

REVUE DE STATISTIQUE APPLIQUÉE

J. J. DENIMAL

Aides à l'interprétation mutuelle de deux hiérarchies construites sur les lignes et les colonnes d'un tableau de contingence

Revue de statistique appliquée, tome 45, n° 4 (1997), p. 93-110

http://www.numdam.org/item?id=RSA_1997__45_4_93_0

© Société française de statistique, 1997, tous droits réservés.

L'accès aux archives de la revue « *Revue de statistique appliquée* » (<http://www.sfds.asso.fr/publicat/rsa.htm>) implique l'accord avec les conditions générales d'utilisation (<http://www.numdam.org/conditions>). Toute utilisation commerciale ou impression systématique est constitutive d'une infraction pénale. Toute copie ou impression de ce fichier doit contenir la présente mention de copyright.

NUMDAM

Article numérisé dans le cadre du programme
Numérisation de documents anciens mathématiques

<http://www.numdam.org/>

AIDES À L'INTERPRÉTATION MUTUELLE DE DEUX HIÉRARCHIES CONSTRUITES SUR LES LIGNES ET LES COLONNES D'UN TABLEAU DE CONTINGENCE

J.J. Denimal

*Laboratoire de Statistique et Probabilités,
Université des Sciences et Technologies de Lille
59655 Villeneuve d'Ascq cedex, France*

RÉSUMÉ

A partir des deux hiérarchies construites à l'aide du critère de Ward sur les lignes et les colonnes d'un tableau de contingence, on se propose de définir des indices permettant d'interpréter chaque hiérarchie en fonction de l'autre. Un logiciel a été rédigé et est appliqué à un tableau de données.

Mots-clés : Critère de Ward, classification hiérarchique, tableau de contingence.

ABSTRACT

Using the two hierarchical clusterings of the rows and the columns of a contingency table, we propose to define new interpretative aids in cluster analysis. An application of the method to a set of data is also given.

Keywords : Ward's criterion, hierarchical clustering, contingency table.

1. Introduction

A partir d'un tableau de contingence k_{IJ} croisant deux ensembles I et J , on considère les 2 hiérarchies H_I et H_J construites respectivement sur I et sur J à l'aide du critère de Ward. A ces deux hiérarchies, seront associées deux bases orthogonales dans leurs espaces respectifs, chaque vecteur de base représentant un nœud de la hiérarchie correspondante. L'utilisation de ces bases pour décomposer les distances séparant le centre de gravité de chaque nœud du centre de gravité général ou séparant l'aîné et le benjamin d'un nœud permettra d'établir des indices qui auront pour objet d'interpréter chaque hiérarchie en fonction de l'autre. Nous verrons, d'autre part, que les hiérarchies H_I et H_J , bien que construites indépendamment l'une de l'autre, peuvent recevoir une interprétation qui se placerait dans le cadre d'une classification hiérarchique croisée. Nous montrerons, également, que l'inertie du tableau k_{IJ} peut se décomposer en une somme de termes doublement indexés par les nœuds des deux hiérarchies. Cette formule sera utilisée pour calculer les inerties de sous-tableaux

emboîtés de k_{IJ} , ces derniers étant définis à partir des nœuds des deux hiérarchies. L'étude de ces inerties pourra permettre une détermination simultanée des coupures de H_I et H_J . Enfin, un logiciel utilisant les résultats de cet article a été écrit. On trouvera, au dernier paragraphe, une application de celui-ci à un tableau de données.

2. Notations

k_{IJ} étant le tableau de contingence analysé, on posera comme à l'habitude :

$$\forall i \in I, k(i) = \sum_{j \in J} k(i, j); \quad \forall j \in J, k(j) = \sum_{i \in I} k(i, j);$$

$$k = \sum_{i \in I} k(i) = \sum_{j \in J} k(j)$$

p étant une classe de I et q une classe de J , on posera :

$$k(p) = \sum_{i \in p} k(i); \quad k(q) = \sum_{j \in q} k(j); \quad k(p, q) = \sum_{i \in p} \sum_{j \in q} k(i, j)$$

Nous nous placerons, dans la suite de l'article, dans le cadre de l'analyse des correspondances et supposerons donc que les espaces $R^{\text{card}(I)}$ et $R^{\text{card}(J)}$ sont munis de leurs métriques du CHI2 respectives, définies à partir du tableau k_{IJ} . La hiérarchie H_I (resp. H_J) est construite à partir du nuage $N(I)$ (resp. $N(J)$) :

$$N(I) = \{(f_I^i, f_i)/i \in I\}, \quad f_I^i = \left(\frac{k(i, j)}{k(i)} \right)_{j \in J}, \quad f_i = \frac{k(i)}{k}$$

$$N(J) = \{(f_J^j, f_j)/j \in J\}, \quad f_J^j = \left(\frac{k(i, j)}{k(j)} \right)_{i \in I}, \quad f_j = \frac{k(j)}{k}$$

Le critère d'agrégation utilisé pour la construction des deux hiérarchies est le critère de Ward. Nous le noterons : $\nu(n)$ pour $n \in H_I$ et $\nu(m)$ pour $m \in H_J$. On rappelle que : $\forall n = (p_1, p_2) \in H_I$

$$\nu(n) = \frac{k(p_1)k(p_2)}{k(p_1) + k(p_2)} \cdot \sum_{j \in J} \frac{1}{k(j)} \left(\frac{k(p_1, j)}{k(p_1)} - \frac{k(p_2, j)}{k(p_2)} \right)^2$$

($\nu(m)$, $m \in H_J$, se définit de manière analogue).

3. Base de $R^{\text{card}(J)}$ (resp. $R^{\text{card}(I)}$) associée à la hiérarchie H_J (resp. H_I)

3.1. Définition

Considérons, par exemple, $R^{\text{card}(J)}$. On pose, alors : $e_j^0 = \frac{k(j)}{k}$, $\forall j \in J$

$$\forall m = (q_1, q_2) \in H_J, e_J^m(j) = \begin{cases} \frac{k(j)}{k(q_1)} & j \in q_1 \\ -\frac{k(j)}{k(q_2)} & j \in q_2 \\ 0 & \text{sinon} \end{cases}$$

Nous noterons de la même façon le nœud $m = (q_1, q_2)$ de H_J et le numéro m qui lui est affecté et qui varie entre 1 et $\text{card}(J) - 1$.

On obtient, ainsi, une suite de $\text{card}(J)$ vecteurs $\{e_J^m / 0 \leq m \leq \text{card}(J) - 1\}$, dont on démontre facilement les propriétés suivantes.

3.2. Propriétés

$R^{\text{card}(J)}$ étant muni de la métrique du CHI2, on montre que :

- a) $\{e_J^m / 0 \leq m \leq \text{card}(J) - 1\}$ est une base orthogonale de $R^{\text{card}(J)}$.
- b) $\forall m \geq 1, m = (q_1, q_2) \in H_J$

$$\|e_J^m\| = \frac{[k(q_1) + k(q_2)]k}{k(q_1) \cdot k(q_2)}$$

3.3. Remarque

Une base $\{e_I^n / 0 \leq n \leq \text{card}(I) - 1\}$ est définie dans $R^{\text{card}(I)}$ de manière analogue.

4. Décomposition de distances carrées dans les bases précédentes

Introduisons tout d'abord les notations suivantes.

4.1. Définitions et notations

On pose :

- a) n étant une classe d'éléments de I et $m = (q_1, q_2)$ un nœud de H_J

$$\Delta(n, (q_1, q_2)) = \frac{k(q_1) \cdot k(q_2)}{k(q_1) + k(q_2)} \cdot k \cdot \left(\frac{k(n, q_1)}{k(n) \cdot k(q_1)} - \frac{k(n, q_2)}{k(n) \cdot k(q_2)} \right)^2$$

b) m étant une classe d'éléments de J et $n = (p_1, p_2)$ un nœud de H_I

$$\Delta((p_1, p_2), m) = \frac{k(p_1) \cdot k(p_2)}{k(p_1) + k(p_2)} \cdot k \cdot \left(\frac{k(p_1, m)}{k(p_1) \cdot k(m)} - \frac{k(p_2, m)}{k(p_2) \cdot k(m)} \right)^2$$

c) $n = (p_1, p_2)$ et $m = (q_1, q_2)$ étant 2 nœuds respectivement de H_I et H_J

$$\begin{aligned} \Delta((p_1, p_2), (q_1, q_2)) &= \frac{k(p_1) \cdot k(p_2)}{k(p_1) + k(p_2)} \cdot \frac{k(q_1) \cdot k(q_2)}{k(q_1) + k(q_2)} \\ &\cdot \left(\frac{k(p_1, q_1)}{k(p_1) \cdot k(q_1)} - \frac{k(p_2, q_1)}{k(p_2) \cdot k(q_1)} + \frac{k(p_2, q_2)}{k(p_2) \cdot k(q_2)} - \frac{k(p_1, q_2)}{k(p_1) \cdot k(q_2)} \right)^2 \end{aligned}$$

d) $m = (q_1, q_2)$ et P_I représentant respectivement un nœud de H_J et une partition de I :

$$\nu(m/P_I) = \frac{k(q_1) \cdot k(q_2)}{k(q_1) + k(q_2)} \sum_{n \in P_I} \frac{1}{k(n)} \left(\frac{k(n, q_1)}{k(q_1)} - \frac{k(n, q_2)}{k(q_2)} \right)^2$$

e) $n = (p_1, p_2)$ et Q_J représentant respectivement un nœud de H_I et une partition de J :

$$\nu(n/Q_J) = \frac{k(p_1) \cdot k(p_2)}{k(p_1) + k(p_2)} \sum_{m \in Q_J} \frac{1}{k(m)} \left(\frac{k(p_1, m)}{k(p_1)} - \frac{k(p_2, m)}{k(p_2)} \right)^2$$

4.2. Remarques

On vérifie facilement que, lorsque les partitions P_I et Q_J sont les partitions des singletons, on a : $\nu(m/P_I) = \nu(m)$ et $\nu(n/Q_J) = \nu(n)$.

4.3. Propriétés

En utilisant les bases orthogonales associées aux hiérarchies H_I et H_J (§ 3), on obtient facilement les résultats suivants :

a) $n, g(n)$, g étant respectivement une classe d'éléments de I , son centre de gravité et le centre de gravité général de $N(I)$:

$$d^2(g(n), g) = \sum_{\substack{m \in H_J \\ m=(q_1, q_2)}} \Delta(n, (q_1, q_2))$$

Démonstration :

$\langle \cdot, \cdot \rangle$ désignant le produit scalaire dans $R^{\text{card}(J)}$, on peut écrire que :

$$d^2(g(n), g) = \sum_{m=0}^{\text{card}(J)-1} \frac{1}{\|e_J^m\|^2} [\langle g(n) - g, e_J^m \rangle]^2$$

On obtient le résultat cherché en remplaçant la norme carrée de e_J^m par sa valeur (voir § 3.2) et en remarquant que :

$$\text{Pour } m = 0, \quad \langle g(n) - g, e_J^m \rangle = 0$$

$$\text{Pour } m \neq 0, \quad \langle g(n) - g, e_J^m \rangle = \langle g(n), e_J^m \rangle = k \cdot \left[\frac{k(n, q_1)}{k(n) \cdot k(q_1)} \frac{k(n, q_2)}{k(n) \cdot k(q_2)} \right]$$

(le numéro du nœud $m = (q_1, q_2)$ étant également noté m).

b) $m, g(m), g$ étant respectivement une classe d'éléments de J , son centre de gravité et le centre de gravité général de $N(J)$:

$$d^2(g(m), g) = \sum_{\substack{n \in HI \\ n=(p_1, p_2)}} \Delta((p_1, p_2), m)$$

c) P_I et $m = (q_1, q_2)$ désignant respectivement une partition de I et un nœud de H_J :

$$\nu(m/P_I) = \sum_{n \in P_I} \frac{k(n)}{k} \cdot \Delta(n, (q_1, q_2))$$

d) Q_J et $n = (p_1, p_2)$ désignant respectivement une partition de J et un nœud de H_I :

$$\nu(n/Q_J) = \sum_{m \in Q_J} \frac{k(m)}{k} \cdot \Delta((p_1, p_2), m)$$

e) $n = (p_1, p_2)$ et $m = (q_1, q_2)$ étant deux nœuds respectivement de H_I et H_J :

$$\nu(n) = \sum_{\substack{m \in H_J \\ m=(q_1, q_2)}} \Delta((p_1, p_2), (q_1, q_2)) \quad \nu(m) = \sum_{\substack{n \in H_I \\ n=(p_1, p_2)}} \Delta((p_1, p_2), (q_1, q_2))$$

f) In désignant l'inertie totale du tableau k_{IJ} :

$$\text{In} = \sum_{\substack{n \in H_I \\ m=(p_1, p_2)}} \sum_{\substack{m \in H_J \\ m=(q_1, q_2)}} \Delta((p_1, p_2), (q_1, q_2))$$

(La démonstration de ces propriétés est analogue à celle de 4.3 a))

4.4. Autres propriétés

a) P_I et P_I^* étant deux partitions de I , P_I^* se déduisant de P_I par l'agrégation des 2 classes p_1 et p_2 de P_I , on a :

$$\forall m = (q_1, q_2) \in H_J \quad \nu(m/P_I) - \nu(m/P_I^*) = \Delta[(p_1, p_2), (q_1, q_2)]$$

b) Q_J et Q_J^* étant deux partitions de J , Q_J^* se déduisant de Q_J par l'agrégation des 2 classes q_1 et q_2 de Q_J , on a :

$$\forall n = (p_1, p_2) \in H_I \quad \nu(n/Q_J) - \nu(n/Q_J^*) = \Delta[(p_1, p_2), (q_1, q_2)]$$

c) P_I étant la partition de I correspondant à une coupure $c(H_I)$ de la hiérarchie H_I , on note $H_{\text{inf}}(P_I)$ et $H_{\text{sup}}(P_I)$ les ensembles des nœuds de H_I situés pour l'un sous la coupure $c(H_I)$ et pour l'autre au dessus. On montre alors que : $\forall m = (q_1, q_2) \in H_J$,

$$\nu(m) = \nu(m/P_I) + \sum_{\substack{n \in H_{\text{inf}}(P_I) \\ n=(p_1, p_2)}} \Delta[(p_1, p_2), (q_1, q_2)]$$

d) Q_J étant la partition de J correspondant à une coupure $c(H_J)$ de la hiérarchie H_J , on note $H_{\text{inf}}(Q_J)$ et $H_{\text{sup}}(Q_J)$ les ensembles des nœuds de H_J situés pour l'un sous la coupure $c(H_J)$ et pour l'autre au-dessus. On montre alors que : $\forall n = (p_1, p_2) \in H_I$,

$$\nu(n) = \nu(n/Q_J) + \sum_{\substack{m \in H_{\text{inf}}(Q_J) \\ m=(q_1, q_2)}} \Delta[(p_1, p_2), (q_1, q_2)]$$

e) Soit P_I et Q_J partitions de I et de J associées respectivement aux coupures $c(H_I)$ et $c(H_J)$ de H_I et H_J . $\text{In}(P_I, Q_J)$ désignant l'inertie du tableau croisant P_I et Q_J déduit par cumul des lignes et colonnes de k_{IJ} , on montre, avec des notations analogues, que :

$$\text{In}(P_I, Q_J) = \sum_{\substack{m \in H_{\text{sup}}(P_I) \\ n=(p_1, p_2)}} \sum_{\substack{m \in H_{\text{sup}}(Q_J) \\ m=(q_1, q_2)}} \Delta[(p_1, p_2), (q_1, q_2)]$$

Démonstration

a) On utilise la propriété c) du § 4.3. On en déduit alors : $\forall m = (q_1, q_2)$

$$\begin{aligned} \nu(m/P_I) - \nu(m/P_I^*) &= \frac{k(p_1)}{k} \cdot \Delta[p_1, (q_1, q_2)] + \frac{k(p_2)}{k} \cdot \Delta[p_2, (q_1, q_2)] \\ &\quad - \frac{k(p_1 \cup p_2)}{k} \Delta[p_1 \cup p_2, (q_1, q_2)] \end{aligned}$$

On obtient, ainsi, $\Delta[(p_1, p_2), (q_1, q_2)]$ après un calcul simple.

b) Démonstration analogue

c) On considère la succession de partitions $P_I^0 = P_I, P_I^1, \dots, P_I^s$ obtenues par une suite de coupures successives de H_I , P_I^s étant la partition composée de la seule classe I . Il est évident que : $\forall m = (q_1, q_2) \in H_J, \nu(m/P_I^s) = 0$. D'autre part, d'après la propriété a) ci-dessus, on a :

$$\forall u \in [0, s-1], \quad \nu(m/P_I^u) - \nu(m/P_I^{u+1}) = \Delta[(p_1^u, p_2^u), (q_1, q_2)]$$

où p_1^u et p_2^u sont les classes de P_I^u qui ont été agrégés pour obtenir P_I^{u+1} . En ajoutant ces égalités membre à membre, on en déduit que :

$$\nu(m/P_I) = \sum_{\substack{n \in H_{\text{sup}}(P_I) \\ n=(p_1, p_2)}} \Delta[(p_1, p_2), (q_1, q_2)]$$

On déduit alors la propriété cherchée de celle du § 4.3 f)

d) démonstration analogue

e) On a classiquement : $\text{In}(P_I, Q_J) = \sum\{\nu(m/P_I)/m \in H_{\text{sup}}(Q_J), m = (q_1, q_2)\}$ ce qui donne le résultat cherché en remplaçant $\nu(m/P_I)$ par sa valeur obtenue ci-dessus.

4.5. Remarques

A partir des nœuds des 2 hiérarchies H_I et H_J , on formera une suite de sous-tableaux emboîtés de k_{IJ} , qui vont du tableau k_{IJ} lui-même au tableau de dimensions 2×2 croisant l'aîné et le benjamin du sommet de H_I avec ceux du sommet de H_J . Cette suite de sous-tableaux s'obtient de la façon suivante : à chaque étape, on forme les deux sous-tableaux de k_{IJ} déduits de l'agrégation de l'aîné et du benjamin du prochain nœud pour l'un de H_I , pour l'autre de H_J .

On choisit alors le sous tableau d'inertie maximum (Le calcul de cette inertie peut se faire par l'application de la formule § 4.4 e)). L'examen de la suite des inerties des sous-tableaux retenus peut ainsi permettre de déterminer conjointement les coupures des hiérarchies H_I et H_J .

5. Classification hiérarchique croisée

Les 2 hiérarchies H_I et H_J ont été construites indépendamment l'une de l'autre à partir de k_{IJ} . Cependant, leurs constructions peuvent recevoir une interprétation dans le cadre d'une classification hiérarchique croisée. En effet, supposons par exemple que des regroupements aient eu lieu sur I et sur J et donnés naissance aux 2 partitions P_I et Q_J de I et J . (p_1, p_2) (resp. (q_1, q_2)) étant deux classes de

P_I (resp. Q_J), on note P_I^* (resp. Q_J^*) la partition qui se déduit de P_I (resp. Q_J) par l'agrégation de p_1 et p_2 (resp. q_1 et q_2). D'après la définition § 4.1e), on a noté $\nu(n/Q_J)$, avec $n = (p_1, p_2)$, l'indice du nœud n calculé sur le tableau croisant P_I et Q_J . La propriété § 4.4 b) montre que cet indice diminue de la quantité positive $\Delta[(p_1, p_2), (q_1, q_2)]$ lorsque l'on passe de la partition Q_J à la partition Q_J^* . Cette quantité $\Delta[(p_1, p_2), (q_1, q_2)]$ représente donc la déformation de l'indice $\nu(n/Q_J)$ avec $n = (p_1, p_2)$ lorsque l'on passe de Q_J à Q_J^* , ou représente (§ 4.4 a)) la déformation de $\nu(m/P_I)$ avec $m = (q_1, q_2)$ lorsque l'on passe de P_I à P_I^* . En conséquence, d'après la propriété 4.4 d), l'indice $\nu(n)$ du nœud $n = (p_1, p_2) \in H_I$ se présente sous la forme de l'addition de $\nu(n/Q_J)$ et de la somme des déformations $\Delta[(p_1, p_2), (q_1, q_2)]$ de tous les nœuds $m = (q_1, q_2)$ de H_J formés avant la constitution de la partition Q_J . L'indice $\nu(m)$, avec $m \in H_J$, reçoit une interprétation analogue (§ 4.4 c)).

6. Aides à l'interprétation conjointe de H_I et H_J

Les programmes d'aides à l'interprétation usuels permettent d'expliquer l'écart entre le centre de gravité de chaque nœud de H_I et le centre de gravité général, ainsi que l'écart entre l'aîné et le benjamin de chacun de ces nœuds, et ceci en fonction des éléments j de J . Les résultats fournis par ces programmes deviennent rapidement fastidieux à dépouiller lorsque le cardinal de J augmente. On se propose, dans cet article, d'expliquer les nœuds de H_I en fonction de ceux de H_J (ou réciproquement) et ceci en se limitant aux nœuds supérieurs de ces deux hiérarchies. Nous utiliserons, pour y parvenir, les formules démontrées aux paragraphes précédents. D'autre part, nous verrons plus loin que même dans le cas de tableaux de dimensions modestes la méthode proposée apporte un plus par rapport aux techniques usuelles.

Un logiciel fournissant des indices calculés à partir des formules démontrées dans cet article a été écrit en Fortran. Nous donnons, ci-dessous, les notations de ces indices ainsi que le numéro des paragraphes où se trouvent les définitions de référence. Une application à un jeu de données sera proposée ensuite.

6.1.

NI = numéro du nœud de HI

AI = numéro de l'aîné de NI

BI = numéro du benjamin de NI

PAI = poids de la classe AI

PBI = poids de la classe BI

(NJ, AJ, BJ, PAJ, PBJ ont des définitions analogues pour la hiérarchie HJ)

6.2. (voir § 4.3 a))

$\forall m = (q_1, q_2) \in H_J, \quad \forall n \in H_I$

$$\text{dnIG}(q_1, q_2) = \frac{\Delta[n, (q_1, q_2)]}{d^2(g(n), g)} * 1\,000$$

$$\text{FAJ}(n; (q_1, q_2)) = \frac{k(n, q_1)}{k(q_1)} * 1\,000$$

$$\text{FBJ}(n; (q_1, q_2)) = \frac{k(n, q_2)}{k(q_2)} * 1\,000$$

6.3. (voir § 4.3 b))

$$\forall n = (p_1, p_2) \in H_I, \quad \forall m \in H_J$$

$$\text{dnJG}(p_1, p_2) = \frac{\Delta[(p_1, p_2), m]}{d^2(g(m), g)} * 1\,000$$

$$\text{FAI}((p_1, p_2); m) = \frac{k(p_1, m)}{k(p_1)} * 1\,000$$

$$\text{FBI}((p_1, p_2); m) = \frac{k(p_2, m)}{k(p_2)} * 1\,000$$

6.4. (voir § 4.3 c) et (voir § 4.3 d))

P_I étant une partition de I liée à une coupure $c(H_I)$ choisie par l'utilisateur,
 $\forall m = (q_1, q_2) \in H_J, \forall n \in P_I$

$$\text{dABJ}(n, (q_1, q_2)) = \frac{\frac{k(n)}{k} \cdot \Delta(n, (q_1, q_2))}{\nu(m/P_I)} * 1\,000$$

Q_J étant une partition de J liée à une coupure $c(H_J)$ choisie par l'utilisateur,
 $\forall n = (p_1, p_2) \in H_I, \forall m \in Q_J$

$$\text{dABI}((p_1, p_2), m) = \frac{\frac{k(m)}{k} \cdot \Delta((p_1, p_2), m)}{\nu(n/Q_J)} * 1\,000$$

6.5. (voir § 4.3 e))

$$\forall n = (p_1, p_2) \in H_I \quad \forall m(q_1, q_2) \in H_J$$

$$\&\text{ABI}(n, m) = \frac{\Delta[(p_1, p_2), (q_1, q_2)]}{\nu(n)} * 1\,000$$

$$\&\text{AI}(n, m) = \left(\frac{k(p_1, q_1)}{k(p_1) \cdot k(q_1)} - \frac{k(p_1, q_2)}{k(p_1) \cdot k(q_2)} \right) * k * 1\,000$$

$$\begin{aligned} \&BI(n, m) &= \left(\frac{k(p_2, q_1)}{k(p_2) \cdot k(q_1)} - \frac{k(p_2, q_2)}{k(p_2) \cdot k(q_2)} \right) * k * 1\,000 \\ \&ABJ(n, m) &= \frac{\Delta[(p_1, p_2), (q_1, q_2)]}{\nu(m)} * 1\,000 \\ \&AJ(n, m) &= \left(\frac{k(p_1, q_1)}{k(p_1) \cdot k(q_1)} - \frac{k(p_2, q_1)}{k(p_2) \cdot k(q_1)} \right) * k * 1\,000 \\ \&BJ(n, m) &= \left(\frac{k(p_1, q_2)}{k(p_1) \cdot k(q_2)} - \frac{k(p_2, q_2)}{k(p_2) \cdot k(q_2)} \right) * k * 1\,000 \end{aligned}$$

6.6. Remarques

P_I et Q_J étant deux partitions de I et de J définies à partir de 2 coupures de H_I et H_J (celles-ci étant fixées par l'utilisateur), le logiciel permet de stocker 3 tableaux :

$$a) \quad (k(n) \cdot \Delta[n, (q_1, q_2)])_{\substack{n \in P_I \\ m = (q_1, q_2) \in H_J}}$$

Ce tableau peut être soumis à l'AFC et permettra de placer d'une part chaque classe n de P_I à proximité des nœuds (q_1, q_2) qui expliquent le mieux l'écart $d^2(g(n), g)$ (§ 4.3 a)) et d'autre part chaque nœud $m = (q_1, q_2)$ de H_J près des classes n de P_I qui expliquent le mieux l'indice $\nu(m/P_I)$ (§ 4.3 c)).

$$b) \quad (k(m) \cdot \Delta[(p_1, p_2), m])_{\substack{n = (p_1, p_2) \in H_I \\ m \in Q_J}}$$

Ce second tableau de même nature peut être également soumis à l'AFC

$$c) \quad (\Delta[p_1, p_2], (q_1, q_2))_{\substack{n = (p_1, p_2) \in H_I \\ m = (q_1, q_2) \in H_J}}$$

Ce dernier tableau rassemble les valeurs décomposant les indices $\nu(n)$, $n \in H_I$ (resp. $\nu(m)$, $m \in H_J$) en fonction des nœuds m de H_J (resp. n de H_I).

7. Application du logiciel à un tableau de données

Un logiciel d'aides à l'interprétation mutuelle de 2 hiérarchies H_I et H_J a donc été écrit en Fortran. Le tableau de contingence choisi pour l'illustrer est tiré du livre «l'analyse des données» de BENZECRI et provient d'une enquête effectuée pour la Régie Française des tabacs sur 100 fumeurs. Ce tableau croise un ensemble I de 12 marques et un ensemble J de 11 attributs et regroupe les valeurs k_{ij} représentant le nombre d'individus estimant que l'attribut j convenait bien à la marque i .

$I = \{\text{Orly(ORLY), Alezan(ALEZ), Corsaire(CORS), Directoire(DIRE), Ducat(DUCA), Fontenoy(FONT), Icare(ICAR), Zodiaque(ZODI), Pavois(PAVI), Cocker(COCK), Escale(ESCA), Hotesse(HOTE)}\}$

$J = \{\text{Vieillot-desuet(VIEU), Nouveau-riche(RICH), Sobre-elegant(ELEG), Cocasse-ridicule (RIDI), Racé(RACE), Mièvre(MIEV), Distingué(DIST),}$

Vulgaire-commun(VULG), Pour un homme (HOMM), Pour une femme(FEMM),
Pour une petite nature(P.NA) }

On trouvera le contenu du tableau k_{IJ} en annexe.

Sur ce tableau ont été construites 2 hiérarchies H_I et H_J à l'aide du critère de Ward et on trouvera, en annexe, le détail des résultats obtenus par notre logiciel d'aides à l'interprétation mutuelle de 2 hiérarchies.

Ce dernier demande d'abord à l'utilisateur de lui fixer 2 coupures $c(H_I)$ et $c(H_J)$ pour les deux hiérarchies H_I et H_J . Puis il calcule, par hiérarchie, 3 types de tableaux. Considérons, par exemple, la hiérarchie H_I . Les deux premiers tableaux décrivent (en particulier à l'aide de l'indice dnIG) pour l'un les nœuds supérieurs de H_I et pour l'autre les classes de la partition P_I associée à $c(H_I)$. Un troisième tableau croisant les nœuds supérieurs retenus pour les deux hiérarchies H_I et H_J est ensuite calculé. On y trouve, en particulier, l'indice &ABI.

Dans les sorties, on vérifiera que les indices dnIG, dnJG,&ABI,&ABJ sont à lire «en lignes», et que les indices dABI, dABJ sont à lire en «colonnes».

Le logiciel détermine également une suite de sous tableaux, à partir des nœuds de H_I et H_J (voir § 4.5). A chaque couple de nœuds (N_I, N_J) , on associe pour chacun d'eux la coupure de l'arbre correspondant se situant sous ce nœud et on calcule l'inertie du sous tableau croisant les deux partitions P_I et P_J déterminées par ces deux coupures.

Ainsi, dans le cas de notre exemple, nous avons retenu 5 nœuds supérieurs pour H_I (représentant 80.3% de l'inertie de k_{IJ}) et 6 nœuds supérieurs pour H_J (représentant 83.1% de l'inertie de k_{IJ}). Le sous tableau obtenu croisant les partitions de I et de J ainsi obtenues a une inertie de 68.8% de celle de k_{IJ} (voir annexe).

Nous avons, à titre d'illustration, procédé à l'interprétation de la hiérarchie H_I à partir des nœuds de H_J . Un travail analogue pourrait être fait pour l'interprétation de H_J en fonction de H_I . Ainsi, pour chaque nœud retenu N_I de H_I , en utilisant l'indice dnIG, nous avons repéré les nœuds (A_J, B_J) de H_J pour lesquels cet indice prend une valeur élevée. Nous avons, ensuite, reporté sur H_I l'une des 2 classes AJ ou BJ, celle pour laquelle le pourcentage FAJ ou FBJ est le plus élevé, autrement dit celle pour laquelle les attributs qu'elle contient sont plus fréquemment associés aux marques de la classe N_I . On trouvera, ci-dessous, les deux hiérarchies restreintes de H_I et H_J associées aux coupures réalisées. Nous avons placé devant chaque nœud l'indice I ou J selon qu'il s'agit de H_I ou H_J . Nous avons, enfin, (à l'aide de l'indice dnIG) résumé dans un tableau puis reporté sur l'arbre de H_I , les classes de H_J qui expliquent chaque nœud retenu de H_I . L'ensemble des résultats ainsi obtenus pourrait être complété par l'examen des indices dABI et &ABI.

L'ensemble des marques (23I) se partage en les deux classes 20I et 22I. Il apparaît, ainsi pour la classe 20I = {CORS, COCK, ICAR, ZODI} en balayant (en ligne) l'ensemble des valeurs de l'indice dnIG que les nœuds $NJ = 21J = (20J, 17J)$ et $NJ = 19J = (1J, 15J)$ prennent pour cet indice les valeurs les plus élevées : à savoir respectivement 0.433 et 0.126 (voir annexe § 3a). Ceci signifie que la distance carrée entre le centre de gravité de la classe 20I et le centre de gravité de l'ensemble de toutes les marques s'explique à 43.3% par le nœud 21J et à 12.6% par le nœud 19J. De manière plus précise, alors que le poids de la classe

20I est de 0.322 (32.2%), on constate (voir annexe § 3a), indices FAJ et FBJ, d'une part concernant le nœud $21J = (20J, 17J)$ que parmi les individus ayant choisis les adjectifs de $20J = \{\text{VIEU, RIDI, P.NA, FEMM, VULG, HOMM}\}$, 40.8% d'entr'eux les ont affectés aux marques de la classe 20I, et d'autre part concernant le nœud $19J = (1J, 15J)$ que parmi les individus ayant choisi les adjectifs de la classe $15J = \{\text{RIDI, MIEV, P.NA}\}$ 48.9% d'entr'eux les ont affectés aux marques de la classe 20I. Ainsi, la classe 20I se trouve plutôt expliquée par les classes 15J et 20J (ce qui est reporté sur le graphique annexe § 1 et précisé dans le résumé annexe § 2). En procédant de même, on vérifie que les classes $3I = \{\text{CORS}\}$ (voir annexe § 3b) et $16I = \{\text{COCK, ICAR, ZODI}\}$ (voir annexe § 3a), dont la réunion redonne 20I, sont expliquées respectivement par les classes $14J = \{\text{VULG, HOMM}\}$ et $15J = \{\text{RIDI, MIEV, P.NA}\}$. On observera, de même, que les classes $17I = \{\text{HOTE, ESCA}\}$, $18I = \{\text{ORLY, PAVO, FONT}\}$ et $15I = \{\text{DIRE, DUCA}\}$ sont respectivement associées aux attributs $10J = \{\text{FEMM}\}$, $2J = \{\text{RICH}\}$ et $1J = \{\text{VIEU}\}$ et qu'enfin la marque $2I = \{\text{ALEZ}\}$ convient bien à l'adjectif $5J = \{\text{RACE}\}$.

Le fait que l'on explique les nœuds d'une hiérarchie en fonction des nœuds de l'autre peut permettre de mettre à jour des nuances que ne permettent pas les aides à l'interprétation classiques qui expliquent les classes de H_I en prenant de manière séparée chaque «colonne» du tableau. Ainsi, dans le cadre de ce petit exemple, on s'aperçoit que les 2 adjectifs «Pour un homme» et «Pour une femme» ne jouent pas des rôles opposés puisqu'ils se retrouvent dans la même classe 18J qui explique la classe $17I = \{\text{HOTE, ESCA}\}$. En examinant les résultats de manière plus précise, on peut voir que ces deux adjectifs sont ceux qui ont été plus fréquemment attribués à la marque ESCA.

8. Conclusion

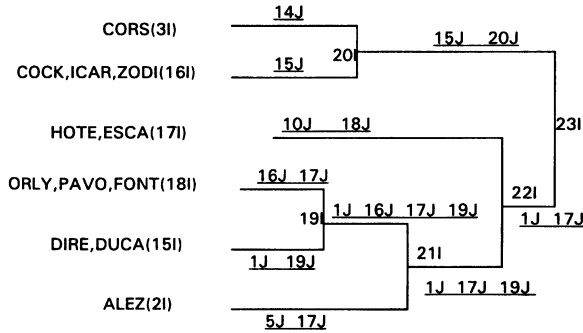
Cette méthode d'aides à l'interprétation mutuelle de deux hiérarchies trouvera son utilité non seulement dans le cas où k_{IJ} est de dimensions élevées, mais aussi pour des tableaux de dimensions modestes puisqu'elle peut mettre en lumière des faits qui peuvent passer inaperçus en utilisant une technique classique. Enfin des généralisations, qui pourront faire l'objet d'un second article, sont envisagées. Elles s'appliqueront dans le cas soit de tableaux de contingence à plusieurs entrées soit de juxtaposition de tableaux de contingence, et permettront d'établir des aides à l'interprétation mutuelle entre r hiérarchies (avec $r \geq 2$).

Références bibliographiques

- BENZECRI, J.P. - LEBEAUX, M.O. - JAMBU, M. (1980) : Aides à l'interprétation en classification automatique. CAD V n°1, p101 à 123
- GOVAERT, G. (1983) : Classification croisée. Thèse d'Etat Paris 6
- GREENACRE, M. (1988) : Clustering the rows and Columns of a Contingency Table. Journal of Classification 5 : p 39 à 51.
- WEISS, M.C. (1978) : Décomposition hiérarchique du khi-deux associé à un tableau de contingence à plusieurs entrées. RSA vol XXVI n°1 p23, 33

Interprétation de H_I en fonction de H_J

1) Représentation de H_I

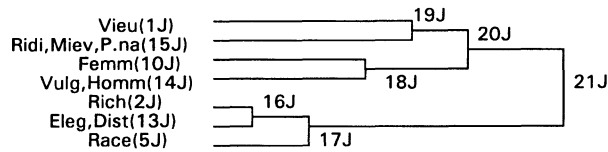


CLASSES DE LA HIERARCHIE H_I		CLASSES DE LA HIERARCHIE H_J	
1I	ORLY	1J	VIEU
2I	ALEZ	2J	RICH
3I	CORS	3J	ELEG
4I	DIRE	4J	RIDI
5I	DUCA	5J	RACE
6I	FONT	6J	MIEV
7I	ICAR	7J	DIST
8I	ZODI	8J	VULG
9I	PAVI	9J	HOMM
10I	COCK	10J	FEMM
11I	ESCA	11J	P.NA
12I	HOTE	12J	MIEV P.NA
13I	ICAR ZODI	13J	ELEG DIST
14I	ORLY PAVO	14J	VULG HOMM
15I	DIRE DUCA	15J	RIDI MIEV P.NA
16I	COCK ICAR ZODI	16J	RICH ELEG DIST
17I	HOTE ESCA	17J	RICH ELEG DIST RACE
18I	ORLY PAVO FONT	18J	FEMM VULG HOMM
19I	ORLY PAVO FONT DIRE DUCA	19J	VIEU RIDI MIEV P.NA
20I	CORS COCK ICAR ZODI	20J	VIEU RIDI MIEV P.NA FEMM VULG HOMM
21I	ORLY PAVO FONT DIRE DUCA ALEZ	21J	VIEU RIDI MIEV P.NA FEMM
22I	HOTE ESCA ORLY PAVO FONT DIRE DUCA ALEZ		VULG HOMM RICH ELEG DIST RACE
23I	CORS COCK ICAR ZODI HOTE ESCA ORLY PAVO FONT DIRE DUCA ALEZ		

2) Tableau résumé

N_I	Poids(%)	: dnIG				
		0	100	200	400	1000
22I	67	21J ⁺⁺⁺ →(20J = 59.1%, 17J = 82.2%); 19J ⁺ →(1J = 73.3%, 15J = 51%)				
21I	48.5	21J ⁺⁺⁺ →(20J = 39.4%, 17J = 63.9%); 19J ⁺ →(1J = 65%, 15J = 38 %)				
20I	32.3	21J ⁺⁺⁺ →(20J = 40.8%, 17J = 17.7%); 19J ⁺ →(1J = 26.6%, 15J = 48.9%)				
19I	38.3	21J ⁺ →(20J = 34.3%, 17J = 46.8%); 20J ⁺ →(19J = 41.3%,18J = 25.3%) 17J ⁺ →(16J = 52.9%, 5J = 26.2%); 19J ⁺⁺ →(1J = 63.3%, 15J = 32.3%)				
2I	9.5	21J ⁺⁺⁺ →(20J = 5.1%,17J = 17%); 17J ⁺⁺ →(16J = 12.2%, 5J = 33.3%)				
3I	9	21J ⁺⁺ →(20J = 13%, 17J = 2.3%); 20J ⁺⁺ →(19J = 7.7%,18J = 19.8%) 18J ⁺⁺ →(10J = 8%,14J = 26%)				
16I	23.2	21J ⁺⁺ →(20J = 27.8%, 17J = 15.4%); 20J ⁺⁺ →(19J = 34.7%,18J = 18.8%) 19J ⁺⁺ →(1J = 15%,15J = 42.8%)				
18I	20.5	21J ⁺⁺⁺ →(20J = 15.3%,17J = 29.5%); 17J ⁺ →(16J = 33.8%, 5J = 15.1%) 16J ⁺ →(2J = 44.5%, 13J = 28.5%)				
15I	18.2	20J ⁺⁺ →(19J = 26.3%,18J = 9.5%); 19J ⁺⁺⁺ →(1J = 46.6%,15J = 18%)				
17I	19	20J ⁺⁺ →(19J = 19.5%, 18J = 30%); 18J ⁺⁺ →(10J = 50%,14J = 19.4%)				

3) Représentation de HJ



1J	VIEU	12J	MIEV P.NA
2J	RICH	13J	ELEG DIST
3J	ELEG	14J	VULG HOMM
4J	RIDI	15J	RIDI MIEV P.NA
5J	RACE	16J	RICH ELEG DIST
6J	MIEV	17J	RICH ELEG DIST RACE
7J	DIST	18J	FEMM VULG HOMM
8J	VULG	19J	VIEU RIDI MIEV P.NA
9J	HOMM	20J	VIEU RIDI MIEV P.NA FEMM VULG HOMM
10J	FEMM	21J	VIEU RIDI MIEV P.NA FEMM VULG HOMM
11J	P.NA		RICH ELEG DIST RACE

ANNEXES

Résultats détaillés du logiciel d'aides à l'interprétation

1) L'inertie du tableau k_{IJ} vaut : 0.55719

	VIEU	RICH	ELEG	RIDI	RACE	MIEV	DIST	VULG	HOMM	FEMM	P.NA
ORLY	1	20	9	1	4	3	11	4	9	9	7
ALEZ	2	9	23	3	33	9	9	4	12	3	5
CORS	14	1	1	15	7	1	1	32	23	9	2
DIRE	38	11	15	15	8	7	17	2	4	8	7
DUCA	18	10	7	6	3	7	4	6	7	4	11
FONT	10	9	11	5	6	5	21	0	13	2	2
ICAR	9	1	6	12	6	12	6	9	5	6	6
ZODI	5	1	2	18	4	9	1	7	5	8	11
PAVI	9	20	7	4	5	6	5	3	10	1	9
COCK	4	9	12	25	15	9	4	10	5	6	24
ESCA	0	7	3	2	5	6	5	12	13	23	10
HOTE	10	12	17	2	3	13	27	7	9	33	5

2) Détermination des sous-tableaux de k_{IJ} et de leurs inerties

NHI	NHJ	INER(%)	HISTOGRAMME DES INERTIES DES SOUS-TABLEAUX
13	12	1000	*****
14	12	983	*****
15	12	963	*****
16	12	938	*****
16	13	910	*****
17	13	892	*****
17	14	860	*****
18	14	816	*****
19	14	772	*****
19	15	730	*****
19	16	688	*****
19	17	666	*****
20	17	586	*****
20	18	486	*****
20	19	387	*****
20	20	303	*****
21	20	218	*****
22	20	184	*****
22	21	121	*****
23	21	102	*****

3) Aides à l'interprétation de H_I en fonction de H_J

a) Explication des écarts $d_2(n, g)$, n étant un nœud supérieur de H_I , en fonction des nœuds supérieurs $m = (AJ, BJ)$ de H_J

NJ					21			20			19			18					
AJ	BJ	20	17	9	18	1	15	10	14	PAJ	PBJ	629	370	353	276	102	251	95	180
NI	AI	BI	PAI	PBI	FAJ	FBJ	dnIG	FAJ	FBJ	dnIG	FAJ	FBJ	dnIG	FAJ	FBJ	dnIG	FAJ	FBJ	dnIG
23	20	22	323	676	1000	1000	0	1000	1000	0	1000	1000	0	1000	1000	0	1000	1000	0
22	17	21	191	485	591	822	433	574	613	7	733	510	126	741	545	84	741	545	84
21	19	2	389	95	394	639	493	458	312	117	650	380	186	241	350	26	241	350	26
20	3	16	90	232	408	177	433	425	386	7	266	489	126	258	454	84	258	454	84
19	18	15	205	183	343	468	154	413	253	164	633	323	293	214	274	9	214	274	9

NJ					17			16					
AJ	BJ	16	5	2	13	PAJ	PBJ	285	84	94	191		
NI	AI	BI	PAI	PBI	FAJ	FBJ	dnIG	FAJ	FBJ	dnIG	FAJ	FBJ	dnIG
23	20	22	323	676	1000	1000	0	1000	1000	0	1000	1000	0
22	17	21	191	485	865	676	81	890	52	3	890	52	3
21	19	2	389	95	652	595	7	718	620	21	718	620	21
20	3	16	90	232	134	323	81	109	147	3	109	147	3
19	18	15	205	183	529	262	195	636	477	66	636	477	66

b) Explication des écarts $d_2(n, g)$, n étant une classe terminale de P_I , en fonction des nœuds supérieurs $m = (AJ, BJ)$ de H_J

NJ					21				20				19			
AJ	BJ	20	17	19	18	1	15	PAJ	PBJ	629	370	353	276	102	251	
NI	AI	BI	PAI	PBI	FAJ	FBJ	dnIG	ABJ	FAJ	FBJ	dnIG	ABJ	FAJ	FBJ	dnIG	ABJ
17	12	11	117	73	196	182	3	2	115	300	358	294	83	129	10	13
2	-	-	-	-	51	170	421	336	45	58	3	2	16	57	15	20
3	-	-	-	-	130	23	293	286	77	198	248	267	116	61	24	40
16	10	13	105	127	278	154	214	148	347	188	236	180	150	428	340	394
18	14	6	134	71	153	295	402	22	149	157	0	0	166	142	3	3
15	4	5	112	70	189	173	6	3	263	95	399	253	466	180	549	527

NJ					18				17				16			
AJ	BJ	10	14	16	5	2	13	PAJ	PBJ	95	180	285	84	94	191	
NI	AI	BI	PAI	PBI	FAJ	FBJ	dnIG	ABJ	FAJ	FBJ	dnIG	ABJ	FAJ	FBJ	dnIG	ABJ
17	12	11	117	73	500	194	397	539	212	80	77	104	172	232	15	98
2	-	-	-	-	26	75	19	27	122	333	367	532	81	142	29	207
3	-	-	-	-	80	260	22	396	8	70	27	48	9	8	0	0
16	10	13	105	127	178	194	0	1	125	252	63	79	100	138	5	33
18	14	6	134	71	107	184	32	32	338	151	194	194	445	285	137	660
15	4	5	112	70	107	90	1	1	191	11	38	40	190	191	0	0

					23				22				21			
					20		22		17		21		19		2	
NI		AI	BI													
PAI	PBI	323	676													
NJ	AJ	BJ	PAJ	PBJ	FAI	FBI	dnJG	ABI	FAI	FBI	dnJG	ABI	FAI	FBI	dnJG	ABI
1	-	-	-	-	84	111	20	15	44	137	155	135	166	17	223	221
15	4	12	92	158	380	189	499	331	169	197	6	4	208	151	15	13
10	10	10	95	95	76	104	26	18	250	47	866	693	52	26	7	7
14	8	9	82	98	253	145	179	148	183	130	26	24	127	142	1	1
5	-	-	-	-	84	84	0	0	35	103	119	88	57	294	815	684
2	-	-	-	-	31	123	362	204	84	139	79	50	153	80	80	58
13	3	7	96	94	87	241	449	280	232	244	1	1	234	285	17	13

					20				19						
					3		16		18		15				
NI		AI	BI												
PAI	PBI	90	232												
NJ	AJ	BJ	PAJ	PBJ	FAI	FBI	dnJG	ABI	FAI	FBI	dnJG	ABI			
1	-	-	-	-	132	66	37	35	82	260	402	613			
15	4	12	92	158	169	463	349	284	174	246	31	41			
10	10	10	95	95	84	73	1	1	49	55	0	0			
14	8	9	82	98	518	150	611	624	161	88	36	59			
5	-	-	-	-	66	91	8	6	62	51	2	2			
2	-	-	-	-	9	40	12	8	203	97	212	237			
13	3	7	96	94	18	113	51	39	265	200	36	44			

c) Explication des écarts $d_2(AJ, BJ)$, $m = (AJ, BJ)$ étant un nœud supérieur de H_J , en fonction des nœuds supérieurs (AI, BI) de H_I

					23				22				21			
					20		22		17		21		19		2	
NI		AI	BI													
PAI	PBI	323	676													
NJ	AJ	BJ	PAJ	PBJ	&AI	&BI	&ABJ	&AI	&BI	&ABJ	&AI	&BI	&ABJ			
21	20	17	629	370	713	-340	467	74	-504	88	-322	-1246	126			
20	19	18	353	276	118	-56	9	-962	301	327	408	-135	33			
19	1	15	102	251	-690	329	222	-239	554	84	795	-429	112			
18	10	14	95	180	-606	289	175	1596	-225	455	-155	-512	9			
17	16	5	285	84	-583	278	169	688	116	46	685	-2199	669			
16	2	13	94	191	-118	56	10	-310	201	58	407	-637	137			

					20				19						
					3		16		18		15				
NI		AI	BI												
PAI	PBI	90	232												
NJ	AJ	BJ	PAJ	PBI	&AI	&BI	&ABJ	&AI	&BI	&ABJ					
21	20	17	629	370	1182	530	53	-690	91	114					
20	19	18	353	276	-1133	683	396	-39	910	130					
19	1	15	102	251	611	-1198	208	115	1558	197					
18	10	14	95	180	-1990	-67	240	-377	93	21					
17	16	5	285	84	-681	-545	1	906	438	22					
16	2	13	94	191	1	-165	2	775	-5	97					