

## SONDAGE SYSTÉMATIQUE (M3, M6)

(16 / 03 / 2020, © Monfort, Dicostat2005, 2005-2020)

Un **sondage systématique** est un **sondage**, aléatoire ou non, effectué selon une « méthode systématique » : partant d'une première **unité statistique** (tirée au hasard ou non), on tire d'autres unités « équidistantes » de la première (dans l'espace ou dans le temps), ou équidistantes entre elles, jusqu'à atteindre une taille d'échantillon  $N$  donnée.

(i) En particulier, on appelle **sondage systématique** le sondage suivant. On étudie un ensemble fini  $\Omega = \{\omega_1, \dots, \omega_M\}$  (eg **population**), dont les éléments (**unités statistiques**) sont dotés d'un « **caractère** » observable  $\eta$ . Autrement dit, il existe une application  $\eta : \Omega \mapsto \mathcal{Y}$ , où  $(\mathcal{Y}, \mathcal{G})$  est un **espace d'observation** donné. On note  $Y = (Y_1, \dots, Y_M)$  les valeurs observables sur la population, avec  $Y_m = \eta(\omega_m)$  ( $m = 1, \dots, M$ ).

$N \in \mathbf{N}^*$  étant donné, soit alors  $F = N / M$  un **taux de sondage** tel que  $L = 1 / F = M / N$  soit un nombre entier (ie  $L \in \mathbf{N}^*$ ).

On appelle **sondage systématique** dans  $\Omega$  un sondage défini par le **plan de sondage** suivant :

(a) on prélève dans  $\Omega$  un élément  $\omega_l$  dont l'indice  $l \in \{1, \dots, L\}$ . Cet élément constitue le premier élément  $a_1$  de l'échantillon  $A$  à extraire de  $\Omega$  ;

(b) on prélève ensuite, systématiquement, les éléments  $a_n \in A$  selon la suite :

$$(1) \quad a_2 = \omega_{L+1}, \dots, a_N = \omega_{(N-1)L+1}.$$

La **partition** des indices relatifs aux unités de  $\Omega$  est la suivante :

$$\{1, \dots, L\}, \{L+1, \dots, 2L\}, \dots, \{(n-1)L + 1, \dots, nL\}, \dots, \{(N-1)L + 1, \dots, M\}.$$

Par suite, le sondage systématique est aléatoire ssi  $l$  est tiré au hasard uniforme dans  $N_L^* = \{1, \dots, L\}$ , ie ssi  $p(\lambda) = P([l = \lambda]) = L^{-1}, \forall \lambda \in N_L^*$ .

Dans ce contexte, on appelle parfois **table de C.H.D. BUYS-BALLOT** le tableau associé à la matrice suivante (cf aussi **composante d'une série temporelle, série temporelle**) :

$$(2) \quad Y = \begin{matrix} & Y_1 & Y_{L+1} & \dots & Y_{(n-1)L+1} & \dots & Y_{(N-1)L+1} \\ & \dots & \dots & \dots & \dots & \dots & \dots \\ Y_1 & Y_1 & Y_{L+1} & \dots & Y_{(n-1)L+1} & \dots & Y_{(N-1)L+1} \\ & \dots & \dots & \dots & \dots & \dots & \dots \\ Y_L & Y_L & Y_{2L} & \dots & Y_{nL} & \dots & Y_{NL} \end{matrix}$$

Le plan de sondage systématique revient donc à tirer, selon une probabilité uniforme, une ligne du tableau Y précédent. On appelle parfois L la **période du tirage**, ou **période du plan**.

(ii) La relation suivante :

$$(3) \quad Y_{ln} = Y_{(n-1)L+l}, \quad l = 1, \dots, L, \quad n = 1, \dots, N,$$

revient à remplacer l'indice double (l, n) par un indice simple (cf **indice, modèle multi-indicé**).

Par suite, l'**échantillon**  $y = (y_1, \dots, y_N)$  des valeurs observées sur l'échantillon A des unités tirées est défini selon  $y_n = \eta(a_n)$  ( $n = 1, \dots, N$ ). Cet échantillon est tq  $y_n = \eta(\omega_{(n-1)L+l})$  ( $\forall n$ ).

On note alors :

$$(4) \quad \begin{aligned} Y_{l.} &= N^{-1} \sum_{n=1}^N Y_{ln} && \text{(moyenne de } \eta \text{ dans la ligne } l), \\ Y_{..} &= M^{-1} \sum_{l=1}^L \sum_{n=1}^N Y_{ln} && \text{(moyenne générale de } \eta \text{ dans le (tableau } Y)). \end{aligned}$$

Le sondage est aléatoire si l est tiré au hasard (uniforme) dans  $N_L^*$ , ie si  $l \sim \mathcal{U}(N_L^*)$  (**loi uniforme discrète**). Par suite, la moyenne  $Y_{l.}$  est une **variable aléatoire** tq :

$$(5) \quad P(Y_{l.} = Y_{k.}) = L^{-1}, \quad \forall k = 1, \dots, L,$$

où  $Y_{k.}$  est une ligne donnée de Y.

De plus :

$$(6) \quad \begin{aligned} E Y_{l.} &= Y_{..} \text{ (estimateur sans biais),} \\ V Y_{l.} &= (N^2 L)^{-1} \cdot \sum_{l=1}^L \left\{ \sum_{n=1}^N (Y_{ln} - Y_{..}) \right\}^2. \end{aligned} \quad \forall l = 1, \dots, L$$

(iii) On établit alors que :

(a) le sondage systématique est aussi efficace qu'un **sondage élémentaire** (bernoullien ou exhaustif) si la **covariance** (théorique) entre deux colonnes du tableau Y précédent est nulle (ie lorsque l'ordre selon lequel les unités sont rangées dans Y est purement aléatoire) ;

(b) par contre, le sondage systématique est moins efficace si cette même covariance est positive. Un cas limite est celui où le caractère  $\eta$  présente des variations périodiques entre les unités, la période étant égale à L : l'unité  $\omega_l$  est identique à l'unité  $\omega_{l+L}$  pour le caractère  $\eta$  considéré ;

(c) lorsque la covariance précédente est négative, le sondage systématique est plus efficace que le sondage élémentaire, les unités « voisines » (pour le

caractère  $\eta$ ) ayant tendance à se ressembler. C'est dans cette situation que ce type de sondage est utilisé.