

Estimations

I Echantillon

Définition.

Un échantillon désigne un sous ensemble d'individus extraits d'une population, sur lesquels vont être étudiées des grandeurs aléatoires.

Soit une expérience aléatoire à laquelle on associe une v. a. r. X , tq. $E(X) = m$ et $V(X) = \sigma^2$. Un échantillon de taille n ou un n -échantillon est une suite de v. a. r. X_1, \dots, X_n indépendantes et de même loi que X .

On appelle moyenne empirique sur cet échantillon, la v. a. r. $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ et on a :

$$E(\bar{X}_n) = m \text{ et } V(\bar{X}_n) = \frac{\sigma^2}{n}$$

On appelle variance empirique sur cet échantillon, la v. a. r. $S_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$ on peut écrire la variance empirique d'une autre façon, en effet :

$$\begin{aligned} \sum_{i=1}^n (X_i - \bar{X}_n)^2 &= \sum_{i=1}^n ((X_i - m) - (\bar{X}_n - m))^2 \\ &= \sum_{i=1}^n (X_i - m)^2 + n(\bar{X}_n - m)^2 - 2(\bar{X}_n - m) \sum_{i=1}^n (X_i - m) \\ &= \sum_{i=1}^n (X_i - m)^2 + n(\bar{X}_n - m)^2 - 2(\bar{X}_n - m)[n\bar{X}_n - nm] \\ &= \sum_{i=1}^n (X_i - m)^2 - n(\bar{X}_n - m)^2 \end{aligned}$$

$$\text{d'où } S_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - m)^2 - (\bar{X}_n - m)^2$$

$$\begin{aligned} \text{et donc } E(S_n^2) &= \frac{1}{n} \sum_{i=1}^n E(X_i - E(X_i))^2 - E(\bar{X}_n - E(\bar{X}_n))^2 \\ &= \frac{1}{n} \cdot n V(X) - V(\bar{X}_n) = V(X) - \frac{\sigma^2}{n} \end{aligned}$$

$$\text{d'où } E(S_n^2) = \frac{n-1}{n} \sigma^2$$

La valeur moyenne de la variance empirique n'est pas exactement égale à la variance théorique, c'est

pourquoi on introduit la variance empirique modifiée :

$$\bar{S}_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = \frac{n}{n-1} \underbrace{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2}_{S_n^2}$$

d'où $E(\bar{S}_n^2) = \sigma^2$.

II Estimation

On considère généralement deux types d'estimation : l'estimation ponctuelle (on cherche à estimer une valeur) et l'estimation par intervalle de confiance où l'on estime la probabilité que la valeur vraie d'un paramètre appartienne à un intervalle donné.

II.1 Estimation ponctuelle

Problème : Soit X une v. a. r. qui suit une loi de probabilité L_θ dépendant du paramètre θ inconnu. On cherche à utiliser l'information connue dans l'échantillon (X_1, \dots, X_n) pour estimer θ .

Définition.

- 1) Etant donné un échantillon (X_1, \dots, X_n) , on appelle statistique toute v. a. r. T_n de la forme $T_n = \phi(X_1, \dots, X_n)$ où ϕ est une fonction mesurable de $\mathbb{R}^n \rightarrow \mathbb{R}$.
- 2) On appelle estimateur de θ , toute statistique T_n dont la valeur observée sur l'échantillon est utilisée pour estimer θ .

Exemple Si le paramètre à estimer est la moyenne théorique de la loi de X c-a-d $\theta = E(X)$ alors l'estimateur naturel est $T_n(X_1, \dots, X_n) = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X}_n$ puisque $E(\bar{X}_n) = \theta$.

II.2 Biais d'un estimateur

Définition.

- 1) Un estimateur T_n est dit sans biais ou non biaisé si :

$$E(T_n) = \theta.$$

- 2) Un estimateur T_n est dit asymptotiquement sans biais si :

$$\lim_{n \rightarrow +\infty} E(T_n) = \theta.$$

Exemple

♣ \bar{X}_n est un estimateur sans biais de la moyenne de X .

♣ $S_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$ est un estimateur asymptotiquement sans biais de la variance théorique de

X ; en effet si on pose $\theta = V(X)$ alors on sait que : $E(S_n^2) = \frac{n-1}{n} \theta$ et donc $\lim_{n \rightarrow +\infty} E(T_n) = \theta$. Par contre

$\bar{S}_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$ est un estimateur sans biais de la variance théorique puisque $E(\bar{S}_n^2) = \theta$.

II.3 Convergence d'un estimateur

Définition.

Un estimateur T_n est convergent s'il tend en probabilité vers θ quand $n \rightarrow +\infty$ c-a-d :

$$\forall \varepsilon > 0, P(|T_n - \theta| < \varepsilon) \xrightarrow[n \rightarrow +\infty]{} 1.$$

ou encore

$$\forall \varepsilon > 0, P(|T_n - \theta| > \varepsilon) \xrightarrow[n \rightarrow +\infty]{} 0.$$

Théorème. Tout estimateur sans biais dont la variance tend vers 0 quand n tend vers $+\infty$ est convergent c-a-d :

$$\begin{cases} E(T_n) = \theta \\ V(T_n) \xrightarrow[n \rightarrow +\infty]{} 0 \end{cases} \implies T_n \xrightarrow[n \rightarrow +\infty]{P} \theta$$

Preuve L'inégalité de Bienaymé-Tchebychev donne :

$$P(|T_n - \theta| > \varepsilon) \leq \frac{V(T_n)}{\varepsilon^2} \xrightarrow[n \rightarrow +\infty]{} 0.$$

II.4 Qualité d'un estimateur

La qualité d'un estimateur se mesure à l'aide d'une distance au paramètre estimé. Pour obtenir un indicateur numérique on retient $E[(T_n - \theta)^2]$ qu'on appelle risque quadratique, on a la relation suivante :

$$\begin{aligned} EQ(T_n) &= E[(T_n - \theta)^2] = E\left[\left[(T_n - E(T_n)) + (E(T_n) - \theta)\right]^2\right] \\ &= E[(T_n - E(T_n))^2] + 2(E(T_n) - \theta) \underbrace{E[T_n - E(T_n)]}_{E(T_n) - E(T_n) = 0} + (E(T_n) - \theta)^2 \\ &\implies EQ(T_n) = V(T_n) + (E(T_n) - \theta)^2 \end{aligned}$$

Pour un estimateur sans biais l'erreur quadratique est égale à la variance.

Pour deux estimateurs T_n et T'_n sans biais, on dira que T_n est plus efficace que T'_n si $V(T_n) \leq V(T'_n)$. La question se pose alors si on peut trouver un troisième estimateur qui serait meilleur que T_n et si oui, il faudrait poursuivre la recherche ce qui nous conduit à essayer d'améliorer indéfiniment un estimateur. Existe-il un meilleur estimateur ? Pour répondre à cette question on a besoin de l'inégalité de Fréchet-Darmonis-Cramer-Rao (FDCR) suivante :

Définition.

On appelle vraisemblance (Likelihood) de l'échantillon (X_1, \dots, X_n) la loi de probabilité de ce n -uplet, notée $L(x_1, \dots, x_n)$ et est définie par $L(x_1, \dots, x_n) = \prod_{i=1}^n P(X_i = x_i | \theta)$ si X est une v. a. r. discrète et par

$$L(x_1, \dots, x_n) = \prod_{i=1}^n f(x_i, \theta) \text{ si } X \text{ est une v. a. r. continue de densité } f(x, \theta).$$

Le théorème suivant va préciser la borne inférieure pour la variance des estimateurs sans biais, sous certaines hypothèses relatives à la loi de X et que nous appellerons hypothèses de Cramer-RAO (ces hypothèses portent sur l'existence de dérivée de f densité de X et la possibilité d'invertir les opérateurs de dérivation et d'intégration).

Définition.

On appelle quantité d'information de Fisher :

$$I_n(\theta) = E\left[\left(\frac{\partial \ln L}{\partial \theta}\right)^2\right] = E\left[-\frac{\partial^2 \ln L(x_1, \dots, x_n)}{\partial \theta^2}\right] \text{ expression plus facile à calculer}$$

↑ si $X(\Omega)$ est indépendant de θ

Remarque

Si les v. a. r. sont indépendantes et de même loi on a $I_n(\theta) = n I_1(\theta)$ si $E \left[\frac{\partial \ln L(x_1, \dots, x_n)}{\partial \theta} \right] = 0$

Preuve

Sous les mêmes hypothèses et puisque :

$$E \left[\frac{\partial \ln L(x_1, \dots, x_n)}{\partial \theta} \right] = 0 \quad \text{alors :}$$

$$\begin{aligned} I_n(\theta) &= E \left[\left(\frac{\partial \ln L}{\partial \theta} \right)^2 \right] - \left(E \left[\frac{\partial \ln L}{\partial \theta} \right] \right)^2 = \text{var} \left(\frac{\partial \ln L}{\partial \theta} \right) \\ &= \text{var} \left[\frac{\partial}{\partial \theta} \ln \left(\prod_{i=1}^n f(x_i, \theta) \right) \right] = \text{var} \left[\frac{\partial}{\partial \theta} \sum_{i=1}^n \ln(f(x_i, \theta)) \right] \\ &= \text{var} \left[\sum_{i=1}^n \frac{\partial}{\partial \theta} \ln(f(x_i, \theta)) \right] = \sum_{i=1}^n \text{var} \left[\frac{\partial}{\partial \theta} \ln(f(x_i, \theta)) \right] \\ &= n \text{var} \left[\frac{\partial}{\partial \theta} \ln(f(x_1, \theta)) \right] = n I_1(\theta) \end{aligned}$$

Théorème. *Sous les hypothèses de Cramer-RAO, en particulier si $X(\Omega)$ est indépendant du paramètre à estimer θ ; pour tout estimateur sans biais T_n de θ on a*

$$V(T_n) \geq \frac{1}{I_n(\theta)}.$$

Ce théorème donne une borne inférieure pour la variance d'un estimateur sans biais, qui peut ou non être atteinte. Si cette borne est effectivement atteinte pour un estimateur, il sera alors le meilleur ou l'optimal.

Définition.

Un estimateur sans biais T_n est dit efficace si sa variance est égale à la borne de FDCR c-a-d

$$V(T_n) = \frac{1}{I_n(\theta)}.$$

III Méthode de construction d'un estimateur

En l'absence d'un estimateur évident, on est amené à construire un estimateur, on dispose de deux méthodes : la méthode du maximum de vraisemblance et celle des moments.

III.1 Méthode du maximum de vraisemblance

Définition. *L'estimation de maximum de vraisemblance de θ est la valeur $\hat{\theta}_n$ de θ qui rend maximale la fonction de vraisemblance $L(x_1, \dots, x_n, \theta)$. L'estimateur du maximum de vraisemblance (emv) de θ est la v. a. r. correspondante.*

Remarque

$\hat{\theta}_n$ est le maximum de la fonction $L(x_1, \dots, x_n, \theta)$ donc $\hat{\theta}_n$ est solution de :

$$(I) \begin{cases} \frac{\partial L}{\partial \theta} = 0 \\ \frac{\partial^2 L}{\partial \theta^2} < 0 \end{cases}$$

mais comme la fonction de vraisemblance $L(x_1, \dots, x_n, \theta)$ se calcule à partir d'un produit on remplace le problème (I) par un problème équivalent pour la log-vraisemblance c-a-d $\hat{\theta}_n$ est solution de :

$$(II) \begin{cases} \frac{\partial \ln L}{\partial \theta} = 0 \\ \frac{\partial^2 \ln L}{\partial \theta^2} < 0 \end{cases}$$

Exemple

Considérons X une v. a. r. de loi exponentielle de paramètre $\frac{1}{\theta}$ où $\theta > 0$, de densité définie pour $x > 0$ par : $f(x, \theta) = \frac{1}{\theta} e^{-\frac{x}{\theta}}$

La fonction de vraisemblance est $L(x_1, \dots, x_n, \theta) = \prod_{i=1}^n f(x_i, \theta) = \frac{1}{\theta^n} e^{-\frac{1}{\theta} \sum_{i=1}^n x_i}$

d'où $\ln L(x_1, \dots, x_n, \theta) = -n \ln(\theta) - \frac{1}{\theta} \sum_{i=1}^n x_i$ et $\frac{\partial \ln L}{\partial \theta} = \frac{-n}{\theta} + \frac{1}{\theta^2} \sum_{i=1}^n x_i$

ainsi $\frac{\partial \ln L}{\partial \theta} = 0 \iff \theta = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}_n$; d'autre part $\frac{\partial^2 \ln L}{\partial \theta^2} = \frac{n}{\theta^2} - \frac{2}{\theta^3} \sum_{i=1}^n x_i = \frac{n}{\theta^3} (\theta - 2\bar{x}_n)$. Pour $\theta = \bar{x}_n$ on a $\left(\frac{\partial^2 \ln L}{\partial \theta^2} \right) (\bar{x}_n) = \frac{-n}{\bar{x}_n^2} < 0$ d'où $\hat{\theta}_n = \bar{x}_n$ et par conséquent $T_n = \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$

III.2 Méthode des moments

Dans le cas où le paramètre à estimer θ est la moyenne théorique de la r. a. r. X , l'estimateur naturel est la moyenne empirique c-a-d $T_n = \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$; de même pour estimer la variance $\theta = V(X)$, on prend comme estimateur la variance empirique $S_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$. Ce principe se généralise aux moments d'ordre supérieur en considérant $T_n = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^k$ pour $k \geq 3$.

IV Intervalles de confiance

Jusqu'à présent on a estimé un paramètre θ par une unique valeur $\hat{\theta}_n$ (estimation ponctuelle). Si l'estimateur $\hat{\theta}_n$ possède de bonnes propriétés (sans biais, de variance minimale, efficace) on peut s'attendre à ce que $\hat{\theta}_n$ soit proche de la vraie valeur de θ . Cependant il est peu probable que ce soit le cas. Il est plus raisonnable de proposer un ensemble de valeurs vraisemblables pour θ plutôt que de l'estimer par une seule valeur. Fournir un intervalle I tq. $\theta \in I$ s'appelle donner une estimation par intervalle de θ ou estimation ensembliste.

Définition. *Un intervalle de confiance de seuil (ou de niveau de signification) $\alpha \in [0, 1]$ pour un paramètre θ , est un intervalle aléatoire I tq. $P(\theta \in I) = 1 - \alpha$.*

Remarque

✱ α est la probabilité que θ n'appartienne pas à I c-a-d c'est la probabilité de se tromper en affirmant que $\theta \in I$.

✱ $1 - \alpha$ est le niveau de confiance que θ appartienne à I .

Le problème à régler est de trouver un procédé pour déterminer un intervalle de confiance pour un paramètre θ . On choisira dans la plupart des cas un intervalle de confiance centré sur un estimateur performant $\hat{\theta}_n$ c-a-d de la forme $I = [\hat{\theta}_n - \varepsilon, \hat{\theta}_n + \varepsilon]$; il reste à trouver ε tq. $P(\theta \in I) = P(\hat{\theta}_n - \varepsilon \leq \theta \leq \hat{\theta}_n + \varepsilon) = 1 - \alpha$.

IV.1 Intervalles associés aux paramètres de la loi normale

IV.1.1 Intervalle pour la moyenne d'une loi normale avec variance connue

Si X_1, \dots, X_n sont des v. a. r. indépendantes et de même loi $\mathcal{N}(m, \sigma^2)$, on sait que $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ est un estimateur sans biais, de variance minimale de m . Cherchons un intervalle de confiance pour m sous la forme $[\bar{X}_n - \varepsilon, \bar{X}_n + \varepsilon]$ le problème revient, pour α fixé, à chercher ε tq. $P(\bar{X}_n - \varepsilon \leq m \leq \bar{X}_n + \varepsilon) = 1 - \alpha$ c-a-d $P(-\varepsilon \leq \bar{X}_n - m \leq \varepsilon) = 1 - \alpha$ ou bien $P\left(-\frac{\sqrt{n}\varepsilon}{\sigma} \leq \sqrt{n} \frac{\bar{X}_n - m}{\sigma} \leq \frac{\sqrt{n}\varepsilon}{\sigma}\right) = 1 - \alpha$

on sait que $\sqrt{n} \frac{\bar{X}_n - m}{\sigma} \xrightarrow[n \rightarrow +\infty]{loi} \mathcal{N}(0, 1)$. Si ϕ est la fonction de répartition de la loi $\mathcal{N}(0, 1)$ alors ε est

$$\begin{aligned} \text{solution de : } 1 - \alpha &= \phi\left(\frac{\sqrt{n}\varepsilon}{\sigma}\right) - \phi\left(-\frac{\sqrt{n}\varepsilon}{\sigma}\right) \\ &= 2\phi\left(\frac{\sqrt{n}\varepsilon}{\sigma}\right) - 1 \end{aligned}$$

$$\text{d'où } \phi\left(\frac{\sqrt{n}\varepsilon}{\sigma}\right) = 1 - \frac{\alpha}{2} \implies \frac{\sqrt{n}\varepsilon}{\sigma} = \phi^{-1}\left(1 - \frac{\alpha}{2}\right)$$

Si on pose $u_\alpha = \phi^{-1}\left(1 - \frac{\alpha}{2}\right) \iff 1 - \frac{\alpha}{2} = \phi(u_\alpha)$, u_α s'appelle le fractile d'ordre $1 - \frac{\alpha}{2}$ de la loi $\mathcal{N}(0, 1)$.

alors $\frac{\sqrt{n}\varepsilon}{\sigma} = u_\alpha \implies \varepsilon = \frac{\sigma}{\sqrt{n}}u_\alpha$ et par conséquent l'intervalle de confiance pour m de niveau de confiance

$$1 - \alpha \text{ est : } I = \left[\bar{X}_n - \frac{\sigma}{\sqrt{n}}u_\alpha, \bar{X}_n + \frac{\sigma}{\sqrt{n}}u_\alpha \right]$$

IV.1.2 Intervalle pour la moyenne d'une loi normale avec variance inconnue

La statistique utilisée dans la situation précédente n'est utilisable que si on connaît la valeur de σ ; or, dans la pratique on ne connaît pas la valeur de σ et va donc devoir être remplacé par un estimateur sans biais de la variance théorique σ^2 : $\bar{S}_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$. On utilise donc comme nouvelle statistique

$\sqrt{n} \frac{\bar{X}_n - m}{\bar{S}_n}$ qui suit une loi de Student à $n - 1$ degrés de liberté.

de $P(\bar{X}_n - \varepsilon \leq m \leq \bar{X}_n + \varepsilon) = 1 - \alpha$ on en déduit que $P\left(-\frac{\sqrt{n}\varepsilon}{\bar{S}_n} \leq \sqrt{n} \frac{\bar{X}_n - m}{\bar{S}_n} \leq \frac{\sqrt{n}\varepsilon}{\bar{S}_n}\right) = 1 - \alpha$

si on appelle F la fonction de répartition de la loi Student $St(n-1)$ à $n - 1$ degrés de liberté on aura $F\left(\frac{\sqrt{n}\varepsilon}{\bar{S}_n}\right) = 1 - \frac{\alpha}{2} \implies \frac{\sqrt{n}\varepsilon}{\bar{S}_n} = F^{-1}\left(1 - \frac{\alpha}{2}\right)$ (*)

Si on pose $t_{n-1, \alpha} = F^{-1}\left(1 - \frac{\alpha}{2}\right)$ c-a-d $F(t_{n-1, \alpha}) = 1 - \frac{\alpha}{2}$, $t_{n-1, \alpha}$ est donc le fractile d'ordre $1 - \frac{\alpha}{2}$ de la

loi de Student à $n - 1$ degrés de liberté; on déduit de (*) que $\varepsilon = \frac{\bar{S}_n}{\sqrt{n}} t_{n-1, \alpha}$; d'où l'intervalle de confiance

$$\text{pour } m \text{ de niveau de confiance } 1 - \alpha \text{ vaut : } I = \left[\bar{X}_n - \frac{\bar{S}_n}{\sqrt{n}} t_{n-1, \alpha}, \bar{X}_n + \frac{\bar{S}_n}{\sqrt{n}} t_{n-1, \alpha} \right]$$

IV.1.3 Intervalle pour la variance d'une loi normale avec espérance connue

$\hat{\sigma}_n = \frac{1}{n} \sum_{i=1}^n (X_i - m)^2$ est le meilleur estimateur de σ^2 de plus $\frac{n\hat{\sigma}_n}{\sigma^2}$ suit la loi du χ_n^2 . On a donc (la loi du χ_n^2 n'est pas symétrique) $\forall a, b \in \mathbb{R}$ tq. $0 < a < b$

$$P\left(a \leq \frac{n\hat{\sigma}_n}{\sigma^2} \leq b\right) = P\left(\frac{n\hat{\sigma}_n}{b} \leq \sigma^2 \leq \frac{n\hat{\sigma}_n}{a}\right) = F_{\chi_n^2}(b) - F_{\chi_n^2}(a)$$

$$1 - \alpha = 1 - \frac{\alpha}{2} - \frac{\alpha}{2}$$

Il y a une infinité de façons de choisir a et b . La façon la plus usuelle de procéder est "d'équilibrer les risques" c-a-d de prendre a et b tq. $F_{\chi_n^2}(b) = 1 - \frac{\alpha}{2}$ et $F_{\chi_n^2}(a) = \frac{\alpha}{2}$ si on obtient $b = \chi_{n, 1-\frac{\alpha}{2}}^2$ et $a = \chi_{n, \frac{\alpha}{2}}^2$ alors l'intervalle de confiance pour σ^2 de niveau de confiance $1 - \alpha$ est :

$$I = \left[\frac{n\hat{\sigma}_n}{\chi_{n, 1-\frac{\alpha}{2}}^2}, \frac{n\hat{\sigma}_n}{\chi_{n, \frac{\alpha}{2}}^2} \right]$$

IV.1.4 Intervalle pour la variance d'une loi normale avec espérance inconnue

Lorsque l'espérance m de la loi normale est inconnue, l'estimateur sans biais et convergent de σ^2 est : $\bar{S}_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$ de plus $(n-1)\frac{\bar{S}_n^2}{\sigma^2}$ suit la loi du χ_{n-1}^2 . De même on cherche les valeurs de a et b tq.

$$P\left(a \leq (n-1)\frac{\bar{S}_n^2}{\sigma^2} \leq b\right) = 1 - \alpha$$

ce qui conduit à :

$$P\left((n-1)\frac{\bar{S}_n^2}{b} \leq \sigma^2 \leq (n-1)\frac{\bar{S}_n^2}{a}\right) = F_{\chi_{n-1}^2}(b) - F_{\chi_{n-1}^2}(a)$$

$$= 1 - \alpha = 1 - \frac{\alpha}{2} - \frac{\alpha}{2}$$

usuellement on choisit $b = \chi_{n-1, 1-\frac{\alpha}{2}}^2$ et $a = \chi_{n-1, \frac{\alpha}{2}}^2$ alors l'intervalle de confiance pour σ^2 de niveau de confiance $1 - \alpha$ est :

$$I = \left[(n-1)\frac{\bar{S}_n^2}{\chi_{n-1, 1-\frac{\alpha}{2}}^2}, (n-1)\frac{\bar{S}_n^2}{\chi_{n-1, \frac{\alpha}{2}}^2} \right]$$