

MODÈLE DES CLASSES LATENTES (B4, K12, O24)

(25 / 04 / 2020, © Monfort, Dicostat2005, 2005-2020)

Le **modèle des classes latentes** est un modèle d'analyse de **structures latentes**, dont l'objet est la décomposition d'une **probabilité** (cf aussi **modèle de structures latentes, loi multivariée**).

(i) Soit $(\Omega, \mathcal{F}, \mathcal{P})$ un **modèle statistique** de base. On considère, $\forall h \in N_H^*$, un « **caractère** », ou une **variable qualitative**, $\kappa_h : \Omega \mapsto \mathcal{K}_h$ comportant un nombre fini H de modalités.

\mathcal{K}_h désigne ainsi l'**ensemble** fini des valeurs de κ_h , ie :

$$(1) \quad \mathcal{K}_h = \{k_{h,1}, \dots, k_{h,M(h)}\}, \quad \forall h \in N_H^*.$$

où l'on note, par commodité, $M(h)$ au lieu de M_h le nombre de modalités de κ_h .

Chaque variable κ_h est un « **attribut polytomique** » observable (cf **observabilité**), appelé **variable manifeste** du modèle (cf **tableau polytomique**).

On définit, $\forall P \in \mathcal{P}$, les **probabilités** :

$$(2) \quad p_l = P([(\kappa_1, \dots, \kappa_H) = (k_{1,m(1)}, \dots, k_{k,m(H)})]),$$

dans lesquelles :

(a) les **indices** sont notés $l = (m_1, \dots, m_H) \in \mathcal{I}$, $\mathcal{I} = \prod_{h=1}^H \mathcal{I}_h$, $\mathcal{I}_h = \{1, \dots, M_h\}$ (en notant $m(h)$ au lieu de m_h l'indice d'une modalité courante) ;

(b) les valeurs d'un **échantillon** sont notées $(K_1, \dots, K_H) : \Omega \mapsto \prod_{h=1}^H \mathcal{K}_h$, où $(k_{1,m(1)}, \dots, k_{H,m(H)})$ est un H -uple parmi tous les H -uples possibles.

Chaque probabilité p_l est celle d'une valeur correspondant à une « case » $(k_{1,m(1)}, \dots, k_{H,m(H)})$ (cf **tableau de contingence**).

(ii) On suppose qu'il existe un **caractère qualitatif** à L modalités, appelées **classes latentes**, ou **classes sous-jacentes**, $\mathcal{Z} = \{z_1, \dots, z_L\}$ définissant une va qualitative $\zeta : \Omega \mapsto \mathcal{Z}$, appelée **variable latente**.

Le **modèle de(s) classes latentes** (de P.F. LAZARSELD) est défini comme un modèle tq le précédent, dans lequel, $\forall P \in \mathcal{P}$, on suppose que les p_l admettent une décomposition de la forme :

$$(3) \quad p_l = \sum_{\lambda=1}^L p_l(\lambda),$$

avec :

$$(4) \quad p_l(\lambda) = q_\lambda \cdot \{\prod_{h=1}^H r_{m(h),\lambda}\}, \quad \forall \lambda \in N_L^*,$$

expression dans laquelle :

$$(5) \quad \begin{aligned} q_\lambda &= P(\zeta = z_\lambda), \quad \forall \lambda \in N_L^*, \\ r_{m(h),\lambda} &= P(K_h = k_{h,m(h)} / \zeta = z_\lambda), \quad \forall m_h \in \mathcal{J}_h, \forall h \in N_H^*, \forall \lambda \in N_L^*. \end{aligned}$$

Autrement dit :

(a) l'équation (3) exprime que les probabilités p_l peuvent se décomposer selon L **catégories** mutuellement exclusives et exhaustives (cf **partition**, **théorème de BAYES**, **théorème des probabilités composées**) ;

(b) l'équation (4) exprime que, dans chaque catégorie $\lambda \in N_L^*$ (donc pour chaque classe latente z_λ), les variables K_h sont indépendantes (cf **indépendance**).

(iii) Les $L^2 \cdot (\prod_{h=1}^H M_h)$ **paramètres** du modèle $\{(3),(4)\}$ peuvent être estimés de plusieurs façons (eg par la **méthode du maximum de vraisemblance**), à l'aide du **N-échantillon** $K = (K_1, \dots, K_H)$ observé pour les variables manifestes κ (les « **observables** »).

Ces paramètres vérifient donc les **contraintes** suivantes :

$$(6) \quad \begin{aligned} \sum_{\lambda=1}^L q_\lambda &= 1, \\ \sum_{m(h)=1}^{M(h)} r_{m(h),\lambda} &= 1, \end{aligned} \quad \forall (h, \lambda) \in N_H^* \times N_L^*.$$

En pratique, les **informations** qui permettent une telle **estimation** sont souvent disposées (ou disponibles) (cf **disposition**) sous la forme :

(a) d'un **tableau de contingence** multidimensionnel (empirique), noté eg $N = (n_l)_{l \in \mathcal{J}}$, dans lequel n_l désigne le nombre d'observations (ou d'**unités statistiques** : individus, etc) élémentaires ayant la **modalité multiple** $(k_{1,m(1)}, \dots, k_{H,m(H)})$;

(b) ou d'un tableau de **fréquences relatives**, défini à partir du précédent :

$$(7) \quad f_l = n_l / \sum_{j \in \mathcal{J}} n_j, \quad \forall l \in \mathcal{J}.$$