t réel, $\mathbb{E}(Xg_t(X))$ est bien défini et : $\mathbb{E}(Xg_t(X)) \geq \mathbb{E}(X)\mathbb{P}([X > t])$.

c) Démontrer, pour tout t réel : $\mathbb{P}([X^* > t]) \geq \mathbb{P}([X > t])$.

On dit que X^* domine stochastiquement X .

9. Soit X_1, \dots, X_n des variables aléatoires positives, indépendantes, non nécessairement de même loi. On suppose qu'elles admettent toutes une espérance strictement positive, et on note $\mu_i = \mathbb{E}(X_i)$.

De plus, on pose : $\mu = \sum_{i=1}^n \mu_i$ et $S_n = \sum_{i=1}^n X_i$.

a) Donner $\mathbb{E}(S_n)$.

b) Soit J une variable aléatoire à valeurs dans $\{1, \dots, n\}$, de loi $\mathbb{P}([J = k]) = \frac{\mu_k}{\mu}$. Quelle est la loi de J si les variables aléatoires X_i sont de même loi ?

On considère X_1^*, \dots, X_n^* des variables aléatoires indépendantes, indépendantes de X_1, \dots, X_n telles que, pour tout entier i tel que $1 \leq i \leq n$, X_i^* suive la loi de X_i biaisée par la taille.

Soit aussi J une variable aléatoire de loi $\mathbb{P}([J = k]) = \frac{\mu_k}{\mu}$, indépendante de $X_1, X_1^*, \dots, X_n, X_n^*$.

On considère la variable aléatoire $X_J = \sum_{j=1}^n X_j \mathbb{1}_{[J=j]}$ et on définit $T_n = S_n - X_J + X_J^*$. Autrement

dit, on choisit un indice aléatoire J et, dans la somme $\sum_{i=1}^n X_i$, on remplace X_J par X_J^* .

c) Soit $h : [0, +\infty[\rightarrow \mathbb{R}$ une fonction bornée et continue sauf éventuellement en un nombre fini de points.

(i) Démontrer :
$$h(T_n) = \sum_{i=1}^n h(T_n) \mathbb{1}_{[J=i]} = \sum_{i=1}^n h(S_n - X_i - X_i^*) \mathbb{1}_{[J=i]}.$$

(ii) En déduire :
$$\mathbb{E}(h(T_n)) = \sum_{i=1}^n \mathbb{P}([J = i]) \mathbb{E}(h(S_n - X_i + X_i^*)).$$

d) Pour $i \in \{1, \dots, n\}$, démontrer, pour tout réel s :
$$\mathbb{E}(h(s + X_i^*)) = \frac{1}{\mu_i} \mathbb{E}(X_i h(s + X_i)).$$

On admettra qu'on en déduit l'égalité :
$$\mathbb{E}(h(S_n - X_i + X_i^*)) = \frac{1}{\mu_i} \mathbb{E}(X_i h(S_n)).$$

e) En déduire :
$$\mathbb{E}(h(T_n)) = \frac{\mathbb{E}(S_n h(S_n))}{\mathbb{E}(S_n)}.$$

f) Conclure que T_n suit la loi de S_n biaisée par la taille.

Troisième partie : Applications en Statistique

On s'intéresse maintenant au cas où le biais par la taille peut être utilisé en statistique, pour construire des estimateurs non biaisés. Une compagnie d'électricité possède n clients où n est un entier naturel non nul donné. Lors de l'année écoulée, le $i^{\text{ème}}$ client a payé x_i euros ($x_i > 0$), mais a en réalité consommé une quantité d'électricité correspondant à y_i euros ($y_i > 0$). La compagnie sait combien ses clients ont payé, et elle souhaite estimer le rapport :

$$\theta = \frac{\sum_{i=1}^n y_i}{\sum_{i=1}^n x_i}$$

pour déterminer à quel point elle a mal facturé ses clients.

10. Soit m un entier fixé tel que $1 \leq m \leq n$. On note \mathcal{P}_m l'ensemble des parties $A \subset \{1, \dots, n\}$ de cardinal m . On considère une variable aléatoire R , à valeurs dans \mathcal{P}_m et de loi uniforme, c'est-à-dire telle que pour toute partie $A \in \mathcal{P}_m$:
$$\mathbb{P}([R = A]) = \frac{1}{\binom{n}{m}}.$$

On souhaite écrire un programme pour choisir l'ensemble R au hasard.

a) On considère la procédure suivante : on prend un premier élément s_1 uniformément dans $\{1, \dots, n\}$, puis un deuxième élément s_2 uniformément dans $\{1, \dots, n\} \setminus \{s_1\}$, etc. puis un $m^{\text{ème}}$ élément s_m uniformément dans $\{1, \dots, n\} \setminus \{s_1, \dots, s_{m-1}\}$. On note $S = (s_1, \dots, s_m)$, qui est un m -uplet aléatoire.

(i) Montrer que pour tout m -uplet (a_1, \dots, a_m) d'entiers distincts de $\{1, \dots, n\}$, on a :

$$\mathbb{P}([S = (a_1, \dots, a_m)]) = \frac{(n-m)!}{n!}$$

(ii) On note $R = \{s_1, \dots, s_m\}$ l'ensemble des entiers tirés lors de la procédure décrite plus haut (l'ordre dans lequel ils ont été tirés n'importe plus). Montrer que pour tout ensemble

$$A = \{a_1, \dots, a_m\} \subset \{1, \dots, n\} \text{ de cardinal } m, \text{ on a : } \mathbb{P}([R = A]) = \frac{m!(n-m)!}{n!}.$$

En déduire que l'ensemble R a été choisi uniformément dans \mathcal{P}_m .

b) Pour un réel x , on note $[x]$ sa partie entière, c'est-à-dire le plus grand entier naturel inférieur ou égal à x . Montrer que si U suit la loi uniforme sur $[0, 1[$, alors $X = 1 + [nU]$ suit la loi uniforme sur $\{1, \dots, n\}$.

- c) On rappelle que la fonction `rand()` renvoie un nombre aléatoire de loi uniforme sur $[0, 1[$, et que `floor(x)` renvoie la partie entière de x . Écrire une fonction `Uniforme` en **Scilab** qui prend en argument un entier n , et renvoie un nombre (aléatoire), uniforme sur $\{1, \dots, n\}$.

```
function x = Uniforme(n)
    ...
endfunction
```

- d) Écrire une fonction `Selection`, qui prend en argument un vecteur V et renvoie un élément x de V pris de manière aléatoire parmi tous les éléments de V , ainsi que le vecteur W , égal au vecteur V auquel on a enlevé l'élément x . L'instruction `length(V)` renvoie le nombre d'éléments du vecteur V .

```
function [x, W] = Selection(V)
    n = length(V)
    ...
endfunction
```

- e) Compléter le programme suivant, qui prend en argument deux entiers n et m avec $m \leq n$, et renvoie un vecteur R de m entiers distincts, pris uniformément dans $\{1, \dots, n\}$:

```
function R = Choix(m, n)
    V = 1:n
    R = []
    for i = 1:m
        ...
    end
endfunction
```

11. Pour une partie $A \in \mathcal{P}_m$, on définit :

$$\bar{x}_A = \frac{1}{m} \sum_{i \in A} x_i, \quad \bar{y}_A = \frac{1}{m} \sum_{i \in A} y_i, \quad \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

La compagnie décide d'utiliser $\theta_R = \frac{\bar{y}_R}{\bar{x}_R}$ comme estimateur de θ .

- a) On définit deux variables aléatoires $X = \bar{x}_R = \frac{1}{m} \sum_{i \in R} x_i$ et $Y = \bar{y}_R = \frac{1}{m} \sum_{i \in R} y_i$, qui correspondent aux montants moyens payés et consommés par les m clients du groupe tiré au hasard.

(i) Démontrer : $\mathbb{E}(X) = \binom{n}{m}^{-1} \sum_{A \in \mathcal{P}_m} \bar{x}_A$.

(ii) Soit $1 \leq i \leq n$ un entier naturel. Calculer le nombre de parties $A \in \mathcal{P}_m$ telles que $i \in A$.

(iii) En déduire :

$$\sum_{A \in \mathcal{P}_m} \sum_{i \in A} x_i = \binom{n-1}{m-1} \sum_{i=1}^n x_i$$

(iv) Conclure : $\mathbb{E}(X) = \bar{x}$. On **admettra** que de même on a : $\mathbb{E}(Y) = \bar{y}$.

(v) Exprimer θ en fonction de $\mathbb{E}(X)$ et $\mathbb{E}(Y)$.

- b) Démontrer : $\mathbb{E}(\theta_R) = \mathbb{E}\left(\frac{Y}{X}\right)$.

c) On donne l'inégalité de Cauchy-Schwarz : si W et Z sont deux variables aléatoires strictement positives, admettant un moment d'ordre deux : $\mathbb{E}(WZ) \leq (\mathbb{E}(W^2))^{\frac{1}{2}} (\mathbb{E}(Z^2))^{\frac{1}{2}}$, avec égalité si et seulement s'il existe un $\alpha > 0$ tel que $W = \alpha Z$.

(i) Démontrer : $\mathbb{E}\left(\frac{1}{X}\right) \geq \frac{1}{\mathbb{E}(X)}$.

(ii) Montrer qu'il y a égalité si et seulement si X est une variable aléatoire constante, c'est-à-dire $X = \mathbb{E}(X) = \bar{x}$.

(iii) Conclure que $\mathbb{E}\left(\frac{1}{X}\right) = \frac{1}{\mathbb{E}(X)}$ si et seulement si pour tout $i \in \llbracket 1, n \rrbracket$, $x_i = \bar{x}$.

d) Si on suppose que X et Y sont indépendantes, montrer que $\mathbb{E}(\theta_R) \geq \theta$, avec égalité si et seulement si pour tout $i \in \llbracket 1, n \rrbracket$, $x_i = \bar{x}$.

Ainsi, $\mathbb{E}(\theta_R)$ n'est pas forcément égal à θ : on dit alors que θ_R est un estimateur *biaisé* de θ .

12. Ce problème peut être résolu en choisissant les m clients non de manière uniforme comme dans la question 10., mais de manière biaisée par la taille. Par analogie avec la construction de T_n dans la question 9., on commence par choisir une variable aléatoire J à valeurs dans $\{1, 2, \dots, n\}$, dont la loi est donnée par : $\mathbb{P}([J = i]) = \frac{x_i}{\sum_{r=1}^n x_r}$. Ensuite, étant donné J , on choisit un groupe V de $m - 1$

clients parmi les $n - 1$ clients différents de J , de manière uniforme. Autrement dit, pour toute partie $A \in \mathcal{P}_m$, et tout $i \in A$, on a :

$$\mathbb{P}_{[J=i]}([V = A \setminus \{i\}]) = \frac{1}{\binom{n-1}{m-1}}$$

Le groupe de clients examiné est alors : $R = V \cup \{J\}$.

a) On commence par déterminer $\mathbb{P}([R = A])$, pour $A \in \mathcal{P}_m$ donné.

(i) Démontrer :

$$\mathbb{P}([R = A]) = \sum_{i \in A} \mathbb{P}([J = i]) \mathbb{P}_{[J=i]}([V = A \setminus \{i\}])$$

(ii) En déduire :

$$\mathbb{P}([R = A]) = \frac{1}{\binom{n}{m}} \frac{\bar{x}_A}{\bar{x}}$$

13. Une fois choisi le groupe de clients R (par la procédure de la question 12.), on définit : $\hat{\theta}_R = \frac{\bar{y}_R}{\bar{x}_R}$.

a) Démontrer :

$$\mathbb{E}(\hat{\theta}_R) = \frac{1}{\binom{n}{m}} \sum_{A \in \mathcal{P}_m} \frac{\bar{y}_A}{\bar{x}}$$

b) Conclure : $\mathbb{E}(\hat{\theta}_R) = \theta$. On a donc construit un estimateur non biaisé de θ .