

BASAVANNEPA TALLUR

Classification et corrélation Application à l'analyse d'un questionnaire parents-élèves sur l'enseignement des mathématiques

Publications de l'Institut de recherche mathématiques de Rennes, 1993-1994, fascicule 3
« Fascicule de didactique des mathématiques », , exp. n° 5, p. 1-18

http://www.numdam.org/item?id=PSMIR_1993-1994__3_A5_0

© Département de mathématiques et informatique, université de Rennes, 1993-1994, tous droits réservés.

L'accès aux archives de la série « Publications mathématiques et informatiques de Rennes » 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/>

Classification et corrélation
Application à l'analyse d'un questionnaire parents-élèves
sur l'enseignement des mathématiques

Basavannepa TALLUR
IRISA - Campus de Beaulieu - 35042 RENNES Cédex

Résumé

La notion de similarité ou de distance est fondamentale dans l'analyse classificatoire des données. Cette notion de similarité ou de distance doit prendre en compte d'une façon fidèle la nature des données à analyser par une représentation mathématique appropriée. Les méthodes de classification hiérarchique basées sur la notion de distance et du critère d'"inertie expliquée" utilisent une représentation euclidienne des données avec une métrique adaptée.

La méthode A.V.L. ("Algorithme de la Vraisemblance des Liens") développée par I.C. Lerman (Lerman 1981) utilise plutôt une notion de similarité entre les éléments à classer. La représentation ensembliste des données et une hypothèse de non association conduit à définir une mesure de similarité. Par cette démarche on retrouve la notion de corrélation comme une mesure de similarité entre variables numériques. Nous avons utilisé cette notion de corrélation pour construire un indice de similarité ainsi qu'un indice d'aggrégation pour la classification des lignes et colonnes d'un tableau des fréquences. Après avoir présenté la construction des indices, on exposera les résultats d'une enquête psycho-pédagogique menée auprès des élèves de 6ème d'une part, et des parents des élèves d'autre part afin d'analyser leur perception de l'enseignement des mathématiques. Ces résultats ont été publiés dans le rapport de l'IRISA numéro 177 (Juillet 1982) dont on reprend quelques passages dans cet article.

Le paragraphe 1 introduit la méthode AVL de I.C.Lerman et présente la démarche générale de construction de l'indice de similarité entre variables descriptives. Ensuite une généralisation de cet indice au cas des tableaux de contingence est présentée au paragraphe 2. Les paragraphes 3 et 4 décrivent respectivement l'algorithme directe basé sur la corrélation (A.B.C.) pour la classification des lignes et/ou colonnes d'un tableau de contingence et l'application aux données de l'enquête psycho-pédagogique. La thèse d'état de B.Tallur (1988) traite en détail les problèmes liés à l'analyse exploratoire de tableaux de contingence par la classification.

1 INTRODUCTION

1.1 Schéma de construction de l'indice de similarité

I.C. Lerman (Lerman, 1981, Chapitre 2) décrit le schéma général conduisant à l'expression de l'indice de proximité entre variables dont on peut distinguer les quatre étapes suivantes :

- 1°) Représentation des variables par des structures mathématiques adéquates.
- 2°) Choix d'un indice "brut" de proximité.

3°) Au couple de structures observées à comparer, on associe un couple de structures aléatoires et donc à l'indice "brut" on associe un indice brut aléatoire sous une hypothèse appropriée d'indépendance (ou d'absence de lien) pour construire l'indice centré et réduit.

4°) Indice de proximité probabiliste (ou la "vraisemblance" d'une similarité).

1. Représentation mathématique.

Une variable descriptive est représentée par une partie d'un ensemble ou par une pondération sur ce dernier. Selon que l'ensemble de représentation soit l'ensemble E des individus ou l'ensemble produit ExE, on distingue deux grandes catégories de variables.

On peut classer dans la première catégorie l'attribut descriptif a qui peut être représenté par la partie E des individus qui le possèdent, et la variable numérique X qui définit une pondération sur E en associant à chaque individu i le nombre X(i), valeur de la variable X pour l'individu i.

Toute variable qualitative discrète qui définit une relation binaire sur E est susceptible d'une représentation dans ExE. Ainsi on peut classer dans la deuxième catégorie :

- la variable "rang" qui définit un ordre total θ sur ExE que l'on représente par son graphe.

- la variable qualitative ordinale qui définit un préordre total w sur E que nous représentons par la partie $R(w)$ de ExE définie par $R(w) = \sum \{E_i \times E_j / i < j\}$ (somme ensembliste) où E_i est la $i^{\text{ème}}$ classe du préordre formée des individus possédant la $i^{\text{ème}}$ modalité de la variable.

- la variable qualitative nominale qui définit une partition π sur E que nous représentons dans l'ensemble $F = P_2(E)$ des parties à deux éléments de E par :

$R(\pi) = \sum P_2(E_i)$ (somme ensembliste)

où $P_2(E_i)$ est l'ensemble des paires réunies dans la $i^{\text{ème}}$ classe.

- d'une manière générale, la variable pondération sur ExE peut être représentée par une matrice carrée

$$\{\mu_{xy} / (x, y) \in ExE\}$$

Soit (a,b) un couple de variables définissant le même type de structure sur E et (α, β) le couple de structures définies sur E par (a, b). A α (resp. β) nous associons l'ensemble A (resp. B) des structures sur E de même type et ayant les mêmes caractéristiques cardinales que α (resp. β).

Le couple (a, b) sera représenté par un couple de parties $(R(\alpha), R(\beta))$ de E ou de ExE selon qu'il s'agisse de variables de la 1ère ou de la 2ème catégorie.

2ème étape. Définition d'un indice "brut" de proximité.

On introduit l'indice brut de similarité entre a et b:

$$s = \text{card}(R(\alpha) \cap R(\beta)) \quad \dots(1-1)$$

3ème étape. Hypothèse d'indépendance.

On associe les deux variables aléatoires duales à l'indice brut de proximité s , sous l'hypothèse d'absence de lien (ou d'indépendance) :

$$S_{\alpha} = \text{card}(R(\alpha) \cap R(\beta'))$$

.....(1-2)

$$S_{\beta} = \text{card}(R(\alpha') \cap R(\beta))$$

où α' (resp. β') est un élément aléatoire dans l'ensemble A (resp. B) muni d'une probabilité uniformément répartie. Les v.a. S_{α} et S_{β} ont la même distribution et on montre (Lerman 1981, chapitre 2) sous des conditions assez générales que cette distribution commune est asymptotiquement normal.

indice centré et réduit :

En désignant par $E(S)$ et $\alpha(S)$ la moyenne et la variance de S_{α} (ou S_{β}) on définit l'indice centré et réduit:

$$Q(\alpha, \beta) = \{s - E(S)\} / \alpha(S)$$

.....(1-3)

La réduction globale des similarités:

Pour n "assez grand" l'indice $Q(\alpha, \beta)$ approché par $\sqrt{\rho(\alpha, \beta)}$, où $\rho(\alpha, \beta)$ désigne le coefficient de corrélation entre les variables associées aux structures α et β , prend des valeurs élevées, en valeur absolue, et l'indice $P(\alpha, \beta)$ défini à l'étape 4 ci-dessous ne permet pas une discrimination suffisante.

En réalité, Lerman (1984) montre - pour tenir compte du contexte - que l'h. a. l. doit directement associer à l'ensemble A des variables observées un ensemble A'' de v. a. indépendantes, et non pas seulement à deux d'entre elles sans référence aux autres, et rapporter l'association entre deux variables à l'ensemble des associations mutuelles. Cela peut se faire notamment au moyen d'une réduction globale des similarités:

$$Q(\alpha, \beta) \rightarrow \{Q(\alpha, \beta) - \text{Moy}(Q)\} / \sqrt{\text{Var}(Q)}$$

où $\text{Moy}(Q)$ et $\text{Var}(Q)$ désignent, respectivement, la moyenne et la variance des valeurs de l'indice $Q(\alpha, \beta)$, pour toutes les paires de structures observées.

Sous certaines conditions on montre (cf. Le Calvé, 1976; Lerman, 1984) la normalité asymptotique de l'indice $Q(\alpha, \beta)$ réduit globalement comme ci-dessus.

4. Indice de similarité probabiliste.

L'indice de similarité définitif utilisé dans l'A.V.L. s'écrit :

$$P(\alpha, \beta) = \text{Pr}\{S < s/N\}$$

.....(1-4)

où S est l'une des deux v.a. duales de même loi S_{α} et S_{β} et N désigne l'hypothèse d'absence de lien (ou d'indépendance) exprimée ci-dessus.

Le degré de ressemblance exprimé par l'indice $P(\alpha, \beta)$, entre les variables a et b est d'autant plus grand que la valeur de s est invraisemblablement grande, relativement à l'hypothèse d'indépendance N .

En utilisant le caractère asymptotiquement normal de la variable S , on a l'approximation suivante :

$$P(\alpha, \beta) = \Phi(Q(\alpha, \beta))$$

.....(1-5)

où Φ est la fonction de répartition de la loi normale $N(0, 1)$

étape 1 : représentation.

Les variables numériques, comme les attributs descriptifs appartiennent à la première catégorie de variables à savoir celles qui peuvent être représentées au niveau de l'ensemble E. La variable numérique est représentée par une mesure sur E alors qu'un attribut est représenté par une partie de E.

Une variable numérique X est définie par sa distribution (X_1, X_2, \dots, X_n) sur l'ensemble d'objets E, où X_i est la mesure de la variable X sur l'objet i.

étape 2 : indice brut.

La base de construction de l'indice de similarité entre deux variables numériques X et Y dont les distributions sur E sont respectivement :

(X_1, X_2, \dots, X_n) et (Y_1, Y_2, \dots, Y_n)

sera

$$s_{XY} = \sum \{X_i Y_i / 1 \leq i \leq n\} \quad \dots\dots(1-6)$$

dont on examinera la distribution sous l'hypothèse d'indépendance.

étape 3 : Sous l'hypothèse d'indépendance entre X et Y, les distributions de ces variables sont fixées mais la position relative de l'une par rapport à l'autre est inconnue.

Les deux variables aléatoires duales associées à s_{XY} sous l'hypothèse d'indépendance sont :

$$S_X = \sum \{X_i Y_{\alpha(i)} / 1 \leq i \leq n\} \quad \text{et} \quad S_Y = \sum \{X_{\alpha(i)} Y_i / 1 \leq i \leq n\} \quad \dots\dots(1-7)$$

où $(\sigma(1), \dots, \sigma(n))$ est un élément aléatoire pris dans l'ensemble de toutes les permutations de $(1, 2, \dots, n)$.

La distribution de S_X est la même que celle de S_Y .

On trouve l'espérance mathématique et la variance de la variable S_X :

$$E[S(X)] = n \mu_X \cdot \mu_Y \quad \text{et} \quad V[S(X)] = (n^2 / (n-1)) \sigma_X^2 \sigma_Y^2 \quad \dots\dots(1-8)$$

où μ_X (resp μ_Y) et σ_X^2 (resp. σ_Y^2) sont la moyenne et la variance de la variable X (resp. Y).

L'indice centré et réduit de proximité entre X et Y s'écrit :

$$Q(X, Y) = (s_{XY} - n \mu_X \mu_Y) / \sqrt{((n^2 / (n-1)) \sigma_X^2 \sigma_Y^2)} \\ = \sqrt{(n-1)} \rho(X, Y) \quad \dots\dots(1-9)$$

où $\rho(X, Y)$ est le coefficient de corrélation linéaire entre X et Y.

Le caractère asymptotiquement normal de la variable $Q(X, Y)$ est démontré sous des conditions assez générales dans le théorème de A.Wald et J. Wolfowitz (1944).

étape 4 : Indice de proximité entre X et Y sera approché par:

$$P(X, Y) = \phi(Q(X, Y))$$

REMARQUE. Nous ne nous étendrons pas ici, sur le problème de construction d'un indice de similarité entre objets ou individus qui a été abordé dans Lerman et Peter (1985) et Lerman (1987).

2 CAS D'UN TABLEAU DE CONTINGENCE

2.1. REPRESENTATION GEOMETRIQUE D'UN TABLEAU DE CONTINGENCE.

Considérons le tableau de contingence

$$K_{I,J} = \{k_{ij} / (i, j) \in I \times J\}$$

----- (2-1)

où les k_{ij} sont des cardinaux des classes $C_i \cap C'_j$ du croisement de deux partitions $\{C_1, C_2, \dots, C_n\}$ et $\{C'_1, C'_2, \dots, C'_m\}$ respectivement en $n = \text{card}(I)$ et $m = \text{card}(J)$ classes.

Pour représenter géométriquement ce tableau des données, on fait jouer à l'un des deux ensembles I ou J, le rôle de l'ensemble des objets et à l'autre le rôle de l'ensemble des variables.

Supposons, pour fixer les idées, que I joue le rôle de l'ensemble des individus et J celui des variables. Nous utiliserons la même représentation de I à travers J, que celle qui est adoptée dans le cadre de l'analyse des correspondances.

Nous noterons, pour tout $(i, j) \in I \times J$,

$$k_{i.} = \sum \{k_{ij} / j \in J\}, \quad k_{.j} = \sum \{k_{ij} / i \in I\}$$

$$k_{..} = \sum \{k_{ij} / (i, j) \in I \times J\}$$

$$f_{ij} = k_{ij} / k_{..}, \quad f_{i.} = k_{i.} / k_{..}, \quad f_{.j} = k_{.j} / k_{..}$$

.....(2-2)

$$f^i_j = (f^i_j / j \in J) \text{ où } f^i_j = f_{ij} / f_{i.}$$

f^i_j : "profil de i à travers J"

On associe à l'ensemble I d'objets le nuage $N(I)$ dans R^m muni de la métrique diagonale $\{(1/f_{.j}) / j \in J\}$:

$$N(I) = \{(f^i_j, f_{i.}) / i \in I\}$$

.....(2-3)

où $f_{i.}$ est le poids affecté au point f^i_j , pour tout $i \in I$.

J.L. Philoche a proposé une représentation plus symétrique (journées de statistique, Paris Juin 1979) de l'ensemble I par le nuage $M(I)$ dans R^m :

$$M(I) = \{(l_{ij}, f_{i.}) / i \in I\}$$

.....(2-4)

ou $l_{ij} = (i_j / j \in J)$, $i_j = f_{ij} / f_{i.}$;

l'objet i de I est représenté par le point l_{ij} de R^m .

L'analyse du nuage $M(I)$ dans R^m , muni de la métrique diagonale $(f_{.j} / j \in J)$ est équivalente à l'analyse du nuage $N(I)$ muni de la métrique diagonale $\{(1/f_{.j}) / j \in J\}$.

Nous verrons que, quelle que soit la représentation géométrique retenue - la représentation classique ou 'symétrique' - nous aboutissons aux mêmes indices de similarités.

Considérons pour fixer les idées, le problème de la classification de l'ensemble J . On peut noter une *différence essentielle* entre notre méthode et la méthode ascendante de classification hiérarchique de J.P. Benzecri et M. Jambu pour ce type de données. Cette dernière est conçue à partir de la représentation de l'ensemble de points-objets J dans R^n et utilise les concepts de distance au carré et d'inertie expliquée au sens de la métrique du χ^2 , tandis que notre méthode se conçoit au niveau de $(R^m)^*$ et utilise le concept de la corrélation, conformément à la représentation définie plus haut.

2.2. INDICE DE PROXIMITE ENTRE COLONNES OU LIGNES D'UN TABLEAU DE CONTINGENCE.

Considérons pour fixer les idées le problème de la classification de l'ensemble J de colonnes. J jouera donc le rôle de variables descriptives et nous allons considérer la représentation de I à travers J définie par le nuage $N(I)$. On va associer à la colonne j , $1 \leq j \leq m$, une variable numérique X_j dont la mesure sur l'individu i , $1 \leq i \leq n$, vaut f_{ij} . La distribution de la variable X_j ($1 \leq j \leq m$) est donnée par la suite des valeurs prises :

$$\{f_{1j}^1, f_{1j}^2, \dots, f_{1j}^i, \dots, f_{1j}^n\}$$

avec les fréquences relatives :

$$\{f_{1.}, f_{2.}, \dots, f_{n.}\}$$

L'indice de proximité entre les éléments j et h de J n'est autre que le coefficient de corrélation $\rho(j, h)$ entre les variables numériques X_j et X_h dont nous allons préciser l'expression.

En effet, la moyenne et la variance de la variable X_j sont respectivement :

$$\mu_j = \sum \{f_{i.} f_{ij} / i \in I\} = f_{.j}$$

$$\sigma_j^2 = \sum \{f_{i.} (f_{ij}^2 - f_{ij}^2) / i \in I\} = \sum \{(f_{ij}^2 / f_{i.}) / i \in I\} - f_{.j}^2 \quad \dots\dots(2-5)$$

D'autre part, la covariance entre X_j et X_h est donnée par :

$$\begin{aligned} \sigma_{jh} &= \sum \{f_{i.} (f_{ij} - f_{.j})(f_{ih} - f_{.h}) / i \in I\} \\ &= \sum \{(f_{ij}f_{ih} / f_{i.}) / i \in I\} - f_{.j}f_{.h} \quad \dots\dots(2-6) \end{aligned}$$

On obtient finalement :

$$\rho(j, h) = \sigma_{jh} / (\sigma_j \sigma_h)$$

$$\rho(j, h) = \frac{\sum\{(f_{ij}f_{ih}/f_{i.}) / i \in I\} - f_{.j}f_{.h}}{[\{\sum\{(f_{ij}^2/f_{i.}) / i \in I\} - f_{.j}^2\}\{\sum\{(f_{ih}^2/f_{i.}) / i \in I\} - f_{.h}^2\}]}^{1/2} \quad \dots\dots(2-7)$$

Cet indice dépend de la représentation euclidienne de I à travers J. Avant d'adopter définitivement l'indice (2-7) ci-dessus, nous avons expérimenté d'autres indices tels que le coefficient de corrélation entre les colonnes et la covariance entre les variables associées qui n'ont pas donné des résultats satisfaisants. On vérifie aisément les deux propriétés suivantes de cet indice :

Propriété 1 : L'indice de proximité $\rho(j, h)$ défini ci-dessus est invariant par rapport à la métrique dont peut être muni R^m pour l'évaluation des distances entre éléments de I.

Propriété 2 : La représentation euclidienne de I à travers J au moyen du nuage M(I) conduit au même indice $\rho(j, h)$ de proximité entre les éléments j et h de J, que (2-7) ci-dessus.

Pour obtenir l'indice de proximité entre éléments de I on considère, de façon analogue la représentation euclidienne de J à travers I au moyen du nuage N(J) ou M(J).

A partir de (2-7) on peut directement écrire l'expression de l'indice de proximité entre les lignes i et i' de I en intervertissant les rôles de I et de J :

$$\rho(i, i') = \frac{\sum\{(f_{ij}f_{i'j}/f_{.j}) / j \in J\} - f_{i.}f_{i' .}}{[\{\sum\{(f_{ij}^2/f_{.j}) / j \in J\} - f_{i.}^2\}\{\sum\{(f_{i'j}^2/f_{.j}) / j \in J\} - f_{i' .}^2\}]}^{1/2} \quad ; 1 \leq i, i' \leq n. \quad \dots\dots(2-8)$$

Les valeurs de l'indice (2-7) (resp. (2-8)) entre colonnes (resp. lignes) joueront le même rôle que celles de l'indice "centré et réduit" Q(a, b) défini dans le paragraphe 1.1 du chapitre 1.

Au lieu de considérer la représentation de l'ensemble I à travers J par le nuage N(I) dans R^m , si on considère la représentation de J à travers I par le nuage N(J) dans R^n muni de la métrique du χ^2 , on peut associer à la colonne j un "individu" défini par son profil $f_{.j}$, ($1 \leq j \leq m$). Lerman et Peter (1985) ont montré que l'indice de similarité entre les colonnes j et h obtenu en tant que le coefficient de corrélation entre les variables X_j et X_h définies ci-dessus est identique à celui défini par le cosinus de l'angle entre les vecteurs $f_{.j}$ et $f_{.h}$ par rapport au centre de gravité du nuage $g_I = \{f_i / i \in I\}$.

Propriété 3 : L'indice de proximité (2-7) (resp. (2-8)) entre colonnes (resp. lignes) vérifie la propriété d'équivalence distributionnelle :

En remplaçant les colonnes (resp. lignes) de "profils" identiques par leur somme, les valeurs de l'indice de proximité entre lignes (resp. colonnes) ne sont pas modifiées.

3 ALGORITHME BASÉ SUR LA CORRÉLATION: A.B.C.

Nous avons suggéré (cf Lerman-Tallur, 1980) une façon d'associer un arbre de classification directement à l'indice de proximité que nous venons d'exposer dans les précédents paragraphes de ce chapitre. La structure des données qui nous intéresse ici est celle qu'on peut assimiler à un tableau de contingence ou une juxtaposition, selon une ou deux directions, des tableaux de contingence.

La méthode que nous proposons ici est une méthode naturelle de construction de l'arbre de classification hiérarchique. Elle est comparable à des méthodes basées sur les carrés des distances entre les centres de gravité des classes, mais au lieu de minimiser un critère basé sur la distance, on maximise un critère basé sur la corrélation.

Soit le tableau de contingence (2-1) :

$K_{IJ} = \{k_{ij} / (i, j) \in I \times J\}$
avec $\text{card}(I) = n$ et $\text{card}(J) = m$

Supposons, pour fixer les idées, que l'ensemble à classifier soit celui des colonnes. Nous avons vu que, pour la définition de l'indice de proximité entre les éléments de J , nous avons associé à chaque élément $j \in J$ une variable numérique X_j en passant par une représentation géométrique du tableau (c'est-à-dire, en nous plaçant dans le nuage $N(I)$ des profils des lignes). Le coefficient de corrélation $\rho(j, j')$ entre les variables X_j et $X_{j'}$ associées respectivement aux colonnes j et j' de J est, à un coefficient multiplicatif près, l'indice centré et réduit $Q(j, j')$. Alors que l'A.V.L. se réfère à une échelle de probabilité fournie par la fonction de répartition $\Phi(Q(j, j'))$ de la loi $N(0,1)$, la méthode proposée ici se réfère directement à l'indice de corrélation.

3.1. Algorithme basé sur la corrélation (A.B.C.)

Le principe de l'A.B.C. consiste à agréger à chaque pas les deux éléments ou deux classes d'éléments les plus proches au sens de l'indice d'agrégation défini de la façon suivante.

Définition. Soient C et D deux classes formées à un certain niveau (pas) de l'algorithme où C est formée par la réunion de p colonnes C_1, C_2, \dots, C_p et D est formée par la réunion de q colonnes D_1, D_2, \dots, D_q . Notons X_{C_i} ($1 \leq i \leq p$) (respectivement X_{D_j} ($1 \leq j \leq q$)) la variable numérique associée à la colonne C_i (resp. D_j), $X_C = \sum\{X_{C_i} / 1 \leq i \leq p\}$ et $X_D = \sum\{X_{D_j} / 1 \leq j \leq q\}$. Alors, l'indice d'agrégation S_{CD} entre les classes C et D sera

$$S_{CD} = \rho(X_C, X_D) / \sqrt{pq} \quad \text{--- (3-1)}$$

où ρ désigne le coefficient de corrélation.

On peut résumer l'algorithme de la façon suivante :

Pas 1. Au départ chaque colonne constitue une classe. On réunit les deux colonnes les plus proches au sens de l'indice d'agrégation ($s-1$) (c-à-d. celles qui réalisent le maximum de cet indice) avec $p = q = 1$:

$$S_{jh} = \rho(X_j, X_h) \text{ pour tout } (j, h) \in J \times J, j \neq h.$$

Pas 2. On associe à la classe formée au pas précédent, la somme des variables associées aux colonnes constituant cette classe. Ainsi, la variable associée à la classe C formée des colonnes C_i ($1 \leq i \leq p$) sera la variable X_C somme des variables X_{C_i} ($1 \leq i \leq p$). Remarquons que cela revient à remplacer les colonnes réunies C_i ($1 \leq i \leq p$) par leur somme et à associer à cette colonne somme une variable numérique conformément au paragraphe 2.2.

Pas 3. On réunit les classes les plus proches au sens de l'indice d'agrégation ($s-1$) (c-à-d. pour lesquelles cet indice est maximum) et on retourne au pas 2 jusqu'à ce que toutes les classes soient réunies.

Définition. La hiérarchie de partitions sera dite "sans inversion" si la valeur maximale de l'indice d'agrégation fondé sur les similarités à un pas s quelconque est inférieure à la valeur maximale de cet indice à tous les pas inférieurs à s . Autrement dit, l'algorithme est sans inversions si la suite des valeurs maximales de l'indice d'agrégation aux niveaux successifs est strictement décroissante, lorsque l'indice d'agrégation est basé sur les similarités.

4. APPLICATION AUX TABLEAUX BINAIRES

Le domaine très vaste d'application de la méthode de l'Analyse des correspondances, qui, on le sait, a été fondamentalement développé pour le traitement des tableaux de contingence montre que les tableaux binaires (présence-absence codées 1-0) peuvent être assimilés à ces derniers pour analyser les nuages des profils. Inspiré par ceci, nous avons essayé notre méthode de classification sur des données de type 0-1 provenant d'une enquête psycho-pédagogique en mathématiques. Les résultats de cette analyse, forts intéressants, sont confrontés à ceux de l'Analyse des correspondances d'une part et à ceux de la classification par l'A.V.L. de l'autre. S'agissant d'éprouver l'efficacité et l'opportunité d'un nouvel outil, il est bien nécessaire de comparer sa performance avec quelques outils standards bien connus. Nous exposerons les résultats de cette application de manière détaillée dans les paragraphes suivants.

4 APPLICATION

4.1 DIDACTIQUE DES MATHÉMATIQUES

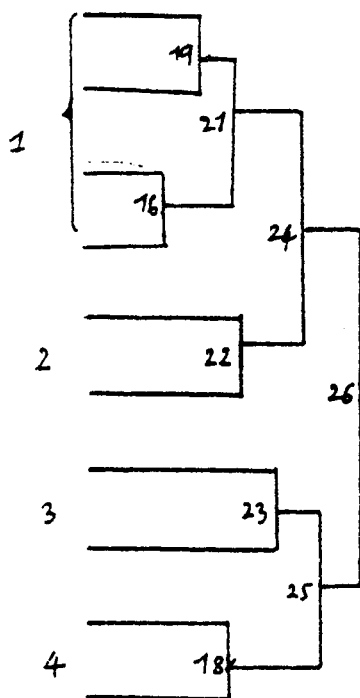
On présentera ci-dessous une application de la méthode de classification exposée dans le paragraphe précédent à une enquête psychopédagogique. Le but de cette enquête était d'analyser de quelle manière l'enseignement des mathématiques en classe de 6^{ème} est perçu par les enfants eux mêmes, les enseignants et les parents.

Une première enquête menée par Jacques Degouys (Université de Haute Bretagne), Régis Gras (IRMAR) et Marcel Postic (Université de Haute Bretagne) en 1978 avait permis l'établissement d'une échelle d'attitude.

L'étude actuelle porte sur un échantillon de 163 élèves. Le questionnaire se compose de 3 parties : un ensemble de questions destinées à l'enfant, un ensemble de questions aux parents et une troisième partie destinée à l'enseignant. Chaque question a deux réponses possibles : oui et non. Nous présenterons ci-dessous l'analyse par classification hiérarchique, selon l'algorithme basé sur la corrélation portant sur deux parties distinctes du questionnaire : réponses des enfants et réponses des parents.

4.1.1 RÉPONSES DES ENFANTS :

Il a été retenu un ensemble de 27 questions auxquelles les réponses sont données par oui ou par non. (voir la liste des questions en annexe). La figure 1 représente l'arbre, condensé à ses noeuds significatifs, de la classification hiérarchique des questions (ou plutôt, il s'agit des attributs définis par des réponses positives aux questions).



Au niveau 22 la statistique globale atteint le maximum et on distingue quatre grandes classes. (voir schéma ci-contre).

Classe 1 - Cette classe regroupe l'ensemble des opinions négatives à l'égard des mathématiques. Cet ensemble des opinions négatives reste marqué par celles de la famille du fait de la jeunesse des enfants. Notons que le trait dominant est l'extériorité des mathématiques par rapport à l'enfant. C'est en situation d'échec que l'enfant se prononce ici. Les mathématiques

Figure 1. Classification des réponses des enfants par la méthode A.B.C.

| | | | |
|---|----------------------|--------------------|--|
| | Math iep famille | >---* | |
| | Situation d'echec | > 6 | |
| | Explication camarade | >---I-----* | |
| | Reussite autres mati | > 2 I I | |
| | Peur d'interrogation | >---* 9-----* | |
| | Sentiment d'etre etr | > I I | |
| 1 | Monde sans math | >-----* 15-----* | |
| | Mauvaise note peur | > I I | |
| | Langage difficile | >-----* 11-----* | |
| | Reussite vie*bon mat | >-----* 16-----* | |
| | Cours*exercice*Compr | > I I *24 | |
| | Utilite aide compreh | >-----* 17-----* | |
| 2 | Confiance *comprehen | >-----* I I | |
| | Meme inelligent echo | >-----* 12-----* | |
| | Math pas Intelligenc | >-----* I I | |
| | Parents aident exo | >-----* 8-----* | |
| | tous aimer math | >-----* I I I--- | |
| 3 | prof copies malaise | >-----* 23-----* | |
| | Exemple*comprehensio | >-----* 10-----* | |
| | Peur de se tromper | >-----* 13-----* | |
| | Envie chercher apres | >---* I I | |
| | Prefere Exercices | >---* 5-----* | |
| | Aime chercher soluti | >-----* I I | |
| 4 | Content du travail | >-----* 25-----* | |
| | Plasir a faire math | > 1 I I | |
| | Langage math exprime | >---I-----* | |
| | Control parents | > *4 I I | |
| | | >---* 7-----* | |
| | | > I I I | |
| | | >---* I I I | |
| | | > 3-----* 14-----* | |
| | | >---* I I I | |
| | | > I *18-----* | |
| | | >-----* I | |
| | | > I | |
| | | >-----* | |

constituent un monde à part, extérieur

à celui de l'enfant. C'est ainsi que d'autres que lui (famille, professeur, camarade) peuvent tenir entre ces mondes des rôles médiateurs. L'enfant ici a des difficultés. La réussite sociale par les mathématiques, l'ésotérisme de leur langage représentent les grands traits du mythe mathématique.

Classe_2 - Cette fois, dans l'ensemble des opinions plutôt tournées vers une représentation positive des mathématiques, la forme générale des opinions est "subordination". Si les maths ou leur enseignement ont telle propriété, alors tout est versé en faveur de la compréhension. Celle-ci, en effet, est successivement subordonnée à la qualité du cours et des exercices, à la confiance que l'on ressent, à l'utilité des mathématiques elles-mêmes.

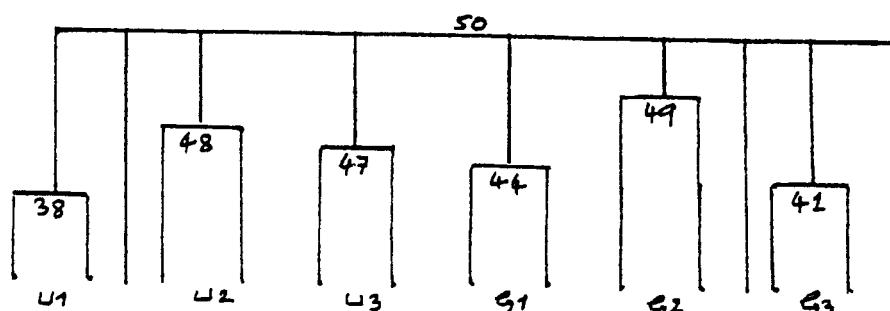
Classe_3 - Cette fois l'élève consent à s'engager sur cette terre où les maths sont apparemment terrifiantes : "tout le monde pourrait aimer les maths" ; il ne s'exclut pas a priori. Pourtant, dans cet engagement, le sentiment de crainte, le coût de l'engagement accompagnent le moment de distribution des copies et celui de prise de parole.

Classe_4 - Cette fois aucune restriction ne vient troubler les opinions très favorables aux mathématiques. L'expression de plaisir apparaît 3 fois dans ces opinions ; de plus la forme personnelle ("je", "mes"...) employée 6 fois sur 7 montre à quel point, l'enfant s'implique dans la réussite, même si ses parents contrôlent ses résultats.

4.1.2 RÉPONSES DES PARENTS

L'ensemble des 58 réponses des parents a été classifié par l'algorithme présenté dans cet article.

Les six grandes classes se situent l'une par rapport à l'autre de la manière suivante :



INTERPRÉTATION (voir figure 2 : l'arbre de classification).

Contrairement à l'enquête faite auprès des enfants, ce n'est pas la réussite en maths qui établit une partition dans l'ensemble des opinions des parents. Elle n'est pas ressentie par eux, comme discriminatrice dans le champ de l'attitude. Par contre, c'est plutôt l'utilité des disciplines, donc les racines de l'avenir et de la réussite sociale qui est en jeu.

Classe U₁ Les mathématiques et les Sciences naturelles constituent les disciplines les plus utiles pour l'avenir. Les résultats étant jugés bons, la meilleure voie de réalisation de ce projet passe par la filière C. Les parents projettent peut-être leur foi dans l'avenir social, sur l'intérêt que l'enfant porterait aux mathématiques, les langues étant jugées comme une gêne dans le projet. Notons que cette classe se renforce à l'ultime niveau, du dernier élément du tryptique scientifique : la physique, dont on sait qu'elle constitue un pilier de la filière C.

Classe U₂ Les résultats en mathématiques sont, cette fois, moyens. Mais refusant de responsabiliser l'enfant, et, peut-être, eux-mêmes pour des raisons héréditaires, le maître et son enseignant sont culpabilisés. La voie "langues vivantes" est, dans ce cas, déclarée la plus utile pour l'avenir.

Classe U₃ Cette fois la planche de salut social est le français ; mais cet intérêt semble l'alibi contre le pouvoir des mathématiques qui se voient accusées des maux pernicioseux : excès de travail, trop grande place, difficulté etc...échec, bien évidemment.

Classe G₁ Cette classe, constituée surtout d'absence de réponses, est l'image d'une population détachée pour des raisons culturelles peut-être, de la préoccupation d'une liaison école-société. Comment sans culpabilisation excessive, pourrait-on subordonner la réussite sociale à la réussite scolaire lorsque, en particulier, l'enfant est très faible en mathématiques ?

Classe G₂ Ici par contre, la réussite est excellente, mais comme naturelle, prise comme héréditaire, les parents aimaient les mathématiques et réussissaient dans cette discipline. A quoi bon dans ce cas où la sérénité sociale a sa place, où l'on peut aimer les mathématiques et les sciences de façon gratuite, discriminer les filles des garçons : l'enseignement des maths a la même importance pour l'un que pour l'autre. Par contre le français semble occuper trop de place dans l'enseignement.

Figure 2. Classification des réponses des parents par la méthode A.B.C.

| | | | |
|----|--------------------------|------|----|
| | maths + utile | > 3 | |
| | langues + travail | > 2 | |
| | langues + difficile | > 19 | |
| | maths + interessant | > 27 | |
| | souhait classe c | > 6 | |
| U1 | capable classe c | > 28 | |
| | bon en maths | > 8 | |
| | sans rep. mauvais res. | > 12 | |
| | sc. nat. + grande place | > 12 | 50 |
| | sc. physiques + travail | > 12 | |
| | sc. physiques + utiles | > 12 | |
| | langues viv. + utile | > 18 | 50 |
| | arts + div. + gde place | > 32 | |
| | n'aime pas prof. | > 44 | |
| | moyen en maths | > 21 | |
| U2 | enseign. inadquat | > 21 | |
| | sc. nat. + utiles | > 43 | 48 |
| | langues v. + interessant | > 45 | |
| | arts + div. + utiles | > 33 | |
| | sc. nat. + travail | > 4 | |
| | sc. nat. + difficile | > 50 | |
| | francais + utile | > 29 | |
| | langues + grande place | > 29 | |
| | francais + interessant | > 9 | 47 |
| | maths + grande place | > 36 | |
| U3 | divers mauvais res. | > 22 | |
| | maths + travail | > 24 | |
| | maths + difficile | > 17 | |
| | faible en maths | > 1 | |
| | sans reponse + utile | > 14 | |
| | sans rep. + difficile | > 15 | |
| | sans rep. + interessant | > 7 | 50 |
| | sans rep. + travail | > 16 | |
| | sans rep. + grande place | > 13 | |
| | sans rep capable c | > 40 | |
| g1 | arts + div. + interess. | > 26 | |
| | sans rep. sexe | > 31 | |
| | arts + div. + difficile | > 30 | 64 |
| | tres faible en maths | > 9 | |
| | difficulte contenu | > 20 | |
| | sans rep. classe c | > 1 | |
| | garcon > fille | > 1 | 50 |
| | sc. physiques + interes. | > 35 | |
| | francais + grande place | > 23 | |
| | tres bon en maths | > 34 | |
| | garcon = fille | > 37 | |
| g2 | francais + travail | > 10 | 42 |
| | francais + difficile | > 10 | |
| | sc. nat. + interessantes | > 16 | 49 |
| | parents aiment maths. | > 1 | |
| | parents bons maths | > 1 | 50 |
| | sc. phys. + grande place | > 11 | |
| | arts + div. + travail | > 11 | 50 |
| | sc. physique + difficile | > 39 | |
| g3 | aide maison | > 25 | 41 |
| | trop tot classe c | > 1 | 50 |
| | insuffisance travail | > 1 | |
| | fille > garcon | > 1 | |

Classe G₃ Quelques opinions peu structurées se retrouvent ici.
C'est une sous-classe des incertitudes où l'enfant recevant une aide à la maison ne manifeste pas ses aptitudes, avec une grande fiabilité.

RECONNAISSANCE

Tous mes remerciements à Régis GRAS de l'aide indispensable qu'il m'a apportée dans l'interprétation des résultats de l'enquête psychopédagogique. C'est lui qui m'a rendu les données accessibles, et les conclusions sont de lui.

5 Annexe

5.1 Questionnaire enfants

| CODE ATTRIBUT | SA SIGNIFICATION |
|------------------------------|--|
| Math imp famille | Si je travaille en maths c'est à cause de l'importance qu'on lui donne dans ma famille |
| Math pas intelligence | Être doué en maths n'est pas une affaire d'intelligence |
| Cours+exercice=comprehension | Quand on comprend les cours et qu'on revoit les exercices on est sûr de comprendre les maths |
| Situation d'échec | En maths on est toujours en situation d'échec |
| Mauvaise note peur | Lorsque j'ai une mauvaise note en maths je n'ose pas le dire à la maison |
| Tous aimer math | Tout le monde pourrait aimer les maths |
| Monde sans math | On peut imaginer un monde sans mathématiques |
| Langage difficile | Ce qui est le plus difficile en mathématiques, c'est le langage |
| Prof copies malaise | On n'est jamais à l'aise quand le prof rend les copies de maths |
| Envie chercher après cours | Après les cours de maths, j'ai envie de continuer à chercher |
| Explication camarade | Je préfère demander une explication à un camarade plutôt qu'à un prof |
| Exemple+compréhension | On est sûr de comprendre quand le prof nous donne un exemple |
| Réussite vie + bon en maths | Pour réussir dans la vie il faut être bon en maths |
| Contrôle parents | Mes parents contrôlent souvent mes résultats en mathématiques |
| Réussite autres matières | On n'a pas besoin de maths si on réussit dans d'autres matières |
| Aime chercher solution | J'aime chercher la solution des problèmes des maths |
| Peur d'interrogation | Lorsqu je sais que je vais être interrogé en maths je n'ai pas envie d'aller en cours |
| Confiance + compréhension | Pour comprendre les maths, il faut se sentir en confiance |
| Peur de se tromper | Quand on répond en math, on a peur de se tromper |
| Préfère exercices | Ce que je préfère en maths, c'est faire des exercices |
| Parents aident en exo | Mes parents m'aident parfois à faire mes exercices de maths |
| Content du travail | Je suis content au travail que je fais en mathématiques |
| Même intelligent échoue | Même s'ils sont intelligents, certains élèves peuvent échouer en mathématiques |
| Langage math exprimer claire | Avec le langage mathématique on peut s'exprimer clairement |
| Sentiment d'être étranger | En cours de mathématiques, on a le sentiment d'être étranger |
| Plaisir à faire math | J'éprouve du plaisir à faire des mathématiques |
| Utilité aide compréhension | On ne peut comprendre les mathématiques que si on voit à quoi elles servent |

5.2 Questionnaire parents

| Question | Réponse | Code Attribut |
|---|---|---|
| Quelle matière vous semble la plus utile pour l'avenir de votre enfant? | Mathématiques Langues vivantes Français Sciences physiques Sciences naturelles Arts + divers Sans reponse | Maths+ utile Langues viv + utile Français + utile Sc. physiques + utiles Sc. nat. + utiles Arts + div. + utiles Sans reponse + utile |
| Quelle matière semble la plus intéressante? | Mathématiques Langues vivantes Français Sciences physiques Sciences naturelles Arts + divers Sans reponse | Maths + intéressant Langues v. + intéressant Français + intéressant Sc. physiques + intéressant Sc. nat. + intéressant Arts + div. + intéressant Sans reponse + intéressant |
| Quel est l'enseignement qui vous paraît occuper la plus grande place en 6ème? | Mathématiques Langues vivantes Français Sciences physiques Sciences naturelles Arts + divers Sans reponse | Maths + grande place Langues viv + grande place Français + grande place Sc. physiques + grande place Sc. nat. + grande place Arts + div. + grande place Sans reponse + grande place |
| Quelle est la matière qui, cette année lui donne le plus de travail? | Mathématiques Langues vivantes Français Sciences physiques Sciences naturelles Arts + divers Sans reponse | Maths + travail Langues viv + travail Français + travail Sc. physiques + travail Sc. nat. + travail Arts + div. + travail Sans reponse + travail |
| Quelle est la matière qui semble présenter le plus de difficultés à votre enfant en 6ème? | Mathématiques Langues vivantes Français Sciences physiques Sciences naturelles Arts + divers Sans reponse | Maths + difficile Langues viv + difficile Français + difficile Sc. physiques + difficile Sc. nat. + difficile Arts + div. + difficile Sans reponse + difficile |
| Comment estimez-vous l'importance de l'enseignement des mathématiques? | Plus important pour les garçons que pour les filles Aussi important pour les garçons que pour les filles Moins important pour les garçons que pour les filles sans opinion | garçon > fille garçon = fille fille > garçon sans rep sexe |
| Comment jugez-vous votre enfant en mathématiques? | très bon en mathématiques bon en mathématiques Moyen en mathématiques Faible en mathématiques Très faible en mathématiques | très bon en maths bon en maths moyen en maths faible en maths très faible en maths |

D'autres attributs du questionnaire parents

| ATTRIBUT | CODE ATTRIBUT |
|---|--|
| Si notre enfant réussit mal en mathématiques c'est parce que : Il ne travaille pas suffisamment La méthode d'enseignement ne lui convient pas le contenu des cours trop difficile pour lui Il n'aime pas l'enseignant de maths cette année Raisons diverses, mauvais résultats Je ne sais pas | insuffisance travail enseign inadéquat difficulté n'aime pas prof divers mauvais rés sans rep mauvais rés |
| L'enfant demande de l'aider dans son travail en maths à la maison | aide maison |
| J'aimais les maths quand j'étais élève | parents aiment maths |
| Je réussissais en maths | parents bon maths |
| Je souhaite que notre enfant suive une classe de type C | souhait classe C |
| Je ne sais pas s'il doit suivre une classe C | sans rep classe C |
| Je pense qu'il est capable de suivre une classe C | capable classe C |
| Je ne sais pas s'il est capable de suivre une classe C | sans rep capable C |
| Il est trop tôt pour savoir s'il est capable de suivre une classe C | trop tôt classe C |

6 BIBLIOGRAPHIE

Lerman I.C., Tallur B. (1980): Classification des éléments constitutifs d'une juxtaposition de tableaux de contingence *Revue de Statistique Appliquée, vol. XXVIII, no. 3*

Lerman I.C. (1981): Classification et analyse ordinaire des données. *Dunod, Paris*

Lerman I.C., Peter P. (1985): Élaboration et logiciel d'un indice de similarité entre objets d'un type quelconque. Application au problème de consensus en classification. *Publication interne n0. 262, IRISA, Université de Rennes 1*

Lerman I.C. (1987): Construction d'un indice de similarité entre objets décrits par des variables d'un type quelconque. Application au problème de consensus en classification. *Revue de Statistique Appliquée, vol. XXXV, no. 2*

Tallur B. (1982): Un nouvel algorithme de classification hiérarchique des éléments constitutifs de tableau de contingence basé sur la corrélation. *Publication interne n0. 177, IRISA, Université de Rennes 1*

Tallur B. (1988): Contribution à l'analyse exploratoire de tableaux de contingence par la classification. *Thèse d'état, Université de Rennes 1*

Wald A., Wolfowitz J. (1944): Statistical tests based on permutations of the observations. *Ann. Math. Stat., vol. 15*