

## STATISTIQUE D'ORDRE (F6)

(17 / 04 / 2020, © Monfort, Dicostat2005, 2005-2020)

La notion de **statistique d'ordre** intervient souvent en **Statistique non paramétrique**, notamment en matière de **tests d'hypothèses**, lorsqu'il est utile de classer des grandeurs (**observations**) par valeurs croissantes (ou, alternativement, décroissantes).

(i) Soit  $(\Omega, \mathcal{F}, P)$  un **espace probabilisé** et  $X = (X_1, \dots, X_N) : \Omega \mapsto \mathbf{R}^N$  un **N-échantillon aléatoire** réel dont la **lp** est notée  $P^X = X(P)$ .

A chaque « réalisation »  $X(\omega) = x$  de  $X$  on peut associer une **permutation**  $\sigma_\omega \in \sigma_N$  (groupe des permutations de  $N_N^* = \{1, \dots, N\}$ ) et un **vecteur aléatoire**  $(Y_1, \dots, Y_N) : \Omega \mapsto \mathbf{R}_{\leq}^N$  tq, à la fois :

$$(1) \quad Y_n(\omega) \leq Y_{n+1}(\omega), \quad \forall n \in \{1, \dots, N-1\},$$

$$(2) \quad Y_n(\omega) = X_{\sigma_\omega(n)}(\omega), \quad \forall n \in \{1, \dots, N\}.$$

où  $\sigma_\omega(n)$  désigne  $\sigma_\omega(n)$ .

Autrement dit,  $Y$  est déterminé selon :

$$(3) \quad \begin{aligned} Y_1 &= \min \{X_n : n \in \{1, \dots, N\}\} = X_{\sigma_\omega(1)}, \\ Y_2 &= \min \{X_n : n \in \{1, \dots, N\} \setminus \{\sigma_\omega(1)\}\} = X_{\sigma_\omega(2)}, \\ &\dots \\ Y_\alpha &= \min \{X_n : n \in \{1, \dots, N\} \setminus \{\sigma_\omega(1), \dots, \sigma_\omega(\alpha-1)\}\} = X_{\sigma_\omega(\alpha)}, \\ &\dots \\ Y_N &= \max \{X_n : n \in \{1, \dots, N\}\} = X_{\sigma_\omega(N)}. \end{aligned}$$

On dit que  $Y$  est l'**échantillon ordonné**, ou l'**échantillon d'ordre**, associé à  $X$  et on le note  $Y = X^{(\cdot)} = (X^{(1)}, \dots, X^{(N)})$ , avec  $X^{(n)} \leq X^{(n+1)}$ ,  $\forall n \in \{1, \dots, N-1\}$  (ordre croissant).

On rencontre aussi les notations  $X_{(\cdot)}$ ,  $X^{(0)}$ ,  $X_{(\cdot)}$ ,  $X^{(\cdot)}$  ou encore  $X_{\cdot} = (X_{1:N}, \dots, X_{N:N})$ , avec  $X_{n:N} \leq X_{n+1:N}$ ,  $\forall n \in \{1, \dots, N-1\}$ .

L'application  $s : \mathbf{R}^N \mapsto \mathbf{R}_{\leq}^N$  définit donc une **statistique**  $S = s(X)$  selon :

$$(4) \quad S = s(X) = X^{(\cdot)} = Y.$$

C'est pourquoi  $X^{(\cdot)}$  est aussi appelé **statistique d'ordre**, ou **statistique ordonnée**, associée à  $X$ .

Toute statistique  $T = t(X^{(\cdot)})$  définie à partir de la statistique d'ordre  $X^{(\cdot)}$  est souvent elle-même appelée **statistique d'ordre**.

Enfin, on adopte parfois l'ordre inverse du précédent  $\leq$ .

(ii) Lorsque plusieurs coordonnées de  $X$  sont égales (**observation multiple**), on convient généralement d'identifier les permutations correspondantes (cf **rang multiple**).

La notion se généralise au cas où  $\mathbf{R}^N$  est remplacé par un ensemble  $\mathcal{X} = \mathcal{X}_0^N$  (puissance cartésienne d'un ensemble  $\mathcal{X}_0$  totalement ordonné) (cf **relation d'ordre**).

D'un point de vue terminologique, on dit que :

(a)  $Y_n = X_{(n)}$  est la  $n$ -ième valeur de l'échantillon  $X_{(\cdot)}$  ;

(b)  $Y_1 = X_{(1)} = \min_n X_n$  est la (ou une) valeur minimum de  $X$ , et  $Y_N = \max_n X_n$  la (ou une) valeur maximum de  $X$ . On dit encore que  $X_{(1)}$  est l'**extrême** (ou l'**extrémité**) inférieur(e) de  $X$  et  $X_{(N)}$  l'**extrême** (ou l'**extrémité**) supérieur(e) de  $X$ . Le couple  $(X_{(1)}, X_{(N)})$  est souvent appelé **couple des (valeurs) extrêmes** de l'échantillon  $X$  ; de même le  $L+M$ -uple  $(X_{(1)}, \dots, X_{(L)}, \dots, X_{(N-M+1)}, \dots, X_{(N)})$ , dans lequel  $1 \leq L + M \leq N$  est appelé **statistique des valeurs extrêmes** (cf aussi **valeur extrême**, **statistique des extrêmes**) ;

(c) la différence  $R_N = X_{(N)} - X_{(1)}$  est l'**étendue empirique** de  $X$  (ou des observations  $X_n$ ).

(iii) La statistique d'ordre possède les propriétés suivantes :

(a) si  $X$  est un **échantillon indépendant** (non nécessairement équidistribué), alors  $Y = X^{(\cdot)}$  n'est plus indépendant, car les inégalités  $Y_1 \leq \dots \leq Y_N$  contraignent le support de la loi  $P^Y$  de  $Y$  :  $\text{Supp } P^Y = \mathbf{R}_{\leq}^N$  (cf **support d'une probabilité**).

(b) si la loi  $P^X$  de  $X$  est une **loi absolument continue** pr à la **mesure de LEBESGUE**  $\lambda_N$  et si  $Y = X^{(\cdot)}$  est la statistique d'ordre de  $X$ , alors la **densité**  $g = dP^Y / d\lambda_N$  de  $Y$  s'exprime en fonction de la densité  $f = dP^X / d\lambda_N$  de  $X$  selon :

$$(5) \quad g(y) = \mathbf{1}_I(y) \cdot \sum_{\sigma \in \sigma(N)} f(\sigma(y)),$$

où l'on note  $I = \mathbf{R}_{<}^N$ ,  $\sigma(y) = (y_{\sigma(1)}, \dots, y_{\sigma(N)})$  et  $\sigma(N) = \sigma_N$  est le **groupe** des permutations de  $N_N^*$  déjà mentionné.

Si, de plus,  $X$  est un **échantillon iid** selon la loi  $P^\xi$  d'une **variable parente**  $\xi : \Omega \mapsto \mathbf{R}$  (ie si  $P^X = (P^\xi)^{\otimes N}$ ), la densité  $g$  de  $Y$  s'écrit :

$$(5)' \quad g(y) = N! \cdot \mathbf{1}_I(y) \cdot \prod_{n=1}^N f_0(y_n),$$

où  $I = \mathbf{R}_{<}^N$ ,  $f_0 = dP^\xi / d\lambda_1$  (densité de  $P^\xi$  pr à  $\lambda_1$ ) et où  $y = (y_1, \dots, y_N)$ .

Si  $F_0$  est la **fonction de répartition** associée à  $P^\xi$  et si l'on pose  $Z_n = F_0(Y_n)$  et  $Z = (Z_1, \dots, Z_N)$ , alors la densité de probabilité  $h$  de  $P^Z$  pr à  $\lambda_N$  ne dépend pas de  $P$  :

$$(6) \quad h(z) = N! \cdot \mathbf{1}_C,$$

où  $C = \{z \in ]0, 1[^N : 0 < z_1 < \dots < z_N < 1\}$ .

La propriété (6) est très utilisée en **Statistique non paramétrique**.

(iv) Si  $X$  est un échantillon iid selon une **loi discrète**  $P^\xi$  (de densité  $f_0$  pr à une **mesure de comptage**), il peut exister des valeurs multiples  $X_n$  dans  $X$ . On peut alors écrire :

$$(7) \quad \begin{aligned} P([X^{(1)} = \dots = X^{(n(1))} = x_1] \cap \dots \cap [X^{(n(1)+\dots+n(M-1)+1)} = \dots = X^{(n(1)+\dots+n(M))} = x_M]) \\ = (n_1! \dots n_M!)^{-1} \cdot n! \cdot \prod_{m=1}^M \{f_0(x_m)\}^{n(m)}, \end{aligned}$$

où l'on note  $n = \sum_{m=1}^M n_m$  et  $n(m)$  pour  $n_m$ .

(v) Si  $X$  est iid selon  $P^\xi$  (dont la fr est notée  $F_0$ ), alors la coordonnée  $X^{(n)}$  admet,  $\forall n \in \{1, \dots, N\}$ , pour fr propre (ie pour **fonction de répartition marginale**) :

$$(8) \quad G_{(n)}(x) = P(X^{(n)} \leq x) = n \cdot C_N^n \cdot \int_0^{F_0(x)} u^{n-1} (1-u)^{N-n} du.$$

Si, de plus,  $F_0$  est dérivable, de dérivée (densité)  $f_0$ , alors la lp de  $X^{(n)}$  admet pour densité :

$$(9) \quad g_{(n)}(x) = n \cdot C_N^n \cdot (F_0(x))^{n-1} \cdot \{1 - F_0(x)\}^{N-n} \cdot f(x).$$

La formule (8) supposait  $P^\xi$  absolument continue pr à  $\lambda_1$ . Dans le cas discret (absolue continuité « discrète »), on obtient :

$$(8)' \quad G_{(n)}(x) = \sum_{\alpha=n}^N C_N^\alpha (F_0(x))^\alpha \cdot \{1 - F_0(x)\}^{N-\alpha}.$$

(vi) Si  $R = (R_1, \dots, R_N)$  est la **statistique de rang** associée à  $X$ , on a (par définition) :

$$(10) \quad X^{(n)} = X_{\sigma(n)} \text{ (ou } = X_p) \Leftrightarrow n = R_{\sigma(n)} \text{ (ou } = R_p).$$

Autrement dit,  $X^{(R(n))} = X_n$ ,  $\forall n \in \{1, \dots, N\}$  (en notant  $R(n)$  pour désigner  $R_n$ ).

(vii) Une statistique d'ordre vérifie des inégalités classiques, eg :

(a) les **inégalités de G. STYAN - H. WOLKOWICZ**, ie,  $\forall n \in \{1, \dots, N\}$  :

$$(11) \quad \begin{aligned} X^{(1)} &\leq \bar{X}_N - (N-1)^{-1/2} S_N, \\ \bar{X}_N - \{(N-n)/n\}^{1/2} \cdot S_N &\leq X^{(n)} \leq \bar{X}_N + \{(N-1)/(N-n+1)\}^{1/2} \cdot S_N, \\ X^{(N)} &\geq \bar{X}_N + (N-1)^{-1/2} S_N, \end{aligned}$$

expressions dans lesquelles  $\bar{X}_N$  est la **moyenne empirique** de l'échantillon  $X$  et  $S_N^2 = N^{-1} \sum_{n=1}^N (X_n - \bar{X}_N)^2 = N^{-1} X' P X$  sa **variance** empirique (non corrigée) ;

(b) l'**inégalité de R.A. GROENEVELD**. Si  $X_n \geq 0, \forall n \in \{1, \dots, N\}$ , on a :

$$(12) \quad 0 \leq X^{(n)} \leq (N - n - 1)^{-1} N \cdot \bar{X}_N, \quad \forall n \in \{1, \dots, N\}.$$

(viii) En Statistique non paramétrique, on utilise souvent la statistique d'ordre  $X^{(.)}$  pour définir des statistiques linéaires de la forme (cf **statistique linéaire de rang**) :

$$(13) \quad L_N = N^{-1} \sum_{n=1}^N s(n / (N+1)) \cdot h(X^{(n)}),$$

dans laquelle la fonction  $s : [0, 1] \mapsto \mathbf{R}$  est appelée (fonction) « score » (cf **fonction score**) et la fonction  $h : \mathbf{R} \mapsto \mathbf{R}$  est une **fonction numérique** mesurable quelconque.