

## LACUNE (C8, G9)

(24 / 03 / 2020, © Monfort, Dicostat2005, 2005-2020)

D'un point de vue conceptuel et terminologique, le terme de **lacune** peut se référer à l'une des deux **situations statistiques** suivantes :

(a) **variable omise**, ou **variable absente** : lorsqu'une **variable** est absente et ne permet pas de modéliser complètement un **phénomène** (ie de définir une **représentation statistique**), on parle d'**omission** ou de **variable omise**. Dans ce cas, c'est la variable qui est manquante lors de la **spécification** du modèle décrivant le phénomène ;

(b) **variable inobservable** : lorsqu'une **variable** n'est pas observée sur une **unité statistique** particulière (cf **observabilité**, **inobservabilité**), on parle plutôt de **lacune** ou d'**observation manquante**. Dans ce cas, c'est l'**observation** de la **variable aléatoire** sur l'unité considérée qui est manquante (ie non apparente ou indisponible). Cette situation peut conduire à l'impossibilité de calcul d'une **statistique** dépendant de cette observation.

C'est le second sens qui est le plus usuel.

(i) Une **lacune**, ou **observation manquante**, peut se rencontrer lors de la mise en oeuvre d'une source d'**information statistique**.

Elle peut résulter d'une absence d'information, partielle ou totale, relative à des **variables** (ou « descripteurs ») décrivant un **phénomène** donné. Autrement dit, le **statisticien** cherche à « mesurer » les « valeurs » de ces variables sur des unités statistiques, mais il ne peut pas observer certaines valeurs sur certaines unités (cf **observabilité**, **inobservable**).

Les raisons principales sont les suivantes :

(a) impossibilité pratique d'observation de certaines variables, ou seulement de certaines données : **système statistique** peu évolué (ne couvrant pas un champ complet souhaitable), instrument de **mesure** imprécis, etc ;

(b) **troncature** de la loi générant ces données ;

(c) **censure** des données, eg :

(c)<sub>1</sub> **unité de sondage** d'un **échantillon** obtenu par **sondage** : perte matérielle (« sondage destructif »), **non réponse** (complète, ou seulement partielle) des unités ;

(c)<sub>2</sub> **unité expérimentale** utilisée pendant une **expérience aléatoire**. Sa **durée de vie** peut se trouver raccourcie pour des raisons expérimentales : expériences avec destruction (au sens physique ou biologique du terme).

(ii) De façon générale, soit  $(\Omega, \mathcal{F}, P)$  un **espace probabilisé** et  $\xi : \Omega \mapsto \mathcal{X}$  une va  $\xi$ , où  $(\mathcal{X}, \mathcal{B})$  désigne l'**espace mesurable** des valeurs possible de  $\xi$  (ie un **espace**

**d'observation**). On note  $\mathcal{L}(\xi)$  la **loi de probabilité** de  $\xi$  et  $X = (X_1, \dots, X_N)$  un **N-échantillon iid** issu de  $\xi$ . L'existence de lacunes signifie que l'on n'observe pas toutes les N « réalisations »  $X_n$  ( $n = 1, \dots, N$ ) de  $\xi$  :

(a) certaines, en nombre  $N^o$ , sont effectivement observées et l'on connaît chaque indice  $n$  tq  $X_n$  est observé. On note alors  $X^o$  la suite de ces réalisations observables ;

(b) les  $N^a = N - N^o$  autres réalisations ne sont pas observées : observations « absentes », ou « inobservations » (cf **inobservable**). On connaît chaque indice  $n$  tq  $X_n$  n'est pas observé. On note alors  $X^a$  la suite de ces réalisations inobservables.

On peut alors structurer  $X$  selon  $X = (X^o, X^a)$ , où  $X^o$  désigne la liste des observations effectives et  $X^a$  celle des lacunes.

Cette **situation statistique** peut interpréter comme suit :

(a) la loi  $\mathcal{L}(\xi)$  de  $\xi$  permet, en principe, de générer une observation quelconque tq  $\xi(\omega) = x$  (eg par **expérimentation** ou par **sondage**) ;

(b) mais il existe une variable dichotomique (**variable binaire**)  $\alpha$  à valeurs dans  $\{0, 1\}$  qui est (implicitement) conjointement tirée selon une loi  $\mathcal{L}(\alpha) = p \cdot \mathbf{1}(\alpha = 1) + (1 - p) \cdot \mathbf{1}(\alpha = 0)$  ou encore selon  $\mathcal{L}(\alpha) = p \cdot \delta_1 + (1 - p) \cdot \delta_0$ , où  $p \in [0, 1]$  désigne la probabilité d'observation de  $\xi$ ,  $1 - p$  celle d'inobservation (lacune) et  $\delta_x$  la **masse de DIRAC** placée en  $x \in \mathcal{X}$ .

Autrement dit,  $\alpha$  « occulte »  $\xi$  au sens où, au lieu d'observer  $\xi$  selon la loi  $\mathcal{L}(\xi)$  elle-même, on observe la valeur de  $\xi$  générée par  $\mathcal{L}(\xi)$  seulement lorsque  $\alpha = 1$ , ie générée par la **loi conditionnelle** de  $\xi$  sachant  $\alpha$ . Cette loi se déduit du **théorème des probabilités composées**, ie (sous forme « symbolique ») :

$$(1) \quad \mathcal{L}(\xi / \alpha) = \mathcal{L}(\xi \cap \alpha) / \mathcal{L}(\alpha),$$

ie, de façon explicite,  $\forall B \in \mathcal{B}$  et  $\forall A \in \mathcal{B}$  :

$$(2) \quad \mathcal{L}([\xi \in B] / [\alpha \in A]) = \mathcal{L}([\xi \in B] \cap [\alpha \in A]) / \mathcal{L}([\alpha \in A]),$$

$\mathcal{L}(\xi)$  est donc la « marge » de la **loi conjointe**  $\mathcal{L}(\xi, \alpha)$  ou  $\mathcal{L}(\xi \cap \alpha)$  (cf **loi marginale, censure**).

D'après (1) :

$$(3) \quad \mathcal{L}(\xi \cap \alpha) = \mathcal{L}(\alpha) \cdot \mathcal{L}(\xi / \alpha),$$

expression dans laquelle :

(a)  $\mathcal{L}(\alpha)$  est connue ;

(b) les réalisations issues de  $\mathcal{L}(\xi / \alpha)$  sont observables, et leur nombre  $N^o$  est connu.

Par suite, on peut estimer :

(a)  $p$  par la **proportion**  $p^{\sim} = N^o / N$ , ce qui conduit à l'estimateur  $\mathcal{L}^{\sim}(\alpha) = p^{\sim} \cdot \delta_1 + (1 - p^{\sim}) \cdot \delta_0$  de  $\mathcal{L}(\alpha)$  ;

(b)  $\mathcal{L}(\xi / \alpha)$  par la **loi empirique** issue de  $X^o$ , et notée  $P_{N^o}$ .

(c) donc  $\mathcal{L}(\xi \cap \alpha)$  par :

(4)  $\mathcal{L}^{\sim}(\xi \cap \alpha) = \mathcal{L}^{\sim}(\alpha) \cdot P_{N^o}$ .

Un estimateur  $\mathcal{L}^{\sim}(\alpha)$  de  $\mathcal{L}(\xi)$  se déduit de (4) par **marginalisation** relativement à  $\alpha$ .

(iii) La description précédente peut s'étendre à plusieurs types de variables (eg endogènes et exogènes) prises en compte dans un modèle.

Soit  $(\Omega, \mathcal{F}, \mathcal{P})$  un **modèle statistique**,  $\zeta = (\xi, \eta) : \Omega \mapsto (\mathcal{X}, \mathcal{Y})$  un **couple aléatoire**, où  $(\mathcal{X}, \mathcal{B})$  désigne l'**espace mesurable** des valeurs possible de  $\xi$  (ie un **espace d'observation** ou un **espace d'états**) et  $(\mathcal{Y}, \mathcal{C})$  celui des valeurs possibles de  $\eta$ . On note  $\mathcal{L}(\xi, \eta)$  la **loi de probabilité** de  $\zeta$  et  $Z = (Z_1, \dots, Z_N)$  un N-échantillon iid issu de  $\zeta$ . On désigne alors par  $X = (X_1, \dots, X_N)$  l'échantillon partiel correspondant à  $\xi$  et par  $Y = (Y_1, \dots, Y_N)$  l'échantillon partiel correspondant à  $\eta$ .

Dans ce contexte, l'existence de lacunes signifie que l'on ne peut observer les N « réalisations »  $(X_n, Y_n)$  ( $n = 1, \dots, N$ ) de  $\zeta$ . Comme ces coordonnées sont associées (réalisations relatives à la même **unité statistique**  $n$ ), on note :

(a)  $N^{oo}$  le nombre de coordonnées tq  $(\xi, \eta)$  est observé selon  $(X_n^o, Y_n^o)$ . On connaît les indices  $n$  correspondants. On note alors  $(X^o, Y^o)$  la suite des réalisations observables simultanément ;

(b)  $N^{oa}$  le nombre de coordonnées tq  $\xi$  seul est observé. On connaît les indices  $n$  correspondants. On note  $(X^o, Y^a)$  la suite des réalisations de ce type ;

(c)  $N^{ao}$  le nombre de coordonnées tq  $\eta$  seul est observé. On connaît les indices  $n$  correspondants. On note  $(X^a, Y^o)$  la suite des réalisations de ce type ;

(d)  $N^{aa} = N - N^{oo} - N^{oa} - N^{ao}$  le nombre de coordonnées dont les observations sont absentes (« inobservations »). On connaît les indice  $n$  tq  $(\xi, \eta)$  n'est pas observé. On note  $(X^a, Y^a)$  la suite des réalisations inobservables simultanément.

On peut alors structurer  $Z$  selon  $Z = \{(X^o, Y^o), (X^o, Y^a), (X^a, Y^o), (X^a, Y^o)\}$ . L'interprétation est similaire à la précédente :

(a) la loi  $\mathcal{L}(\xi, \eta)$  de  $(\xi, \eta)$  doit normalement donner lieu à des observations  $(\xi(\omega), \eta(\omega)) = (x, y)$  de  $(\xi, \eta)$  ;

(b) le couple de variables dichotomiques  $(\alpha, \beta)$  à valeurs dans  $\{0, 1\}^2$  est conjointement tiré,  $\alpha$  selon une loi  $\mathcal{L}(\alpha / (\xi, \eta)) = p \cdot \delta_1 + (1 - p) \cdot \delta_0$ , où  $p \in [0, 1]$ , et  $\beta$  selon une loi  $\mathcal{L}(\beta / (\xi, \eta)) = r \cdot \delta_1 + (1 - r) \cdot \delta_0$ , où  $r \in [0, 1]$ .

Autrement dit,  $(\alpha, \beta)$  « occulte »  $(\xi, \eta)$ , au sens où, au lieu d'observer  $(\xi, \eta)$  selon la loi  $\mathcal{L}(\xi, \eta)$  elle-même, on observe seulement la valeur de  $(\xi, \eta)$  générée par la loi conditionnelle :

$$(5) \quad \mathcal{L}((\xi, \eta) / (\alpha, \beta)) = \mathcal{L}((\xi, \eta) \cap (\alpha, \beta)) / \mathcal{L}(\alpha, \beta),$$

La loi  $\mathcal{L}(\xi, \eta)$  est la loi marginale de la loi jointe  $\mathcal{L}((\xi, \eta) \cap (\alpha, \beta))$ .

D'après (5) :

$$(6) \quad \mathcal{L}((\xi, \eta) \cap (\alpha, \beta)) = \mathcal{L}(\alpha, \beta) \cdot \mathcal{L}((\xi, \eta) / (\alpha, \beta)),$$

expression dans laquelle :

(a)  $\mathcal{L}(\alpha, \beta)$  est connue, du moins si l'on admet l'hypothèse selon laquelle  $\alpha$  et  $\beta$  sont indépendantes entre elles (l'occultation agit séparément sur  $\xi$  et sur  $\eta$ ) :

$$(7) \quad \mathcal{L}(\alpha, \beta) = \{p \cdot \delta_1 + (1 - p) \cdot \delta_0\} \otimes \{r \cdot \delta_1 + (1 - r) \cdot \delta_0\}.$$

(b) les réalisations issues de  $\mathcal{L}((\xi, \eta) / (\alpha, \beta))$  sont observables, ainsi que les nombres  $N^{oo}$ ,  $N^{ao}$ ,  $N^{oa}$ , et  $N^{aa}$ .

(iv) En présence de lacunes, les **procédures statistiques** usuelles (estimation des paramètres, tests, prévisions, classifications, etc) doivent donc faire l'objet de **modification** ou d'adaptation.

Il en est de même lorsque les échantillons considérés ne sont pas iid.

(v) Ainsi, dans le cas d'un **problème d'estimation**, on étudie le **modèle statistique**  $(\mathcal{X}, \mathcal{B}, (P_\theta^X)_{\theta \in \Theta})$ , dans lequel  $X$  représente un **échantillon** (ou une **statistique** d'intérêt) et  $P_\theta^X$  une **loi de probabilité** possible de  $X$ , dominée par une **mesure positive**  $\mu$  (loi à **densité**). On note alors :

$$(8) \quad (dP_\theta^X / d\mu)(x) = L(x, \theta)$$

la **vraisemblance** associée au modèle.

Par suite, si  $X_p$  (resp  $x_p$ ) désigne l'ensemble des valeurs « présentes » (ie observées) de  $X$  (resp dans la densité  $\mathcal{L}$ ) et  $X_a$  (resp  $x_a$ ) l'ensemble des valeurs « absentes » (ie non observées), plusieurs méthodes d'estimation de  $\theta$  sont possibles.

(a) **méthode de « faux » maximum de vraisemblance (R.L. ANDERSON - M.S. BARTLETT)** dans laquelle le « **paramètre** » considéré, à estimer, est le couple  $(x_a, \theta)$ . Il s'agit alors de résoudre le **problème d'optimisation** suivant :

$$(9) \quad \sup_{(x_a, \theta)} \mathcal{L}(x_p, x_a, \theta),$$

(où  $x_a$  désigne  $x_a$ ).

Il ne s'agit pas d'un « vrai » **maximum de vraisemblance** car, en général, le « paramètre partiel »  $x_a$ , parfois appelé **paramètre incident**, varie avec  $N$  (taille de l'échantillon  $X$ ) ;

(b) **méthode de maximisation de la vraisemblance partiellement intégrée (M.H. de GROOT - K. GOEL)**. Si l'on note  $\mathcal{X} = \mathcal{X}_p \times \mathcal{X}_a$  (ou, selon le cas,  $\mathcal{X} = \mathcal{X}_p \cup \mathcal{X}_a$ ), l'ensemble des valeurs de  $X$ , représenté selon la présence et l'absence, ie si  $x = (x_p, x_a)$ , et si  $\mu / \mathcal{X}_a$  est la **restriction** de  $\mu$  à l'ensemble des valeurs absentes (on note  $\mathcal{X}_a$  pour désigner  $\mathcal{X}_a$ ), la méthode consiste à maximiser pr à  $\theta$  la « vraisemblance » (partiellement) intégrée :

$$(10) \quad \mathcal{L}(x_p, \theta) = \int L(x_p, x_a, \theta) d\mu / \mathcal{X}_a(x_a) ;$$

(c) **méthode bayésienne (S.J. PRESS - A.J. SCOTT)** consistant à définir une mesure de **probabilité a priori**  $\Pi_{(x_a, \theta)}$  relative au couple  $(x_a, \theta)$ , auquel on fait jouer le rôle de paramètre ;

(d) **méthode en trois temps** :

(d)<sub>1</sub> on considère d'abord un modèle « restreint » dans lequel l'échantillon est  $X_p$  et le paramètre  $\theta$ . On estime alors provisoirement  $\theta$  à partir de la vraisemblance associée à ce modèle (ie à partir de la seule donnée informative observable  $x_p$ ). Cette étape n'est pas toujours réalisable ;

(d)<sub>2</sub> on revient au modèle « complet » et l'on estime  $x_a$  à l'aide d'une **caractéristique de centralité** (eg la **moyenne**) de la **loi** de ce modèle (caractéristique déjà estimée) ;

(d)<sub>3</sub> on applique la **méthode du maximum de vraisemblance** avec la vraisemblance complétée  $\mathcal{L}(x_p, x_a, \theta)$ , ce qui fournit l'estimateur recherché de  $\theta$ .

(vi) Ainsi, l'analyse statistique des lacunes revient à considérer qu'il n'est possible d'observer que dans certaines parties de  $\mathcal{X}$  (cf aussi **censure**). Autrement dit, si  $s : \mathcal{X} \mapsto \mathcal{Y}$  est une **statistique** surjective (cf **application surjective**), on observe seulement les données incomplètes  $Y = s(X) : \Omega \mapsto \mathcal{Y}$  au lieu d'observer complètement  $X : \Omega \mapsto \mathcal{X}$ .

Comme dans le cas d'une **aberration**, la notion de lacune peut donc aussi être présentée à partir de celle de **mélange de lois**. La situation est cependant plus favorable, car on connaît généralement les **indices** (ou « identifiants ») aussi bien que les observations des unités observées, mais aussi les indices (ou identifiants) des unités non observées.

Cette approche permet alors de définir directement la loi de  $\xi$  à l'aide des lois composant ce mélange :

(a) en effet, si  $\xi$  est la **variable aléatoire** considérée, si  $\xi_p$  est une variable **observable** (ie dont les observations sont toutes présentes) et si  $\xi_u$  est une variable **inobservable** (ie dont les observations sont toutes absentes) (« *unobservable* »), on peut définir la loi de  $\xi$  comme **combinaison linéaire convexe** de ces dernières, ie :

$$(11) \quad P^\xi = \alpha_o \mathcal{L}(\xi_o) + \alpha_u \mathcal{L}(\xi_u),$$

où  $(\alpha_o, \alpha_u) \in S_2$  (**simplexe** de  $\mathbf{R}^2$ ),  $P^\xi$  désigne la loi de  $\xi$  et  $\mathcal{L}(v)$  la loi de  $v$  (avec  $v = \xi_o$  ou  $\xi_u$ ) ;

(b) disposant d'observations  $X_n$  de  $\xi$ , avec  $n = 1, \dots, N$ , et en notant  $X = (X_1, \dots, X_N)$  l'**échantillon**, on peut représenter la **loi empirique**  $P_N$  ou  $\mathcal{L}_N(\xi)$  de  $\xi$  selon une équation de la forme :

$$(12) \quad \mathcal{L}_N(\xi) = \alpha_o \cdot \mathcal{L}_N(\xi_o) + \alpha_u \cdot \mathcal{L}_N(\xi_u).$$

Cette équation n'est cependant pas une **équation estimante** car elle ne permet pas, en elle-même, d'estimer chaque « **paramètre** »  $\alpha$ .

Lorsque  $N_o$  observations sont présentes (observables) et  $N_u$  absentes (unobservables), avec  $N_o + N_u = N$ , et que  $X$  est organisé de façon conforme, selon  $(X^o, X^u)$ , les données du problème sont  $(N_o, X^o, N_u)$ .

En général, on peut estimer resp  $\alpha_o$  et  $\alpha_u$  à l'aide des **proportions** :

$$(13) \quad \begin{aligned} \alpha_o \tilde{\phantom{\alpha}} &= N_o / N, \\ \alpha_u \tilde{\phantom{\alpha}} &= N_u / N. \end{aligned}$$

D'autre part,  $\mathcal{L}(\xi_o)$  peut être estimée à l'aide de la fre  $\mathcal{L}_{N_o}^o(\xi_o)$  définie à partir du sous-échantillon  $X^o$  de  $X$  ( $N_o$  désignant, par commodité,  $N_o$ ).