

## ESTIMATEUR DE HUBER (F6, F8, H4, H5, J)

(14 / 01 / 2020, © Monfort, Dicostat2005, 2005-2020)

(i) Soit  $(\mathcal{X}, \mathcal{B}, P_\theta^X)_{\theta \in \Theta}$  un **modèle statistique** dans lequel :

(a)  $\Theta = \mathbf{R}$  et  $\theta$  représente un **paramètre de position** caractérisant la **lp**  $P_\theta^X$  générant l'**échantillon**  $X = (X_1, \dots, X_N)$  ;

(b)  $\mathcal{X} = \mathbf{R}^N$  et  $X : \Omega \mapsto \mathbf{R}^N$  est un **échantillon iid** constitué de  $N$  **copies**  $X_n$  dont la **variable parente**  $\xi$  admet l'une des lois  $P_\theta^X$  ( $\theta \in \Theta$ ) pour **loi de probabilité**. Autrement dit,  $P_\theta^X = (P_\theta^\xi)^{\otimes N}$ ,  $\forall \theta \in \Theta$  ;

(c) il existe une **densité**  $f_0$  sur  $\mathbf{R}$  tq (**dérivée de NIKODYM-RADON**) :

$$(1) \quad f(x, \theta) = (dP_\theta^\xi / d\lambda)(x) = f_0(x - \theta), \quad \forall (x, \theta) \in \mathcal{X} \times \Theta.$$

où  $\lambda$  désigne la **mesure de LEBESGUE** sur  $\mathbf{R}$ .

Soit  $\rho : \mathbf{R} \mapsto \mathbf{R}$  une fonction (non constante) donnée.

On appelle **estimateur (de type M) de P.J. HUBER**, ou **M-estimateur de P.J. HUBER**, de  $\theta$  une **statistique**  $T_N = t_N(X)$  qui rend minimum pr à  $\theta$  la **va** (cf **équation estimante**) :

$$(3) \quad s_\rho(\theta) = \sum_{n=1}^N \rho(X_n - \theta).$$

(ii) A titre d'exemples :

(a) si  $\rho(u) = u^2$ , on obtient la **moyenne empirique**  $T_N = \bar{X}_N$  ;

(b) si  $\rho(u) = |u|$ , on obtient la **médiane empirique**  $T_N = q_{1/2} X$  ;

(c) si  $\rho(u) = -\text{Log } f_0(u)$ , on obtient l'**estimateur du maximum de vraisemblance** de  $\theta$ . Cet exemple relie ainsi l'estimateur de type M (« maximum » de vraisemblance) à celui du maximum de vraisemblance.

(iii) Sous certaines hypothèses, on montre que la solution  $T_N$  de (3) est unique et que :

(a)  $T_N$  converge presque-sûrement vers  $\theta$  (cf **convergence presque sûre**) :

$$(4) \quad T_N \xrightarrow{N \rightarrow +\infty} \theta \text{ (P}_\theta\text{-p.s.)}, \quad \forall \theta \in \Theta ;$$

(b)  $T_N$  est asymptotiquement gaussien :

$$(5) \quad N^{1/2} (T_N - \theta^\#) \xrightarrow{\mathcal{L}}_{N \rightarrow +\infty} \mathcal{N}(0, \sigma^2) \text{ (loi normale réduite),}$$

où  $\theta^\#$  est la solution en  $\theta$  de  $s_\rho(\theta) = 0$  et  $\sigma^2 = (D s_\rho(\theta^\#))^{-2} \sigma_0^2$ , avec  $\sigma_0^2 = \int \{\rho(x - \theta^\#)\}^2 dF_N(x)$ ,  $F_N$  étant la **fonction de répartition empirique** associée à  $X$ .

(iv) La notion de M-estimateur généralise celle de **moyenne de WINSOR** (cf **transformation de WINSOR**). En effet, si :

$$(6) \quad W_N(L, M) = N^{-1} \{L X^{(L+1)} + X^{(L+1)} + \dots + X^{(N-M)} + M X^{(N-M)}\}$$

est la **moyenne (arithmétique) de WINSOR** d'ordre  $(L, M)$  et  $X^{(\cdot)}$  l'échantillon ordonné associé à  $X$  qui permet de la calculer (cf **statistique ordonnée**), alors l'estimateur de HUBER  $T_N$  est de la forme :

$$(7) \quad T_N = N^{-1} \{L U + X^{(L+1)} + X^{(L+2)} + \dots + X^{(N-M)} + M V\},$$

pour des couples  $(L, M)$  et  $(U, V)$  tq :

$$(8) \quad X^{(L)} \leq U \leq X^{(L+1)} \quad \text{et} \quad X^{(N-M)} \leq V \leq X^{(N-M+1)}.$$

(v) Un M-estimateur peut être utilisé pour détecter une **aberration** éventuelle, ie un « point aberrant » éventuel  $X_n$ , car il est robuste (eg pr à un **mélange de lois**) (cf **robustesse**).

En effet, si  $F_\theta$  est la **fr** associée à la lp  $P_{\theta^\xi}$  ( $\forall \theta \in \Theta$ ), et si cette fonction correspond à une loi normale contaminée (cf **contamination des lois**), ie :

$$(9) \quad F_\theta(x) = (1 - \varepsilon) \Phi(x) + \varepsilon G(x), \quad \forall x \in \mathbf{R},$$

où  $\Phi$  est la fr de la loi  $\mathcal{N}(0, 1)$ ,  $\varepsilon \in [0, 1[$  est donné ( $\varepsilon \leq 1$ ) et  $G$  est une fr inconnue, on montre que le choix de la fonction :

$$(10) \quad \rho(u) = \begin{cases} u^2 / 2 & \text{si } |u| < c, \\ c |u| - (c^2 / 2) & \text{si } |u| \geq c, \end{cases}$$

dans laquelle  $c > 0$  dépend de  $\varepsilon$  selon la relation :

$$(11) \quad 1 - \varepsilon = \{\Phi(c) - \Phi(-c)\} + c^{-1} (2 / \pi)^{1/2} \exp(-c^2 / 2),$$

fournit l'**estimateur de centralité** le plus robuste (cf **robustesse**) au sens où sa **variance** asymptotique est minimum lorsque  $G$  parcourt la classe des fr dont les lp associées sont des **lois symétriques**.

(vi) Les hypothèses portant sur  $\rho$  sont diverses. Parmi les plus usuelles, on peut noter les suivantes :

(a)  $\rho$  est paire et convexe sur  $\mathbf{R}$  (avec  $\lim_{|u| \rightarrow +\infty} \rho(u) = +\infty$ ) ;

(b)  $\rho$  est à valeurs dans  $\mathbf{R}_+$  et non décroissante sur  $\mathbf{R}$  ;

(c)  $\rho$  est symétrique par à 0 et tq  $\rho(0) = 0$  ;

(d)  $\rho$  est différentiable et  $\rho' = D\rho$  est continue (sauf en un nombre fini de points de  $\mathbf{R}$ ).

(vii) Il existe des notions de **M-estimateur généralisé**, notamment au cas où  $\rho$  est définie sur  $\mathbf{R}^K$  et à valeurs dans  $\mathbf{R}$ , et où la somme  $s_\rho(\theta)$  est remplacée par :

$$(12) \quad \sum_{(n(1), \dots, n(K))} \rho(X_{n(1)} - \theta, \dots, X_{n(K)} - \theta),$$

avec  $\{n_1, \dots, n_K\} \subset \{1, \dots, N\}$  et  $\beta \neq \alpha \Rightarrow n_\beta \neq n_\alpha$  (où les  $n(1), \dots, n(K)$  désignent, alternativement, les  $n_1, \dots, n_K$ ).

(viii) Le modèle initial s'étend directement au cas d'un **paramètre d'échelle et de position**. On considère alors un modèle dans lequel :

(a)  $\Theta = \mathbf{R} \times \mathbf{R}_+^*$  et  $\theta = (\alpha, \beta)$  représente un couple constitué d'un **paramètre de position**  $\alpha$  et d'un **paramètre d'échelle**  $\beta$  (cf **paramètre d'échelle et de position**);

(b)  $\mathcal{X} = \mathbf{R}^N$  et  $X$  est défini comme précédemment ;

(c) il existe une densité  $f_0$  sur  $\mathbf{R}$  tq la **dérivée de NIKODYM-RADON** de  $P_{\theta^\xi}$  par à la mesure de LEBESGUE  $\lambda$  est de la forme :

$$(13) \quad f(x, \theta) = \beta^{-1} f_0\{(x - \alpha) / \beta\}, \quad \forall (x, \theta) \in \mathcal{X} \times \Theta.$$

On appelle alors **estimateur de type M de HUBER**, ou **M-estimateur de HUBER**, de  $\theta$  toute solution  $T_N = (A_N, B_N)$  du système en  $(\alpha, \beta)$  :

$$(14) \quad \begin{aligned} \sum_{n=1}^N \varphi\{(X_n - \alpha) / \beta\} &= 0, \\ \sum_{n=1}^N \psi\{(X_n - \alpha) / \beta\} &= 0, \end{aligned}$$

dans lequel le choix des fonctions  $\varphi$  et  $\psi$  dépend des propriétés recherchées pour  $T_N$  (notamment sa robustesse).

Un estimateur de HUBER particulier est tq :

$$(15) \quad \varphi(u) = \begin{cases} -c & \text{si } u \leq -c, \\ u & \text{si } -c \leq u \leq +c, \\ +c & \text{si } +c \leq u, \end{cases}$$

et :

$$(16) \quad \psi(u) = \{\varphi(u)\}^2 - E\varphi^2(U),$$

où  $E\varphi^2(U) = (2\pi)^{-1/2} \int \{\varphi(u)\}^2 \exp(-u^2/2) du$ .

(ix) Si  $X$  est un **échantillon iid** comme  $\xi$ , si  $P_{\theta}^{\xi}$  est une **loi symétrique** ( $\forall \theta$ ), si  $\varphi$  est impaire et  $\psi$  paire, on montre que l'estimateur de HUBER (ie la solution  $T_N = (A_N, B_N)$  de (14)) vérifie les propriétés suivantes :

(a)  $R_N (A_N - \alpha) \xrightarrow{\mathcal{L}}_{N \rightarrow +\infty} \mathcal{N}(0, \beta^2 \cdot \Sigma_A)$  (**loi gaussienne** centrée), avec  $\Sigma_A = E \varphi^2(\zeta) E \{\varphi'(\zeta)\}^2$ , en notant  $\zeta = (\xi - \alpha) / \beta$  ;

(b)  $R_N (B_N - \beta) \xrightarrow{\mathcal{L}}_{N \rightarrow +\infty} \mathcal{N}(0, \beta^2 \Sigma_B)$ , avec  $\Sigma_B = E \psi^2(\zeta) \cdot E (\zeta \cdot \psi'(\zeta))^2$  et  $\zeta = (\xi - \alpha) / \beta$  ;

(c)  $A_N$  et  $B_N$  sont asymptotiquement indépendants.