

MÉTHODE DE LAMBERT-WILLIAMS (K9)

(16 / 12 / 2019, © Monfort, Dicostat2005, 2005-2019)

La méthode de LAMBERT - WILLIAMS est une méthode de **classification automatique** descendante, qui procède par **dichotomies successives** de l'ensemble d'**unités** à analyser (**population, échantillon**).

(i) Soit $X \in M_{NK} (\{0, 1\})$ une **matrice** constituée de N **observations** (en lignes) relatives à H **variables qualitatives** η_h (en colonnes), la h-ième variable η_h possédant les modalités $\{1, \dots, M_h\}$, $\forall h \in N_H^*$, avec $\sum_{h=1}^H M_h = K$ (nombre total de colonnes).

Autrement dit, X est un (N,K)-**tableau statistique** mis sous **forme disjonctive complète**, son élément général x_{nk} étant défini selon :

$$(1) \quad x_{nk} = \begin{cases} 1 & \text{si l'observation d'indice } n \text{ de la variable } n^\circ h \text{ possède la modalité } m_h, \\ 0 & \text{sinon,} \end{cases}$$

avec $k = M_1 + \dots + M_{h-1} + m_h$.

La matrice X peut être représentée par **blocs** disposés en lignes, le h-ième bloc X (h) étant relatif à la variable η_h :

$$(2) \quad X = [X(1) \text{ /// } X(H)], \quad \text{avec } X(h) \in M_{N, M(h)} (\{0, 1\}),$$

où M(h) désigne M_h et /// des changements de lignes.

Par suite, $\forall (g, h)$, la matrice $C(g, h) = X(g)' X(h)$ est un **tableau de contingence** de dimensions (M_g, M_h) dont le terme général $c_{m(g), m(h)}(g, h)$ est le nombre d'observations n possédant à la fois la modalité m_g de g et la modalité m_h de h. Il en est de même de la matrice d'ensemble $X' X$, dont les blocs sont les $C(g, h)$, $\forall (g, h)$.

(ii) La **méthode de J.M. LAMBERT - W.T. WILLIAMS** procède par étapes :

(a) calcul, $\forall (g, h)$, des **statistiques du chi-deux** empiriques :

$$(3) \quad \mathcal{X}_{gh}^2 = c_{\cdot, \cdot}(g, h) \cdot \{(\sum_{m(g)=1}^{M(g)} \sum_{m(h)=1}^{M(h)} R_{gh}^2) - 1\},$$

dans lesquelles $R_{gh}^2 = (C_{g, \cdot} \cdot C_{\cdot, h})^{-1} \cdot C_{gh}^2$, avec $C_{gh} = c_{m(g), m(h)}(g, h)$, $C_{g, \cdot} = c_{m(g), \cdot}(g, h)$ et $C_{\cdot, h} = c_{\cdot, m(h)}(g, h)$;

(b) calcul des **statistiques** :

$$(4) \quad I(g) = \sum_{h \neq g} \mathcal{X}_{gh}^2, \quad \forall g.$$

On retient la variable η_0 d'indice h_0 tq : $I(h_0) = \max_g I(g)$ (**distance** maximale) ;

(c) partitionnement en deux classes de l'ensemble des modalités $\{1, \dots, M_{h(0)}\}$ de la variable η_0 (où $h(0)$ désigne h_0). On peut, eg :

(c)₁ soit considérer l'ensemble des **partitions** en deux classes de cet ensemble. Il existe $\exp_2 (M_{h(0)} - 1) - 1$ partitions de ce type (avec $\exp_2 x = 2^x$) ;

(c)₂ soit ne considérer que les partitions en deux classes consécutives de cet ensemble, en conservant l'ordre $\{1, \dots, M_{h(0)}\}$ des modalités de η_0 (cas d'une **variable ordinale** η_0). Il n'existe alors que $M_{h(0)} - 1$ partitions en deux classes ;

En notant $\{\mathcal{J}_{h(0)'}, \mathcal{J}_{h(0)''}\}$ une partition quelconque du type précédent, on choisit celle pour laquelle $I(h_0)$ est maximum ;

(d) en notant $\{(\mathcal{J}_{h(0)'})^{\sim}, (\mathcal{J}_{h(0)''})^{\sim}\}$ cette dernière partition, on assimile ses deux classes à deux modalités d'une nouvelle variable qualitative (qui remplace alors η_0). On définit alors deux (N,1)-sous-matrices de X (h_0) : dans la première (resp la seconde), on réunit les observations pour lesquelles la variable prend la modalité $(\mathcal{J}_{h(0)'})^{\sim}$ (resp $(\mathcal{J}_{h(0)''})^{\sim}$) ;

(e) on itère le procédé pour chacun des tableaux considérés. Chaque itération aboutit ainsi à partitionner progressivement l'ensemble E constitué des N observations.