

ESTIMATION PONCTUELLE

L'estimation ponctuelle d'une grandeur (caractéristique) statistique d'une population est le calcul d'une estimée de la valeur de cette grandeur. Cette estimée est associée à un estimateur. En règle générale, on cherche à calculer l'estimée la plus précise, c-à-d. l'estimateur le meilleur.

5.1

Estimation et estimateur

On s'intéresse à la caractéristique X d'une population (éventuellement à un vecteur de caractéristiques), dont la loi dépend d'un paramètre inconnu $\Theta \in \mathcal{G} \subset \mathbb{R}^m$.

On note par

- $P_{\Theta}(X = x)$ la loi de X au point x si la v.a. X est discrète ;
- $f_{\Theta}(x)$ la densité de la loi de X au point x si la v.a. X est continue.

On effectue un tirage de n articles d'une population de taille potentiellement infinie, noté (x_1, \dots, x_n) . On note (X_1, \dots, X_n) l'échantillon aléatoire associé à ce tirage (il s'agit d'un vecteur aléatoire dont une réalisation particulière est (x_1, \dots, x_n)).

Avoir une estimation de la valeur θ de la caractéristique Θ consiste à avoir une valeur approchée $\hat{\theta}$ de cette caractéristique en utilisant le tirage.

DÉFINITION 5.0.7 Soit $\Theta \in \mathcal{G} \subset \mathbb{R}^m$ la caractéristique de valeur inconnue de la loi qu'elle suit la v.a. X .

Un estimateur $\hat{\Theta}_n$ de Θ est une statistique de l'échantillon

$$\hat{\Theta}_n = g(X_1, \dots, X_n); \quad g : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

telle que pour chaque échantillon (x_1, \dots, x_n) la valeur $\hat{\theta}_n = g(x_1, \dots, x_n)$ de $\hat{\Theta}_n$ est proche de la valeur θ de Θ .

$\hat{\theta}_n$ est une estimation de θ .

À titre d'exemple considérons la moyenne μ_X d'une population. On a donc $\Theta = \mu_X$. Soit (X_1, \dots, X_n) un échantillon de la population. On prendra comme estimateur de la moyenne μ_X , la moyenne empirique de l'échantillon $\hat{\Theta}_n = \bar{X}_n = \frac{1}{n} \sum_{k=1}^n X_k$. Si (x_1, \dots, x_n) est une réalisation de l'échantillon, alors l'estimation a comme valeur $\hat{\theta}_n = \bar{x}_n = \frac{1}{n} \sum_{k=1}^n x_k$.

5.2

Critères de qualité d'un estimateur

Pour choisir le meilleur estimateur il faut se doter de critères de qualité pour les estimateurs.

Estimateur non biaisé

Lorsqu'on prélève plusieurs échantillons de taille n de la population, on obtient plusieurs valeurs différentes de l'estimation $\hat{\theta}_n$ de la caractéristique statistique que nous voulons estimer. Ces valeurs sont des valeurs aléatoires, dans la mesure où l'échantillonnage se fait de manière aléatoire. On peut donc calculer son espérance. Si on souhaite que les estimations nous fournissent la valeur exacte θ de cette caractéristique, on peut imposer que l'estimateur soit *sans biais*, c-à-d. que $E(\hat{\Theta}_n) = \Theta$.

Estimateur efficace

Une estimation $\hat{\theta}_n$ est d'autant plus fiable que sa variance $V(\hat{\theta}_n)$ (c-à-d. la dispersion de ces valeurs autour de $E(\hat{\theta}_n)$) est petite. Un estimateur est *efficace* si la variance de ses estimations est minimale.

Estimateur convergent

Un estimateur est *convergent* si sa variance tend vers zéro lorsque le nombre d'échantillons tend vers l'infini : $\lim_{n \rightarrow \infty} V(\hat{\Theta}_n) = 0$.

5.3

Précision d'un estimateur

On mesure l'erreur quadratique moyenne (EQM) ou risque d'une estimation $\hat{\theta}_n$ à l'aide de la formule

$$R_{\Theta}(\hat{\theta}_n) = E \left[(\hat{\theta}_n - \theta)^2 \right] = V(\hat{\theta}_n) + (E(\hat{\theta}_n) - \theta)^2$$

Nous savons que

$$B_{\Theta}(\hat{\theta}_n) = E(\hat{\theta}_n) - \theta$$

représente le biais de l'estimation. Donc l'EQM de l'estimateur est la variance de l'estimation plus son biais au carré :

$$R_{\Theta}(\hat{\theta}_n) = V(\hat{\theta}_n) + B_{\Theta}^2(\hat{\theta}_n)$$

Par conséquent un estimateur est d'autant plus précis qu'il est de variance minimale est sans biais. On appellera cet estimateur *estimateur sans biais de variance minimale*. Si cet estimateur existe, il est unique presque sûrement.

En règle générale on cherche un estimateur d'EQM minimale. Dans la plupart de cas il faut faire un compromis entre la valeur la variance de l'estimateur et son biais.

5.4

Exemples d'estimateurs

Nous donnons dans la suite une liste d'estimateurs des principales caractéristiques (grandeurs) statistiques.

Espérance $E(X) = \mu_X$ **de la v.a. X** C'est la moyenne empirique $\bar{X}_n = \frac{1}{n} \sum_{k=1}^n X_k$,

c-à-d. nous avons pour l'estimateur $\hat{\Theta}_n = \bar{X}_n$. Pour un échantillon donné

(x_1, \dots, x_n) on a $\bar{x}_n = \frac{1}{n} \sum_{k=1}^n x_k$ et donc l'estimation est $\hat{\theta}_n = \bar{x}_n$.

Nous avons $E(\bar{X}_n) = E(X) = \mu_X$ et $V(\bar{X}_n) = \frac{V(X)}{n} = \frac{\sigma_X^2}{n}$.

Variance $V(X)$ **de la v.a. X lorsque $E(X) = \mu_X$ est connue** C'est la quantité

$T_n^2 = \frac{1}{n} \sum_{k=1}^n (X_k - \mu_X)^2$, c-à-d. nous avons pour l'estimateur $\hat{\Theta}_n = T_n^2$. Pour un

échantillon donné (x_1, \dots, x_n) on a $t_n^2 = \frac{1}{n} \sum_{k=1}^n (x_k - \mu_X)^2$ et donc l'estimation de la variance est $\hat{\theta}_n = t_n^2$.

Nous avons $E(T_n^2) = \sigma_X^2$ et $V(T_n^2) = \frac{\mu_X^4 - \sigma_X^4}{n}$.

Variance $V(X)$ de la v.a. X lorsque $E(X) = \mu_X$ est inconnue C'est la variance

empirique $S_n^2 = \frac{1}{n} \sum_{k=1}^n (X_k - \bar{X}_n)^2$, c-à-d. nous avons pour l'estimateur $\hat{\Theta}_n = S_n^2$.

Pour un échantillon donné (x_1, \dots, x_n) on a $s_n^2 = \frac{1}{n} \sum_{k=1}^n (x_k - \bar{x}_n)^2$ et donc l'estimation de la variance est $\hat{\theta}_n = s_n^2$.

Nous avons $E(S_n^2) = \frac{n-1}{n} \sigma_X^2$ et $V(T_n^2) = \frac{n-1}{n^3} ((n-1)\mu_X^4 - (n-3)\sigma_X^4)$.

Si la variance empirique est corrigée (c-à-d. elle est sans biais), alors $S_n^2 =$

$\frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X}_n)^2$ et dans ce cas $E(S_n^2) = \sigma_X^2$ et $V(T_n^2) = \frac{1}{n} (\mu_X^4 - \frac{n-3}{n-1} \sigma_X^4)$.

Covariance $C(X, Y)$ entre deux v.a. X, Y lorsque $E(X) = \mu_X$ et $E(Y) = \mu_Y$ sont connues

$\hat{\Theta}_n = cov_n(X, Y) = \frac{1}{n} \sum_{k=1}^n (X_k - \mu_X)(Y_k - \mu_Y)$ est un estimateur de la covariance.

Pour deux échantillons donnés (x_1, \dots, x_n) et (y_1, \dots, y_n) on a que $cov_n(x, y) =$

$\frac{1}{n} \sum_{k=1}^n (x_k - \mu_X)(y_k - \mu_Y)$ est la valeur de l'estimation et donc $\hat{\theta}_n = cov_n(x, y)$.

D'après ce que nous venons de voir les estimateurs de l'espérance et celui de la variance corrigée sont sans biais. Par contre la variance de l'estimateur de la variance est plus petite que la variance de l'estimateur de la variance corrigée.

Covariance $C(X, Y)$ entre deux v.a. X, Y lorsque $E(X) = \mu_X$ et $E(Y) = \mu_Y$ sont inconnues

$\hat{\Theta}_n = \widehat{cov}_n(X, Y) = \frac{1}{n} \sum_{k=1}^n (X_k - \bar{X}_n)(Y_k - \bar{Y}_n)$. Pour deux échantillons donnés

(x_1, \dots, x_n) et (y_1, \dots, y_n) on a que $\widehat{cov}_n(x, y) = \frac{1}{n} \sum_{k=1}^n (x_k - \bar{x})(y_k - \bar{y})$ est l'esti-

mation et donc $\hat{\theta}_n = \widehat{cov}_n(x, y)$.

Corrélation $\rho(X, Y)$ entre deux v.a. X, Y $\hat{\Theta}_n = r_n(X, Y) = \frac{\widehat{cov}(X, Y)}{S_n(X)S_n(Y)}$. Pour deux

échantillons donnés (x_1, \dots, x_n) et (y_1, \dots, y_n) on a que $r_n(x, y) = \frac{\sum_{k=1}^n (x_k - \bar{x})(y_k - \bar{y})}{\sum_{k=1}^n (x_k - \bar{x})^2 \sum_{k=1}^n (y_k - \bar{y}_n)^2}$

est la valeur de l'estimation et donc $\hat{\theta}_n = r_n(x, y)$. Notons que $r_n(x, y) \in [-1, 1]$.

5.5

Estimateur des moindres carrés

On cherche à calculer l'estimation $\hat{\theta}_n$ sur la base d'une réalisation (x_1, \dots, x_n) de v.a. (X_1, \dots, X_n) . La méthode des moindres carrés consiste à minimiser la somme des carrés des écarts $x_k - \hat{\theta}_n$ entre les valeurs observées x_k et l'estimation $\hat{\theta}_n$ de l'estimateur, c-à-d. on cherche à minimiser la quantité

$$Q(\hat{\theta}_n) = \sum_{k=1}^n (x_k - \hat{\theta}_n)^2$$

Au minimum nous avons

$$\frac{dQ}{d\hat{\theta}_n} \equiv -2 \sum_{k=1}^n (x_k - \hat{\theta}_n) = 0$$

d'où on peut calculer $\hat{\theta}_n$. De plus la valeur de $\hat{\theta}_n$ doit vérifier la relation

$$\frac{d^2Q}{(d\hat{\theta}_n)^2} > 0$$

Cette méthode ne donne pas de bons résultats si n est modéré. De plus ces estimateurs ne sont pas convergents.

5.6

Exercices

EXERCICE 5.1 À l'occasion d'un test qualité, on a examiné 600 articles et on a trouvé 78 défectueux. On cherche à évaluer la probabilité p qu'un article soit défectueux.

- (1) Proposez un estimateur pour p et justifiez votre réponse. Calculez l'estimation qui correspond au résultat du test.
- (2) Trouver la loi suivie par cet estimateur. Est-il possible d'utiliser l'approximation de la loi normale pour effectuer les calculs numériques ?

EXERCICE 5.2 Lors d'un test qualité on trouve 40% d'articles hors normes. Calculer la probabilité que dans un nouveau test portant sur 400 articles, au moins 150 articles soient hors normes.

EXERCICE 5.3 On considère un échantillon X_1, X_2, X_n tiré d'une population de moyenne μ_X et de variance σ_X^2 . On propose pour la moyenne les deux estimateurs suivants :

$$- \bar{x}_1 = \frac{x_1 + 2x_2 + 3x_3}{6}$$

$$- \bar{x}_2 = \frac{x_1 + x_2 + x_3}{3}$$

Comparer ces deux estimateurs du point de vue de leur biais et de leur efficacité.

EXERCICE 5.4 Soit $X \sim \mathcal{N}(\mu, \sigma)$; $\mu = \sigma = \theta$ inconnu. On veut estimer θ .

Soient $T = X$ et $U = \frac{X}{2}$ deux estimateurs.

Comparer ces deux estimateurs.

EXERCICE 5.5 Soit X une v.a. de densité

$$f(x; \theta) = \begin{cases} \frac{2}{\theta} \left(1 - \frac{x}{\theta}\right) & \text{si } 0 \leq x \leq \theta, \\ 0 & \text{sinon} \end{cases}$$

où θ est un paramètre inconnu strictement positif.

- (1) Calculer $E(X)$ et $V(X)$
- (2) On cherche à estimer θ à partir d'un échantillon (X_1, \dots, X_n) de v.a. indépendantes, de même loi que X . On utilise $T_n = \bar{X}$ comme estimateur.
 - (a) Calculer le biais de cet estimateur. Pouvez-vous lui éliminer son biais ?
 - (b) Établir que T_n modifié sans biais, est un estimateur de θ convergent.
 - (c) Calculer $V(T_n)$ et en conclure que T_n converge en moyenne quadratique vers θ .
 - (d) Déterminer la loi asymptotique de T_n .

EXERCICE 5.6 Nous avons un échantillon de n tirages X_1, \dots, X_n de v.a.i.i.d. qui suivent la loi normale $\mathcal{N}(\mu_X, \sigma_X^2)$. On voudrait estimer μ_X et prévoir le prochain tirage X_{n+1} . La variance est connue et égale à 4.

- (1) Montrer que l'estimateur sans biais le plus précis de μ_X est la moyenne empirique $\bar{X}_n = \frac{1}{n} \sum_{k=1}^n X_k$ et préciser sa loi.
- (2) Le prochain tirage X_{n+1} est supposé indépendant des X_1, \dots, X_n et de même loi. Donner la meilleure prévision de X_{n+1} .
- (3) Trouver la loi de la v.a. $Z_{n+1} = X_{n+1} - \bar{X}_n$ qui représente l'erreur de la prévision.