

RÉGIS GRAS

ANNIE LARHER

**L'implication statistique : une nouvelle méthode d'analyse  
de données didactiques**

*Publications de l'Institut de recherche mathématiques de Rennes, 1990-1991, fascicule 5  
« Didactique des mathématiques », , exp. n° 4, p. 1-30*

[http://www.numdam.org/item?id=PSMIR\\_1990-1991\\_\\_5\\_A4\\_0](http://www.numdam.org/item?id=PSMIR_1990-1991__5_A4_0)

© Département de mathématiques et informatique, université de Rennes,  
1990-1991, 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/>

**L'IMPLICATION STATISTIQUE :  
UNE NOUVELLE METHODE D'ANALYSE DE DONNEES DIDACTIQUES**

**Régis GRAS et Annie LARHER**

Laboratoire de Didactique de l'IRMAR

**Résumé.** Si le problème de la concomitance de deux évènements  $a$  et  $b$  trouve une partie de sa réponse dans l'étude symétrique de la corrélation ou dans celle de la similarité, celui de l'implication (si  $a$  alors  $b$ ) passe, en revanche, par l'examen d'une relation dissymétrique. Par rapport à une problématique psychologique de complexité, R. GRAS, dans sa thèse, a apporté une contribution qui a permis de nombreuses applications de ce type de relation dans des travaux de recherche en psychologie génétique et en didactique des mathématiques, domaines non exclusifs d'autres champs d'application. Mais les variables considérées dans sa recherche se limitent aux variables binaires, présence-absence d'un caractère chez un individu donné. Il s'agit ici d'étendre l'étude de l'implication statistique (ou quasi-implication) à d'autres types de variables et, surtout, à des classes de telles variables. Cette extension nous permet de construire un arbre de classes orientées. Ces développements et cette construction originale résultent, principalement, de la thèse d'Annie LARHER.

**INTRODUCTION.** La notion de similarité entre attributs  $a$  et  $b$  - ou variables binaires - est essentiellement symétrique. Mais s'agissant de répondre à l'étude de "si  $a$  alors  $b$ ", R. GRAS [GRAS 1979] puis I.C. LERMAN, R. GRAS, H. ROSTAM [LERMAN, GRAS, ROSTAM, 1981] ont défini et précisé la notion de quasi-implication mesurée par une intensité d'implication et la notion de graphe d'implication, image de la relation de préordre partiel qui en découle. L'intensité d'implication s'exprime, dans le cas général, par une intégrale gaussienne, qui rend compte du caractère invraisemblable de l'observation faite du nombre de cas où l'implication  $a \Rightarrow b$  se trouve défailante, dans une hypothèse d'absence de lien implicatif a priori. Des propriétés de cette intensité, rappelées ou étudiées dans la thèse d'A. LARHER [LARHER, 1991], sous la direction de R. GRAS, seront présentées ici. Mais l'originalité de cette thèse et des travaux menés depuis consiste principalement, toujours de façon dissymétrique :

(1) L'essentiel de ce texte est soumis à la revue : "Mathématiques, Informatique et Sciences Humaines".

- d'une part, dans l'extension à d'autres types de variables (modales et fréquentielles) de la notion d'intensité d'implication et de son uniformisation,

- d'autre part :

- . dans la construction nouvelle d'un indice d'implication entre classes  $\mathcal{A}$  et  $\mathcal{B}$  de variables, étendant l'évaluation et l'implication  $a \Rightarrow b$  à celle de  $\mathcal{A} \Rightarrow \mathcal{B}$ , sur la base d'une cohésion implicative suffisante des classes examinées et d'une relation implicative maximale entre leurs éléments respectifs,

- . enfin dans l'organisation en structure arborescente de l'ensemble de classes, organisation empruntée à la classification hiérarchique de I.C. LERMAN. Mais ici, le lien entre deux noeuds de l'arbre est orienté.

Nous examinerons plus en détail ces deux volets.

Les concepts théoriques qui s'y rattachent ne doivent pas être considérés comme des spéculations détachées de la volonté de construire des modèles de comportements ou de processus de pensée. Au contraire, ils s'inscrivent en réponse, provisoire sans doute, à quelques questions levées par notre problématique actuelle en didactique des mathématiques, questions qui serviront de support intuitif à la modélisation puis la formalisation qui suivront :

- . à un niveau de cursus donné, peut-on, dans une situation-problème donnée, déterminer une hiérarchie partiellement ordonnée de procédures de résolution de problèmes de mathématiques, signes d'une connaissance en voie de constitution ?

- . à un niveau de cursus donné, peut-on définir à partir de classes ordonnées de procédures, des conceptions homogènes et résistantes relativement à un certain savoir <sup>(1)</sup> ?

etc.

Comme nous en verrons un exemple dans le § 3.1 nous avons utilisé les modèles statistiques élaborés pour des variables (§ 1), puis des classes des variables (§ 2) afin d'outiller le didacticien dans l'approche de questions telles que ci-dessus.

Notre problématique croise celle de certains chercheurs en intelligence artificielle dont nous reparlerons plus loin et dont le souci est :

- d'une part, la considération modale de la relation d'attribution d'un descripteur déterminé à un objet (ou un sujet),
- d'autre part, la reconnaissance de formes.

(1) Citons, par exemple, la conception "entiers naturels" que les jeunes enfants ont des décimaux, compte tenu des travaux scolaires sur la mesure des grandeurs, décimaux qui ne seraient que des entiers par un changement convenable d'unité.

Cependant, comme nous venons de le voir, bien que la première problématique ne soit pas orientée par les problèmes d'apprentissage (au sens de l'I.A.) ou de structure de base de connaissances, nos points de vue pourraient se rejoindre dans la nécessité de placer les sujets, ayant conduit à une classification hiérarchique ou implicite, par rapport aux ensembles de variables classifiées.

## § 1 - IMPLICATION ENTRE VARIABLES BINAIRES ET EXTENSION.

### 1.1. Modélisation.

Dans le cas binaire, la situation générique est la suivante. Croisant une population  $E$  et un ensemble de variables  $V$  et du fait de l'observation exceptionnelle de l'implication stricte de la variable  $a$  sur la variable  $b$ , on veut donner un sens statistique à une implication non stricte :  $a \Rightarrow b$ . En termes ensemblistes,  $A$  et  $B$  représentant les sous-populations possédant respectivement  $a$  et  $b$ , il y a équivalence à mesurer l'inclusion non stricte de  $A$  dans  $B$ .

Par suite, s'inspirant de la méthode de I.C. LERMAN [81] pour définir la similarité, R. GRAS [79] axiomatise la notion d'implication statistique de la façon suivante :

Soit  $X$  et  $Y$  deux parties aléatoires quelconques de  $E$  et de mêmes cardinaux respectifs que  $A$  et  $B$ , et  $\bar{Y}$  et  $\bar{B}$  les complémentaires respectifs de  $Y$  et de  $B$ .

$(a \Rightarrow b)$  est admissible au niveau de la confiance  $\alpha$  si et seulement si, dans une hypothèse d'indépendance (ou d'absence de lien a priori),  $\Pr[\text{card}(X \cap \bar{Y}) \leq \text{card}(A \cap \bar{B})] \leq 1 - \alpha$ .

Intuitivement et qualitativement, ceci signifie que l'implication  $a \Rightarrow b$  sera admissible à l'issue d'une expérience si le nombre d'individus de  $E$  la contredisant dans l'expérience est invraisemblablement petit par rapport au nombre d'individus attendu dans une hypothèse d'absence de lien.

La modélisation probabiliste que nous retenons de façon privilégiée est décrite par un processus de tirage aléatoire en 3 étapes (cf. I.C. LERMAN, R. GRAS, H. ROSTAM [81]). Notons  $n_a, n_b, n_{\bar{b}}, n_{a\bar{b}}, n_{a\bar{b}}$ , les cardinaux respectifs de  $A, B, \bar{B}, A \cap \bar{B}, A \cup \bar{B}$  :

. on considère le référentiel  $E$  comme la réalisation d'un référentiel aléatoire  $\mathcal{E}$  dont le cardinal  $\mathcal{N}$  serait une variable aléatoire de Poisson de paramètre le cardinal  $n$  de  $E$  observé :

$$\Pr[\mathcal{N} = m] = \frac{n^m}{m!} e^{-n}$$

. le choix aléatoire d'une partie quelconque (par exemple  $X$ ) de cardinal aléatoire  $K$  pour une distribution uniforme de probabilité sur les éléments de  $\mathcal{E}$  et égale à  $\frac{d}{m}$  (dans le cas de  $X, d = n_a$ ) est de type binomial :

$$\Pr[K = k/\mathcal{N} = m] = \binom{m}{k} \left(\frac{d}{m}\right)^k \left(1 - \frac{d}{m}\right)^{m-k} \quad (\text{pour } k \leq m)$$

$X$  et  $\bar{Y}$  étant deux parties quelconques choisies de façon indépendante parmi les parties ayant respectivement pour cardinaux  $n_a$  et  $n_{\bar{b}}$ , la probabilité qu'un élément de  $\mathcal{E}$  appartienne à  $X \cap \bar{Y}$  est :  $p(a) p(\bar{b})$  où  $p(a) = \frac{n_a}{m}$  et  $p(\bar{b}) = \frac{n_{\bar{b}}}{m}$ .

La loi de probabilité du cardinal de  $X \cap \bar{Y}$  est binomiale de paramètres  $m$  et  $\pi = p(a) p(\bar{b})$

$$\Pr[\text{Card}(X \cap \bar{Y}) = s/\mathcal{N} = m] = \binom{m}{s} \pi^s (1 - \pi)^{m-s}$$

pour  $s \leq n_{a\bar{b}}$  et  $m \geq n_{a\bar{b}}$ .

Par suite, en faisant varier le conditionnement de  $\mathcal{N}$ , on obtient :

$$\begin{aligned} \Pr[\text{Card}(X \cap \bar{Y}) = s] &= \sum_{m \geq s} \Pr[\text{Card}(X \cap \bar{Y}) = s/\mathcal{N} = m] \times \Pr[\mathcal{N} = m] \\ &= \frac{(n\pi)^s}{s!} e^{-n\pi}. \end{aligned}$$

La variable aléatoire  $\text{Card}(X \cap \bar{Y})$  suit donc la loi de Poisson de paramètre  $n\pi = n p(a) p(\bar{b})$  (de moyenne et de variance  $n\pi$ ).

Comme I.C. LERMAN, nous réduisons et centrons cette variable de Poisson en la variable :

$$Q(a, \bar{b}) = \frac{\text{Card}(X \cap \bar{Y}) - n p(a) p(\bar{b})}{\sqrt{n p(a) p(\bar{b})}} = \frac{\text{Card}(X \cap \bar{Y}) - \frac{n_a n_{\bar{b}}}{n}}{\sqrt{\frac{n_a n_{\bar{b}}}{n}}}.$$

Dans l'expérience, la valeur observée de  $Q(a, \bar{b})$  est  $q(a, \bar{b}) = \frac{n_{a\bar{b}} - \frac{n_a n_{\bar{b}}}{n}}{\sqrt{\frac{n_a n_{\bar{b}}}{n}}}$ .

Dans les cas légitimant convenablement l'approximation, la variable  $Q(a, \bar{b})$  suit la loi normale centrée réduite. L'intensité d'implication, qualité de l'admissibilité de  $a \Rightarrow b$ , pour  $n_a \leq n_b$ , est alors définie à partir de l'indice  $q(a, \bar{b})$  par :

$$\begin{aligned} \varphi(a, \bar{b}) &= 1 - \Pr[Q(a, \bar{b}) \leq q(a, \bar{b})] \\ &= \frac{1}{\sqrt{2\pi}} \int_{q(a, \bar{b})}^{\infty} e^{-\frac{t^2}{2}} dt \end{aligned}$$

Ainsi, pour un niveau de confiance  $\alpha$ , l'implication  $a \Rightarrow b$  sera admissible si et seulement si :

$$\varphi(a, \bar{b}) \geq \alpha.$$

Par exemple, à  $\alpha = 0,95$  correspond la valeur de l'indice :  $q(a, \bar{b}) = -1,65$ . De même, si  $\alpha = 0,5$  alors  $q(a, \bar{b}) = 0$  et  $\alpha \geq 0,5$  équivaut à  $q(a, \bar{b}) \leq 0$ .

On notera que l'approche ci-dessus peut s'exprimer en termes de test d'hypothèse. Cependant, nous ne retenons pas cette voie qui, du fait de sa visée de prise de décision, limiterait les considérations qui vont suivre. Mais cette possibilité marque bien la différence avec l'approche de J. LOEVINGER (LOEVINGER [1947] qui définissait la quasi-implication de  $a$  sur  $b$  par l'indice :

$$H(a, b) = 1 - \frac{n_{a\bar{b}}}{n_a n_{\bar{b}}}.$$

Cet indice présente l'inconvénient, ne se référant pas à une échelle de probabilité, de ne pas fournir de seuil de vraisemblance et d'être invariant dans toute dilatation de  $E, A, B$  et  $A \cap \bar{B}$ . Cette limitation apparaît dans l'approche de J. PEARL (PEARL [1988]), de S. ACID et als (ACID [1991]) et A. GAMMERMAN, Z. LUO (GAMMERMAN A et LUO Z. [1991]). Chez ces derniers chercheurs, c'est l'écart entre la distribution conjointe entre  $a$  et  $b$  (et non  $a$  et  $\bar{b}$ ) et la distribution produit qui tient lieu de critère comparatif.

## 1.2. Quelques propriétés du modèle de l'implication statistique.

A. LARHER dans sa thèse [1991] étudie différentes propriétés de  $q$  et  $\varphi$ . Retenons ici celles qui nous semblent les plus importantes :

- si,  $n_a$  étant fixé et  $A$  inclus dans  $B$ ,  $n_b$  tend vers  $n$  ( $B$  croît vers  $E$ ), alors  $\varphi(a, \bar{b})$  tend vers 0,5. Un prolongement par continuité nous permet donc de définir : si  $B = E$  alors  $\varphi(a, \bar{b}) = 0,5$ .

- L'indice  $q(a, \bar{b}) = \frac{n_{a\bar{b}} - \frac{n_a n_{\bar{b}}}{n}}{\sqrt{\frac{n_a n_{\bar{b}}}{n}}}$  est la racine carrée de la contribution de la case  $(a, \bar{b})$  à la

statistique du  $\chi^2$  attachée au croisement des 2 variables  $a$  et  $b$ .

- On démontre aisément que pour toute variable  $a$  :

$$0,95 \leq \varphi(a, \bar{a}) \leq 1 \Leftrightarrow n_a \in \left[ \frac{n - \sqrt{n(n-11)}}{2} ; \frac{n + \sqrt{n(n-11)}}{2} \right]. \quad (1)$$

Or, pour  $n$  assez grand, l'intervalle ci-dessus est voisin de  $[0, n]$ , ce qui permet d'affirmer que l'implication statistique  $a \Rightarrow a$  a un sens, tout en réservant une limite de confiance à la stabilité du caractère reproductible de la variable  $a$ .

• La relation  $\mathcal{R}$  sur  $V^2$  définie par :

$$\forall (a, b) \in V^2 \quad a \mathcal{R} b \text{ dès que } \varphi(a, \bar{b}) \geq 0,95 \text{ et que } n_a \text{ vérifie (1)}$$

est donc réflexive mais ni symétrique, ni antisymétrique, ni transitive. Pour lui associer un graphe valué, sans cycle et transitif, on en prendra la restriction  $\mathcal{R}'$  aux variables vérifiant la condition : si  $a \mathcal{R}' b$  et  $b \mathcal{R}' c$ , alors l'arc  $(a, c)$  appartient au graphe seulement si  $\varphi(a, \bar{c}) \geq 0,5$ .  $\mathcal{R}'$  définit alors un préordre partiel et permet une représentation claire de la relation d'implication statistique (cf. GRAS [1979]). Ce seuil de 0,5 permet, en outre, la satisfaction d'un objectif d'accroissement informationnel. En effet, si  $I$  est l'incertitude associée aux variables  $a$  et  $b$ , il permet de vérifier  $I(a|b) \leq I(a)$ .

Notre approche diffère également de celle de S. AMARGER et als (AMARGER S., DUBOIS D. et PRADE H. [1991]) qui, à partir d'une certaine inférence (une probabilité conditionnelle telle que  $p(b|a)$  appartient à un intervalle, sans être parfaitement connue), induisent transitivement, de proche en proche, des probabilités conditionnelles sur un graphe incomplet, et cela sans la contrainte d'un seuil. Notre problématique, à l'opposé, vise l'analyse d'un tableau donné, sans ambition inductive a priori, mais en imposant un seuil de transitivité.

• Comparons le coefficient de corrélation  $\rho(a, b)$  et l'indice  $q(a, \bar{b})$ . On note que :

$$q(a, \bar{b}) = \frac{n_{a\bar{b}} - \frac{n_a n_{\bar{b}}}{n}}{\sqrt{\frac{n_a n_{\bar{b}}}{n}}} = \frac{n_a n_b - n n_{a\bar{b}}}{\sqrt{n n_a n_{\bar{b}}}}$$

or  $\rho(a, b) = \frac{n \cdot n_{a\bar{b}} - n_a n_b}{\sqrt{n_a n_b n_{\bar{a}} n_{\bar{b}}}}$ . Ainsi  $q(a, \bar{b}) = 0 \Leftrightarrow \rho(a, \bar{b}) = 0$ .

Supposons  $q(a, \bar{b}) \neq 0$ . Alors :  $\frac{\rho(a, b)}{q(a, \bar{b})} = -\sqrt{\frac{n}{n_b n_{\bar{a}}}}$ .

Ainsi  $\rho(a, b) \geq 0$  est équivalent à  $\varphi(a, \bar{b}) \geq 0,5$ . Ceci signifie que implication et corrélation linéaire vont plutôt "dans le même sens". Cependant, on peut observer une croissance de l'implication en même temps qu'une décroissance de la corrélation, ce qui montre bien, qu'outre la dépendance aux effectifs  $n$ ,  $n_{\bar{a}}$  et  $n_b$ , le rapport  $\frac{\rho}{q}$  indique la non-coïncidence des deux concepts.

• Etudions la sensibilité de  $q$  aux faibles variations d'effectif. Supposons pour cela que, par exemple,  $n_{a\bar{b}}$  devienne  $n'_{a\bar{b}} = n_{a\bar{b}} + k$  où  $k \in \mathbb{Z}$  sans que varient  $n_a$  et  $n_{\bar{b}}$ . Alors :

$$q'(a, \bar{b}) = q(a, \bar{b}) + \frac{k}{\sqrt{\frac{n_a n_{\bar{b}}}{n}}}.$$

L'effet de l'erreur de mesure  $k$  est atténué en  $\sqrt{\frac{n}{n_a n_{\bar{b}}}}$ , ce qui permet dans la plupart des cas de maintenir la validité d'une implication au même seuil ou à un seuil voisin.

### 1.3. Extension à des variables numériques.

Dans sa thèse, A. LARHER considère deux autres types de variables :

. **variables modales** associées à une logique multivalente où les valeurs de vérité ne sont plus seulement "vrai (1)" ou "faux (2)" mais toute valeur de  $[0,1]$ . On retrouve ici un type de variable que E. DIDAY (DIDAY E. [1991]) axiomatise et étudie en profondeur dans le cadre de l'analyse des connaissances (au sens de l'I.A.). Mais, notre théorisation sera plus élémentaire ;

. **variables quantitatives** spécifiant un certain nombre de fois où, par exemple, la modalité  $a$  est présente chez le sujet  $i$  (individu, classe d'individus, ...).

En fait, une suggestion de I.C. LERMAN confortant notre propre approche, permet d'interpréter ces variables en terme de **variables fréquentielles** définies de la façon suivante :

Soit  $(a, b) \in V^2$  et  $i \in E$ ,  $\alpha$  (resp.  $\beta$ ) la valeur maximum de  $a$  (resp.  $b$ ) sur  $E$ .  $\alpha_i$  (resp.  $\beta_i$ ) la valeur observée de  $a$  (resp.  $b$ ) sur  $i$ . Posons :

$$\alpha'_i = \frac{\alpha_i}{\alpha}, \beta'_i = \frac{\beta_i}{\beta}, \nu_a = \sum_{i \in E} \alpha'_i, \nu_b = \sum_{i \in E} \beta'_i, \nu_{a\bar{b}} = \sum_{i \in E} \alpha'_i (1 - \beta'_i) \text{ et } \nu_{\bar{b}} = n - \nu_b.$$

Ces valeurs pour  $\alpha = \beta = 1$ ,  $\alpha_i$  et  $\beta_i = 1$  ou  $0$  étendent alors le cas binaire où les effectifs fréquentiels sont  $n_a, n_b$  et  $n_{a\bar{b}}$ . En effet, si  $1_A$  et  $1_B$  sont les fonctions indicatrices des sous-ensembles d'individus possédant respectivement les caractères  $a$  et  $b$ , l'effectif  $n_{a\bar{b}}$  s'écrit :

$$n_{a\bar{b}} = \sum_{i \in E} 1_A(i) [1 - 1_B(i)].$$

Donc posons, comme pour les variables binaires :

$$q(a, \bar{b}) = \frac{\nu_{a\bar{b}} - \frac{\nu_a \nu_b}{n}}{\sqrt{\frac{\nu_a \nu_{\bar{b}}}{n}}}$$

qui sera, pour  $\nu_a \leq \nu_b$ , l'indicateur retenu pour l'implication statistique de  $a$  sur  $b$ .



## § 2 - IMPLICATION ENTRE CLASSES DE VARIABLES.

Elle ne prend véritablement son sens qu'à condition qu'à l'intérieur de chaque classe dont on examine la relation avec d'autres, existe une certaine "cohésion" entre les variables qui la constituent. Cette cohésion, généralement nourrie de cohérence sémantique ou, dans le cas de la didactique, de conditions psychologiques, cognitives, situationnelles, etc., doit se traduire ici par une mesure (quantitative). On pourrait penser qu'un ensemble d'indices de similarité assez élevés entre les éléments de la classe serait un bon indicateur de cohésion. Nous ne retenons pas cette approche qui ne rendrait compte qu'une d'une cohésion de profils symétriquement comparables, ne restituant pas une dynamique interne orientée (donc non symétrique). Or nous disposons avec les intensités d'implication entre variables d'un instrument de mesure d'un emboîtement de deux parties d'une population  $E$ . C'est donc cette voie que nous choisissons pour une cohésion implicative donc orientée, comme peut l'être une filiation procédurale ou une genèse. Nous verrons ensuite quel indicateur permettrait de rendre compte d'une extension aux classes, de la notion d'implication.

### 2.1. Cohésion implicative.

Afin d'en améliorer l'intuition, nous la définirons progressivement pour 2, puis 3 et  $r$  éléments de classe.

La cohésion, se voulant indicateur d'ordre implicatif au sein d'une classe d'attributs, s'oppose en cela au "désordre" dont rend compte l'entropie d'une expérience aléatoire. Rappelons au sujet de celle-ci que,  $X$  étant une variable aléatoire prenant ses valeurs dans  $S = \{m_1, m_2, \dots, m_k\}$  muni de la loi  $\{p_1, p_2, \dots, p_k\}$ , l'entropie est l'espérance mathématique de la variable  $I(X)$  prenant les valeurs  $I(m_1), I(m_2), \dots, I(m_k)$ ;  $I(m_j)$  est l'incertitude sur  $\{m_j\}$  ou information apportée par la réalisation de  $\{m_j\}$ . Ainsi :

$$\mathcal{E}[I(X)] = \sum_{j=1}^k -p_j \log_2 p_j$$

est l'entropie de l'expérience.

#### . Cas de 2 éléments : classe $(a, b)$ .

Supposons  $n_a < n_b$ . Nous allons définir la cohésion du couple  $(a, b)$ .

Soit  $\chi$  la variable aléatoire indicatrice de l'évènement  $[Q(a, \bar{b}) \geq q(a, \bar{b})]$ . Alors :

$$\Pr(\chi = 1) = \varphi(a, \bar{b}) = p$$

et 
$$\Pr(\chi = 0) = 1 - \varphi(a, \bar{b}) = 1 - p.$$

L'entropie ou incertitude de cette expérience est alors :

$$\mathcal{E}[I(\chi)] = p \log_2 p - (1-p) \log_2 (1-p) .$$

Par exemple, si  $\varphi(a, \bar{b}) = p = 0,95$ , alors :

$$\mathcal{E}[I(X)] = \frac{-0,95 \ell n 0,95 - 0,05 \ell n 0,05}{\ell n 2} = 0,286 .$$

Notons que si  $\varphi(a, \bar{b}) = 1$ , alors  $\mathcal{E}[I(\chi)] = 0$  en convenant que  $0 \ell n 0 = 0$

si  $\varphi(a, \bar{b}) = 0,5$  alors  $\mathcal{E}[I(\chi)] = 1$  (entropie maximale)

et si  $\varphi'(a, \bar{b}) = 1 - \varphi(a, \bar{b})$ , alors  $\mathcal{E}[I(\chi)] = \mathcal{E}[I(\chi')]$ .

Plus précisément, en posant :  $\mathcal{E} = f(p) = -p \log_2 p - (1-p) \log_2 (1-p)$ , on a :

$$f(1-p) = f(p) : \text{symétrie par rapport à } p = 0,5 \text{ et } \frac{df}{dp} = \log_2 \frac{1-p}{p} \quad (0 < p < 1)$$

donc  $\mathcal{E}$  croît de 0 à 1 sur  $]0;0,5]$  et décroît de 1 à 0 sur  $[0,5;1[$ .

La propriété de symétrie de  $\mathcal{E}$  allant à l'encontre de la dissymétrie de la quasi-implication, nous retiendrons finalement comme indicateur de cohésion l'application  $c$  définie sur l'ensemble  $V$  des variables :

$$\begin{aligned} \text{si } \varphi(a, \bar{b}) = p \geq 0,5 \quad , \quad c(a, b) &= [1 - [p \log_2 p + (1-p) \log_2 (1-p)]^2]^{1/2} \\ &\text{(racine carrée du complémentaire à 1 du carré de l'entropie)} \\ \text{et si } \varphi(a, \bar{b}) = p < 0,5 \quad , \quad c(a, b) &= 0 \text{ (absence de cohésion).} \end{aligned}$$

La fonction "carré de l'entropie" est choisie pour des raisons d'analyse et de contraste. Nous prenons la racine carrée de son complément à 1 pour donner à la cohésion la dimension de l'entropie et pour accroître sa valeur numérique (en effet pour  $x \in [0,1]$ ,  $\sqrt{1-x^2} \geq 1-x$ ).

$$\text{Ainsi } c(a, b) = \sqrt{1-\mathcal{E}^2} .$$

$$\text{Posons : } c(a, b) = g(\mathcal{E}) = g[f(p)]$$

$$\frac{dg}{d\mathcal{E}} = - \frac{\mathcal{E}}{\sqrt{1-\mathcal{E}^2}} \quad \text{et} \quad \frac{dc}{dp} = \frac{-\mathcal{E}}{\sqrt{1-\mathcal{E}^2}} \times \log_2 \frac{1-p}{p} .$$

$c(a, b)$  croît donc de 0 à 1 quand  $p = \varphi(a, \bar{b})$  croît de 0,5 à 1. La fonction  $c$  de  $p$  est continue en  $p = \frac{1}{2}$ .

Nous prolongeons par continuité en prenant :

$c(a, b) = 1$  lorsque  $p = 1$  (c'est-à-dire lorsque l'implication  $a \Rightarrow b$  est stricte).

L'hypothèse initiale  $n_a \leq n_b$  n'est plus à formuler si pour tout paire d'éléments  $(a,b)$  de  $V$ , nous nous intéressons au  $\max[\varphi(a,\bar{b}), \varphi(b,\bar{a})]$  dont nous savons qu'il respecte l'ordre des cardinaux dans le cas où l'un d'entre les 2 nombres  $\varphi(a,\bar{b})$  et  $\varphi(b,\bar{a})$  est négatif.

Rappelons en effet le résultat montré dans la thèse d'A. LARHER :

si  $n_a \leq n_b$  et  $\varphi(a,\bar{b}) \leq 0$ , alors  $\varphi(a,\bar{b}) \geq \varphi(b,\bar{a})$  ; autrement dit si  $n_a \leq n_b$  (resp.  $n_b \leq n_a$ ) et  $\varphi(a,\bar{b}) \leq 0$  (resp.  $\varphi(b,\bar{a}) \leq 0$ ),  $\max [\varphi(a,\bar{b}), \varphi(b,\bar{a})] = \varphi(a,\bar{b})$  (resp.  $\varphi(b,\bar{a})$ ). Aussi la cohésion de la classe  $(a,b)$  est définie sans équivoque à partir de la plus grande des 2 valeurs de vérité des énoncés :

$$[a \Rightarrow b] \text{ et } [b \Rightarrow a]$$

par :

$$\begin{aligned} c(a,b) &= (1 - \varepsilon^2)^{1/2} \text{ où } \varepsilon = -p \log_2 p - (1-p) \log_2 (1-p), \\ &\text{si } p = \max [\varphi(a,\bar{b}), \varphi(b,\bar{a})] \geq 0,5 ; \\ c(a,b) &= 0 \text{ si } p \leq 0,5. \end{aligned}$$

Lorsque la classe est réduite à un seul élément, la cohésion de cette classe est  $c(a,a)$  ; nous avons vu (cf. p. 6) que la relation  $\mathcal{R}$  est réflexive et :

$$0,95 \leq \varphi(a,\bar{a}) \leq 1 \quad \text{pour} \quad n_a \in \left[ \frac{n - \sqrt{n(n-1)}}{2} ; \frac{n + \sqrt{n(n-1)}}{2} \right]. \quad (1)$$

Soit  $\varepsilon$  l'entropie :  $\varepsilon = -p \log_2 p - (1-p) \log_2 (1-p)$

$$\lim_{p \rightarrow 1^-} \varepsilon = 0.$$

Soit  $c$  la cohésion :  $c = \sqrt{1 - \varepsilon^2}$

$$\lim_{p \rightarrow 1^-} c = 1.$$

Nous prendrons donc  $c = 1$  pour tout couple d'éléments égaux ou donc pour toute classe réduite à un seul élément qui satisfait la relation (1).

**Exemple.**  $n = 100$

$$\text{si } n_a \in [3;97] \quad , \quad c(a,a) = 1.$$

**. Cas de 3 éléments  $a, b$  et  $c$ .**

Six valeurs d'intensité correspondent a priori à l'ensemble  $A = \{a,b,c\}$  :

$$\varphi(a,\bar{b}), \varphi(a,\bar{c}), \varphi(b,\bar{a}), \varphi(b,\bar{c}), \varphi(c,\bar{a}) \text{ et } \varphi(c,\bar{b}).$$

L'indice de cohésion implicative doit contenir l'information révélée par les relations implicatives binaires entre tous les éléments de l'ensemble  $A$ . Mais, en même temps, pour

conserver la dynamique dissymétrique de l'implication, seule la relation la plus puissante entre deux éléments quelconques reste pertinente par rapport à notre objectif. Par suite, parmi toutes les associations 3 à 3 ne faisant intervenir qu'une fois chaque couple d'éléments de  $\{a,b,c\}$  et restituant au mieux la puissance de certaines implications, nous retenons :

$$\max [\varphi(a,\bar{b}), \varphi(b,\bar{a})], \max [\varphi(a,\bar{c}), \varphi(c,\bar{a})] \text{ et } \max [\varphi(b,\bar{c}), \varphi(c,\bar{b})].$$

Comme précédemment, dans le cas où ils sont supérieurs ou égaux à 0,5, les maxima obtenus sont compatibles avec l'ordre des effectifs  $n_a, n_b$  et  $n_c$ . Par exemple, si  $n_a \leq n_b \leq n_c$ , les trois maxima sont :  $\varphi(a,\bar{b}), \varphi(a,\bar{c})$  et  $\varphi(b,\bar{c})$ .

Le couple  $\mathcal{A} = ((a,b),c)$  sera alors appelé **classe** et sa cohésion implicative sera définie ainsi :

$$C(\mathcal{A}) = [c(a,b) \times c(b,c) \times c(a,c)]^{1/3}$$

moyenne géométrique des cohésions des couples

La préférence accordée à la moyenne géométrique plutôt qu'à la moyenne arithmétique tient à notre volonté, d'une part d'obtenir une cohésion nulle pour une classe dès que la cohésion d'un de ses couples est nulle, c'est-à-dire dès que les implications mutuelles sont inférieures ou égales à 0,5, d'autre part de "ramener"  $C(\mathcal{A})$  au voisinage de 1 lorsque les cohésions des couples sont assez fortes.

### . Cas de $r$ éléments $a_1, a_2, \dots, a_r$ .

Nous opérons comme ci-dessus, c'est-à-dire en retenant les maxima des intensités d'implication entre 2 éléments quelconques de l'ensemble  $\dot{A} = \{a_1, a_2, \dots, a_r\}$ . A ces maxima sont associés les cohésions implicatives des couples et l'ordre induit sur  $\dot{A}$  par les effectifs  $n_{a_1}, n_{a_2}, \dots, n_{a_r}$ . Par exemple, si  $n_{a_1} \leq n_{a_2} \leq \dots \leq n_{a_r}$ , nous appellerons **classe** le couple  $\mathcal{A} = ((a_1, a_2), a_3, \dots, a_r)$ , et comme il y a  $\frac{r(r-1)}{2}$  paires, sa cohésion implicative sera :

$$C(\mathcal{A}) = \left[ \prod_{\substack{i \in \{1, \dots, r-1\} \\ j > i}} c(a_i, a_j) \right]^{\frac{2}{r(r-1)}}$$

## 2.2. Implication entre classes.

Nous souhaitons que l'implication entre deux classes se constitue à partir des informations suivantes :

- les cohésions respectives des 2 classes,
- une intensité d'implication extrême des éléments d'une classe sur les éléments de l'autre,
- les cardinaux respectifs des 2 classes.

Chacune de ces informations crédite l'indice que nous retiendrons si :

- l'indice croît avec les cohésions de chaque classe et s'annule lorsque la cohésion de l'une d'entre elles est nulle,
- l'indice croît avec la liaison extrême (minimale si l'on vise un degré d'exigence élevé, maximale si l'on recherche une souplesse réaliste),
- l'indice décroît avec les cardinaux des classes, eu égard à la prise en compte d'une liaison maximale.

Posons : .  $\dot{A}$  et  $\dot{B}$  deux parties disjointes :  $\dot{A} = \{a_1, \dots, a_r\}$  et  $\dot{B} = \{b_1, \dots, b_s\}$  ,  
 .  $\mathcal{A}$  et  $\mathcal{B}$  les classes qui leur sont respectivement associées ,  
 .  $C(\mathcal{A})$  et  $C(\mathcal{B})$  leurs cohésions respectives.

Conformément aux lois de probabilité des sup. de variables aléatoires, a priori uniformément distribuées, nous définissons l'indice d'implication  $\psi(\mathcal{A}, \mathcal{B})$  de la classe  $\mathcal{A}$  vers la classe  $\mathcal{B}$  par :

$$\psi(\mathcal{A}, \mathcal{B}) = \left\{ \sup_{\substack{i \in \{1, \dots, r\} \\ j \in \{1, \dots, s\}}} \varphi(a_i, \bar{b}_j) \right\}^{rs} \times [C(\mathcal{A}) \times C(\mathcal{B})]^{\frac{1}{2}} .$$

L'expression  $[C(\mathcal{A}) C(\mathcal{B})]^{\frac{1}{2}}$  représente la cohésion moyenne (géométrique) de  $\mathcal{A}$  et  $\mathcal{B}$  ; elle intègre les informations de cohésivité des 2 classes en jeu ; de plus, cette expression est telle que si  $C(\mathcal{A})$  et  $C(\mathcal{B})$  sont simultanément multipliées par  $k$ , alors  $\psi(\mathcal{A}, \mathcal{B})$  est multipliée par  $k$ .

#### Remarques.

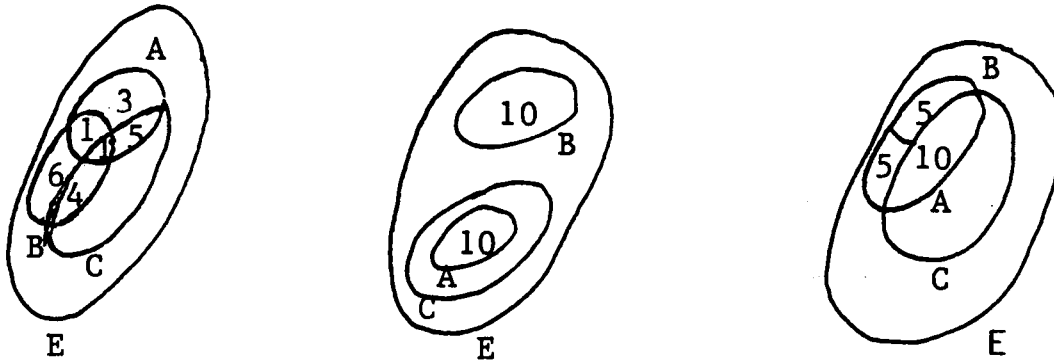
1°) Si  $\dot{A}$  et  $\dot{B}$  sont réduits à des singletons, et donc  $\mathcal{A}$  et  $\mathcal{B}$  à des couples d'éléments égaux – cette formule vérifie la définition de l'implication entre variables puisque dans ce cas  $C(\mathcal{A})=C(\mathcal{B})=1$  ; l'indice choisi coïncide alors avec l'intensité d'implication, soit :  $\psi((a,a), (b,b)) = \varphi(a, \bar{b})$ .

2°) La méthode de construction de l'indice d'implication ci-dessus montre à l'évidence qu'elle est indépendante de la nature des variables traitées, que celles-ci soient binaires, non binaires ou numériques.

3°) Un critère d'implication entre classes basé, dans le cas des variables binaires (ou attributs), sur la réunion ensembliste des ensembles d'individus possédant ces attributs n'est pas pertinent. D'une part, il privilégie les seules variables de cette nature, d'autre part il ne rend pas compte de la cohésion des classes.

Par exemple, si  $a, b$  et  $c$  sont des attributs représentés respectivement par les parties  $A, B$  et  $C$  et si, de plus,  $\text{Card}(A \cup B)$ ,  $\text{Card } C$ , et  $\text{Card}[(A \cup B) \cap \bar{C}]$  restent constants, on peut obtenir des

situations très différentes s'opposant à l'idée intuitive de l'implication entre classes telle que nous la concevons. La cohésion de  $(a,b)$  et les positions respectives de  $A, B$  et  $C$  dans 3 cas différents montrent de façon évidente l'inadéquation du choix de la réunion pour signifier l'implication :



$$\text{Card}(A \cup B) = 20$$

$$\text{Card}[(A \cup B) \cap \bar{C}] = 10$$

4°) Les moyens de calcul à développer pour examiner tous les indices d'implication de classes dans un ensemble de  $x$  variables (définies sur une population de  $n$  individus) croissent exponentiellement avec  $x$  ; en effet, l'ensemble des parties à considérer contient  $2^x$  éléments. Aussi, dans la réalité, nous pensons raisonnable de procéder en didactique de la façon suivante :

- calculer les implications entre variables,
- construire puis analyser le graphe d'implication,
- émettre des hypothèses quant à la cohérence de formation de classes par rapport à la problématique didactique,
- constituer de telles classes de variables (items ou modalités de réponse dans le cas d'un questionnaire) et évaluer les valeurs des indices d'implication entre ces classes prises 2 à 2 ; éventuellement, modifier sensiblement les contenus de ces classes pour améliorer les valeurs des indices d'implication. Un indicateur statistique, qui reste à déterminer, pourrait servir à définir un test d'arrêt de modification.

Par exemple, dans le cas où les variables seraient des procédures (disjointes ou non, mais identifiables en présence-absence chez chacun des individus), on pourrait obtenir des relations entre classes de procédures, gommant les effets microscopiques de procédures trop parcellisées et rendant compte de démarches générales, voire de conceptions, ayant une certaine stabilité. Notons à ce sujet que la condition de disjonction de  $A$  et  $B$  exprimée plus haut peut être levée sans affecter radicalement les objectifs déclarés à propos de l'implication entre classes : en effet, la cohésion de classes est un garant d'homogénéité "dirigée" qui relativise l'incidence qu'aurait le

premier facteur du produit définissant l'implication  $\psi(\mathcal{A}, \mathcal{B})$  où  $\mathcal{A}$  et  $\mathcal{B}$  sont les classes respectivement associées aux ensembles  $\dot{A}$  et  $\dot{B}$ .

Pour terminer ce paragraphe, soulignons l'intérêt des notions de cohésion et d'implication entre classes pour l'analyse implicative.

Soit une classe  $\mathcal{A}$  de variables de cohésion non nulle. Si l'une de ces variables,  $a$ , admet une implication inférieure à 0,5 sur une variable d'une classe  $\mathcal{B}$ , l'arbre implicatif ne pourra rendre compte de cette implication en raison de l'exigence d'une intensité supérieure ou égale à 0,5 pour la fermeture transitive. Par contre, l'implication entre classes ne sollicitant que le sup des implications dirigées de  $\mathcal{A}$  vers  $\mathcal{B}$ , l'attribut  $a$  va apparaître dans l'implication  $\mathcal{A} \Rightarrow \mathcal{B}$ . Ainsi, nous conservons une information qui aurait disparu si l'on s'était contenté d'envisager l'implication de variables seules.

### 2.3. Etude de la vraisemblance de $\psi(\mathcal{A}, \mathcal{B})$ .

Dans ce paragraphe, nous cherchons des critères permettant, à partir d'une valeur observée dans une situation, toujours dans l'hypothèse d'une absence de lien a priori, d'accorder un sens implicatif à la liaison non symétrique de l'ensemble de variables  $\dot{A}$  (classe associée :  $\mathcal{A}$ ) sur l'ensemble de variables  $\dot{B}$  (classe associée :  $\mathcal{B}$ ). En fait, à travers le questionnement de la valeur de

$$\psi(\mathcal{A}, \mathcal{B}) = [\sup_{i,j} \varphi(a_i, \bar{b}_j)]^{rs} \times [C(\mathcal{A}) C(\mathcal{B})]^{\frac{1}{2}}, \text{ nous avons la même attitude que celle}$$

que nous avons eue à l'égard de  $\varphi(a, \bar{b}) = 1 - Pr\{Q(a, \bar{b}) \leq q(a, \bar{b})\}$ . Pour cela, il nous faut examiner les lois régissant les variations des variables en jeu : intensités d'implication et cohésions implicatives. La complexité des calculs nous contraint à des études limitées mais déjà significatives des "degrés d'étonnement" des valeurs observées.

#### 2.3.1. Loi de variation de $c(a, b)$ .

Nous menons cette étude dans le cas des variables attributs, mais cette étude serait très proche dans les autres cas.

Lorsque la cohésion d'une classe  $\mathcal{A}$  est faible, la liaison implicative entre les éléments est cependant loin d'être négligeable. En effet, cette valeur faible reste cependant le témoin qu'il existe toujours entre 2 éléments quelconques de l'ensemble associé  $\dot{A}$  une implication unidirectionnelle au moins égale à 0,5. Précisons cela dans le cas où  $\dot{A} = \{a, b\}$ .

a) Deux attributs  $a$  et  $b$  étant donnés, considérons la variable aléatoire  $\Phi$  dont la valeur observée,  $n_{\alpha\bar{b}}$  étant connu, est  $\varphi(a, \bar{b})$ . Ainsi,  $[\Phi \geq 1 - \alpha]$  si et seulement si  $Q(a, \bar{b}) \leq \varphi'^{-1}(\alpha)$  où  $\varphi'^{-1}$  est la fonction réciproque de l'intégrale gaussienne :

$$\alpha = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\varphi(a, \bar{b})} e^{-\frac{t^2}{2}} dt.$$

Revenant à la nature "poissonnienne" de la variable  $\text{Card}(X \cap \bar{Y})$  définie au § 1.1, représentant le cardinal aléatoire de  $E_a \cap E_{\bar{b}}$  de valeur observée  $n_{\alpha\bar{b}}$ , on obtient :

$$\Pr[\Phi \geq 1 - \alpha] = \Pr[\text{Card}(X \cap \bar{Y}) \leq k] \text{ où } \varphi' \left( \frac{k - \frac{n_a n_{\bar{b}}}{n}}{\sqrt{\frac{n_a n_{\bar{b}}}{n}}} \right) = \alpha.$$

$$\text{Soit } \Pr[\Phi \geq 1 - \alpha] = \sum_{\ell=0}^k e^{-\lambda} \frac{\lambda^\ell}{\ell!} \text{ avec } \lambda = \frac{n_a n_{\bar{b}}}{n}.$$

Bien entendu, pour  $n$  assez grand, nous retrouvons :  $\Pr[\Phi \geq 1 - \alpha] \approx \alpha$ , donc  $\Phi$  est uniformément distribuée sur  $[0,1]$ .

Par exemple, si  $n = 100$ ,  $n_a = 15$ ,  $n_b = 20$  :

$$\Pr[\Phi \geq 0,9] = 0,086 \quad (k < 7,57)$$

$$\Pr[\Phi \geq 0,95] = 0,043 \quad (k < 6,28).$$

b) Soit  $\alpha' = 1 - \alpha \geq 0$ . Alors, revenant sur les variations de  $c(a,b)$  en fonction de  $\varphi(a, b^-)$ , examinons la loi de la variable aléatoire  $\Gamma(a,b)$ , de réalisation  $c(a,b)$  en fonction de celle de  $\Phi$ .

$$\begin{aligned} \Pr[\Gamma(a,b) > \alpha'] &= \Pr\{1 - [\Phi \log_2 \Phi + (1 - \Phi) \log_2 (1 - \Phi)]^2 > \alpha'^2 \text{ et } \Phi \geq 0,5\} \\ &= \Pr\left\{1 - \frac{1}{(\ln 2)^2} [\Phi \ln \Phi + (1 - \Phi) \ln (1 - \Phi)]^2 > \alpha'^2 \text{ et } \Phi \geq 0,5\right\} \\ &= \Pr\{(1 - \alpha'^2)(\ln 2)^2 > [\Phi \ln \Phi + (1 - \Phi) \ln (1 - \Phi)]^2 \text{ et } \Phi \geq 0,5\}. \end{aligned}$$

Supposant  $\Phi \geq 0,5$  réalisé :

$$\begin{aligned} \Pr[\Gamma(a,b) > \alpha'] &= \Pr\left\{ \frac{|\Phi \ln \Phi + (1 - \Phi) \ln (1 - \Phi)| < \ln 2 \sqrt{1 - \alpha'^2}}{\leq 0} \right\} \\ &= \Pr\{\Phi \ln \Phi + (1 - \Phi) \ln (1 - \Phi) > -\ln 2 \sqrt{1 - \alpha'^2}\}. \end{aligned}$$



Mais,  $\Phi \ln \Phi + (1 - \Phi) \ln (1 - \Phi) = \ln \Phi^\Phi (1 - \Phi)^{1-\Phi}$

d'où  $\Pr[\Gamma(a,b) > \alpha'] = \Pr [\Phi^\Phi (1 - \Phi)^{1-\Phi} > \frac{1}{2\sqrt{1-\alpha'^2}} ]$ .

Par exemple :

\* si  $\alpha' = 0,95$ ,  $\frac{1}{2\sqrt{1-\alpha'^2}} = 0,8054$ .

D'où  $\Pr [\Gamma(a,b) > 0,95] = \Pr [\Phi > 0,945] = 0,055$ .

\* Réciproquement,  $\Pr [\Phi > 0,95] = \Pr [\Phi^\Phi (1 - \Phi)^{1-\Phi} > 0,820]$

$$= \Pr [\Gamma(a,b) > 0,958]$$

en effet  $\frac{1}{2\sqrt{1-\alpha'^2}} > 0,820$

pour  $2\sqrt{1-\alpha'^2} < 1,2195$

soit  $\sqrt{1-\alpha'^2} < 0,28629$

ou encore  $\alpha' > 0,958$ .

\* On trouve également :

$$\Pr [\Gamma(a,b) > 0,9] = 0,09$$

$$\Pr [\Gamma(a,b) > 0,8] = 0,15.$$

### 2.3.2. Etude de la cohésion $C(\mathcal{A})$ .

$\mathcal{A}$  est la classe associée à l'ensemble  $A = \{a_i\}_{i=1,\dots,r}$  ; rappelons que,  $c(a_i, a_j) = (1 - \mathcal{E}^2)^{\frac{1}{2}}$  étant la meilleure des 2 cohésions  $c(a_i, a_j)$  et  $c(a_j, a_i)$  et  $\mathcal{E}$  l'entropie associée à l'implication  $a_i \Rightarrow a_j$ , alors

$$C(\mathcal{A}) = \left[ \prod_{i,j} c(a_i, a_j) \right]^{\frac{2}{r(r-1)}}$$

soit encore :  $\ln C(\mathcal{A}) = \frac{2}{r(r-1)} \sum_{i,j} \ln c(a_i, a_j)$ .

Dans une hypothèse d'absence de lien, les variables aléatoires  $\Gamma(a_i, a_j)$ , dont les cohésions  $c(a_i, a_j)$  sont des réalisations, sont indépendantes et  $\ln \Gamma(\mathcal{A})$  a pour loi le produit de convolution des lois de ces variables  $\Gamma(a_i, a_j)$ , au coefficient  $\frac{2}{r(r-1)}$  près. Or nous avons vu la difficulté d'accéder à

la loi de  $\Gamma(a_i, a_j)$ . Il nous sera, par contre, loisible d'effectuer des simulations pour estimer la loi de  $\Gamma(\mathcal{A})$ . Cette tâche est une composante de la thèse d'André TOTOHASINA, à Rennes.

Notons cependant que :

$$\text{si } \forall i, \forall j, \quad c(a_i, a_j) > 0,95, \quad \text{alors } C(\mathcal{A}) > 0,95.$$

### 2.3.3. Etude du terme $[\sup_{i,j} \varphi(a_i, \bar{b}_j)]^{r^s}$ .

Comme nous l'avons fait pour l'intensité d'implication entre attributs, nous considérons chaque probabilité  $\varphi(a_i, \bar{b}_j)$  comme la réalisation d'une variable aléatoire  $\Phi_{ij}$  uniformément distribuée sur  $[0,1]$ . Dans l'hypothèse d'absence de lien a priori entre les attributs, les variables aléatoires  $\Phi_{ij}$  sont indépendantes. Par suite, si  $S = \sup_{i,j} \Phi_{ij}$ ,

$$\text{alors} \quad \Pr[S \geq 1 - \alpha] = \alpha.$$

### 2.3.4. Etude de $\psi(\mathcal{A}, \mathcal{B})$ , dans le cas où $\mathcal{A}$ et $\mathcal{B}$ ne contiennent que 2 éléments au plus.

$$\text{Rappelons que : } \psi(\mathcal{A}, \mathcal{B}) = [\sup_{i,j} \varphi(a_i, \bar{b}_j)]^{r^s} (C(\mathcal{A}) C(\mathcal{B}))^{\frac{1}{2}}$$

$$\text{donc : } \Pr[\Phi > \alpha'] \leq \Pr[(\sup_{i,j} \Phi_{ij})^{r^s} > \alpha' \text{ et } C(\mathcal{A}) > \alpha'^2 \text{ et } C(\mathcal{B}) > \alpha'^2]$$

$$\leq \Pr[(\sup_{i,j} \Phi_{ij})^{r^s} > \alpha'] \cdot \Pr[C(\mathcal{A}) > \alpha'^2] \cdot \Pr[C(\mathcal{B}) > \alpha'^2].$$

dans une hypothèse d'indépendance a priori.

Par exemple,

1°) si nous souhaitons un "degré d'étonnement" ou indice d'implication de classes  $\psi(\mathcal{A}, \mathcal{B})$  qui soit supérieur ou égal à 0,97, il suffit de prendre  $\alpha' > 0,7$  car alors :  $C(\mathcal{A})$  et  $C(\mathcal{B})$  sont supérieures à 0,49 donc pour chaque classe  $\mathcal{A}$  et  $\mathcal{B}$  on a :  $\varphi^\varphi(1 - \varphi)^{1-\varphi} > 0,546$ , soit  $\varphi > 0,7$ .

$$\text{Par suite, } \Pr[\Gamma(\mathcal{A}) > 0,49] = \Pr[\varphi^\varphi(1 - \varphi)^{1-\varphi} > 0,546] = \Pr[\varphi > 0,7] \leq 0,3.$$

Dans ce cas,  $\Pr[(\sup_{i,j} \Phi_{ij})^{r^s} > 0,7] \times \Pr[\Gamma(\mathcal{A}) > 0,49] \times \Pr[\Gamma(\mathcal{B}) > 0,49]$  est inférieure ou égale à  $0,3 \times 0,3 \times 0,3 \leq 0,03$  et  $\Pr[\Phi > \alpha'] \leq 0,03$ , soit un "degré d'étonnement" supérieur ou égal à 0,97.

Ceci signifie qu'une observation  $\psi(\mathcal{A}, \mathcal{B}) > 0,7$ , pour  $\mathcal{A}$  et  $\mathcal{B}$  constituées de couples, constitue un événement d'intensité d'implication de classes supérieure à 0,97.

2°) Pour  $\psi(a,b) \geq 0,95$ , il suffit de prendre  $\alpha' \geq 0,63$  car alors :  $C(\mathcal{A}) > 0,4$  et  $C(\mathcal{B}) > 0,4$  et donc pour chaque classe  $\mathcal{A}$  et  $\mathcal{B}$ , on a :  $\varphi^\varphi(1 - \varphi)^{1-\varphi} > 0,53$  soit  $\varphi > 0,65$  et  $\Pr[\Gamma(\mathcal{A}) > 0,4] \leq 0,35$ .

$$\begin{aligned} \text{Dans ce cas : } \Pr[(\sup_{i,j} \Phi_{ij})^{r^s} > 0,63] \times \Pr[\Gamma(\mathcal{A}) > 0,4] \times \Pr[\Gamma(\mathcal{B}) > 0,4] \\ \leq 0,37 \times 0,35 \times 0,35 \leq 0,05. \end{aligned}$$

Une observation  $\psi(\mathcal{A}, \mathcal{B}) \geq 0,63$  constitue un événement d'intensité d'implication de classes supérieure à 0,95.

Notons que cette intensité d'implication de la classe  $\mathcal{A}$  sur la classe  $\mathcal{B}$  joue par rapport à la valeur observée, indice d'implication  $\psi(\mathcal{A}, \mathcal{B})$ , le même rôle que celui joué par l'intensité  $\varphi(a, \bar{b})$ , dans le cas de 2 attributs, par rapport à l'indice  $q(a, \bar{b})$ . Cette intensité mesure la qualité de l'"étonnement" que nous avons, face à une observation dans l'hypothèse d'absence de lien.

## 2.4. Agrégations successives des classes.

### 2.4.1. Algorithme.

Selon l'objectif classique des méthodes hiérarchiques, nous allons définir un algorithme d'agrégations successives des classes.

**1<sup>ère</sup> étape :** considérant toutes les paires de variables, on détermine le couple  $(a_i, a_j)$  conduisant au max  $(a_k, a_\ell)$  ; celui-ci correspond au max  $\varphi(a_k, \bar{a}_\ell)$ . En cas de solutions multiples, on choisit le couple  $(a_i, a_j)$  tel que  $n_{a_i}$  soit le plus petit effectif parmi tous les effectifs des attributs figurant en 1<sup>ère</sup> place des couples ; si à nouveau il y a solutions multiples, on choisit  $a_j$  tel que  $n_{a_j}$  soit l'effectif minimum parmi tous les effectifs des variables figurant en 2<sup>ème</sup> place des couples. En cas de nouvelle équivoque, on choisit l'un quelconque des couples restants.

Ainsi,  $\forall (k, \ell), \varphi(a_i, \bar{a}_j) \geq \varphi(a_k, \bar{a}_\ell)$ , et par suite :  $\forall (k, \ell), c(a_i, a_j) \geq c(a_k, a_\ell)$  et on agrège l'ensemble  $\{a_i, a_j\}$  en la classe  $(a_i, a_j)$  orientée de  $a_i$  vers  $a_j$ . Notons que, bien que  $n_{a_i} \leq n_{a_j}$ ,  $n_{a_i}$  n'est peut-être pas l'effectif minimum.

**2<sup>ème</sup> étape :** on compare les cohésions obtenues selon les phases  $\mathcal{O}_2$  et  $\mathcal{O}_3$  suivantes et on conserve la classe correspondant à la plus forte.

$\mathcal{O}_2$  : on détermine le couple (ou l'un des couples) optimal dont la cohésion est immédiatement inférieure ou égale à la cohésion retenue à l'étape précédente ;

$\mathcal{O}_3$  : on détermine l'ensemble  $\{a_i, a_j, a_k\}$  puis la classe associée dont la cohésion est meilleure que l'une quelconque des cohésions des couples restant après l'étape précédente.

Notons alors que  $c(a_i, a_j) \geq C((a_i, a_j), a_k)$ .

Cette propriété tient à la définition de la cohésion comme moyenne géométrique et à la construction précédente de la classe  $(a_i, a_j)$  :

$$c(a_i, a_j) \geq [c(a_i, a_j) c(a_j, a_k) c(a_i, a_k)]^{\frac{1}{3}}.$$

Finalement, lors de cette étape, on agrège une paire ou une classe à 3 éléments, mais simultanément, la cohésion de la nouvelle classe n'est pas meilleure que lors de l'étape précédente.

**3<sup>ème</sup> étape** : on compare les cohésions obtenues selon les phases  $\mathcal{O}_2$ ,  $\mathcal{O}_3$  et  $\mathcal{O}_4$  suivantes et on conserve la classe correspondant à la plus forte.

$\mathcal{O}_2$  : on détermine un couple "**maximal**" parmi les couples restants ;

$\mathcal{O}_3$  : on compare la cohésion de ce couple à celle des classes à 3 éléments non encore tous réunis, mais dont 2 d'entre eux l'ont déjà été ;

$\mathcal{O}_4$  : on compare la cohésion maximale des 2 phases à celle des classes à 4 éléments élargissant la classe à 3 éléments éventuellement constituée lors de la 2<sup>ème</sup> étape en réunissant les 2 classes à 2 éléments déjà formés. Si tel est le cas avec  $\{a_i, a_j, a_k, a_\ell\}$ , alors :  $C((a_i, a_j), a_k) \geq C(((a_i, a_j), a_k), a_\ell)$  et a fortiori :  $c(a_i, a_j) \geq C(((a_i, a_j), a_k), a_\ell)$ .

L'argument est le même que précédemment (cf. moyenne géométrique). La cohésion obtenue à l'issue des phases  $\mathcal{O}_2$ ,  $\mathcal{O}_3$  et  $\mathcal{O}_4$  n'est pas meilleure que la précédente.

**pi<sup>ème</sup> étape**: avant cette étape, ont été agrégées éventuellement des classes à :

- . 2 éléments, construites pas à pas selon une cohésion décroissante ;
- . 3 éléments " " " " " " " " " "
- et non meilleure que celle des couples les ayant générées ;
- . 4 éléments, construites pas à pas selon une cohésion décroissante et non meilleure que celle des classes à 3 éléments les ayant générées ;
- . .....
- . p éléments selon les mêmes critères.

Selon les mêmes phases que précédemment, si une classe issue d'un ensemble à  $(m+1)$  éléments est constituable à cette étape, c'est qu'elle réalise une cohésion maximale parmi toutes celles à  $q \leq m$  éléments constituables à cette étape, compte tenu de celles déjà constituées. De plus, la cohésion de cette classe à  $(m+1)$  éléments n'est pas meilleure que celle de la classe à  $m$  éléments qu'on veut étendre. (Si une classe à  $m+r$  éléments est constituable, l'argument est valide a fortiori).

En conclusion, d'une étape à l'autre, la cohésion des classes formées est toujours décroissante au sens large. De plus, le processus est fini puisqu'à chaque étape, on accroît l'effectif d'une classe d'au moins un élément qui ne sera plus isolé. Ce processus prend fin lorsque toute nouvelle classe constituable admet une cohésion nulle ou bien lorsque la classe ultime est constituée de tous les éléments.

Un étudiant en cours de thèse à Rennes, Saddo AG ALMOULOUUD élabore sur P.C. et Macintosh, et à partir de réalisations informatiques déjà opérationnelles sur P.C., des programmes d'analyses de données incluant :

- . le calcul d'intensités d'implication entre attributs ou variables binaires et entre variables non binaires et variables numériques ;

. le calcul des cohésions et d'indices d'implication entre classes selon l'algorithme précédent ;

. la construction de l'arbre d'implication entre classes sur le modèle algorithmique utilisé par R. GRAS pour l'arbre de classification hiérarchique selon l'A.V.L. de I.C. LERMAN.

Ces programmes ont servi aux traitements de données dont il sera question ultérieurement. Nous renvoyons à la thèse pour de plus amples informations sur les logiciels en question.

### 3.4.2. Consistance de classe et minimalité.

On dira qu'à un niveau donné (lors d'une étape donnée), une classe  $\mathcal{C}$  **consistante** vient de se former par agrégation de la classe  $\mathcal{C}_2$  à  $\mathcal{C}_1$  (éventuellement singleton) s'il existe au moins une classe  $\mathcal{C}_0$  telle que la relation implicative entre classes soit améliorée :

$$\mathcal{C} = \mathcal{C}_1 \cup \mathcal{C}_2$$

$$\psi(\mathcal{C}_1, \mathcal{C}_0) < \psi(\mathcal{C}_1 \cup \mathcal{C}_2, \mathcal{C}_0) \text{ ou } \psi(\mathcal{C}_0, \mathcal{C}_1) < \psi(\mathcal{C}_0, \mathcal{C}_1 \cup \mathcal{C}_2).$$

De plus, soit  $\mathcal{C}'$  une extension quelconque de  $\mathcal{C}$ .

$$\text{Si } \forall \mathcal{C}_0 \subset \mathcal{C}, \psi(\mathcal{C}_0, \mathcal{C}') \leq \psi(\mathcal{C}_0, \mathcal{C})$$

$$\text{ou } \psi(\mathcal{C}', \mathcal{C}_0) \leq \psi(\mathcal{C}, \mathcal{C}_0),$$

$\mathcal{C}$  est dite **consistante minimale** (elle est minimale par rapport à son effectif d'éléments constituants).

### 3.4.3. Exemple.

$$\left\{ \begin{array}{l} N=100, \text{Card } A=n_a=7, \text{Card } B=n_b=16 \\ \text{Card } C=n_c=20, \text{Card } D=n_d=25 \\ \text{Card } E=n_e=40, \text{Card } F=n_f=95 \\ n_a < n_b < n_c < n_d < n_e < n_f \end{array} \right.$$

et

$$\text{Card}(A \cap \bar{B}) = 2, \text{Card}(A \cap \bar{C}) = 7, \text{Card}(A \cap \bar{D}) = 7, \text{Card}(A \cap \bar{E}) = 4, \text{Card}(A \cap \bar{F}) = 0$$

$$\text{Card}(B \cap \bar{C}) = 16, \text{Card}(B \cap \bar{D}) = 15, \text{Card}(B \cap \bar{E}) = 12, \text{Card}(B \cap \bar{F}) = 0$$

$$\text{Card}(C \cap \bar{D}) = 20, \text{Card}(C \cap \bar{E}) = 5, \text{Card}(C \cap \bar{F}) = 0$$

$$\text{Card}(D \cap \bar{E}) = 25, \text{Card}(D \cap \bar{F}) = 1$$

$$\text{Card}(E \cap \bar{F}) = 0.$$

En conclusion, l'implication statistique entre variables, prolongées en implication entre classes de variables nous fournit un outil d'étude de caractères de nature très variée. Elle semble utile au didacticien, au psychologue et de façon générale à tout chercheur disposant de données où les liaisons dans une population sont floues mais structurables en arbre puis en classes orientées. Des questions en termes de genèse, de complexité, de conduites nécessaires, etc., peuvent y trouver réponse. Nous cherchons à le montrer dans le paragraphe suivant.

### **§ 3 - APPLICATIONS DE L'IMPLICATION STATISTIQUE A L'ANALYSE DIDACTIQUE DE QUESTIONNAIRES.**

L'étude présentée ici vise à connaître l'origine et le nature des erreurs les plus fréquentes, à analyser les procédures utilisées par des élèves de 1<sup>er</sup> cycle mis en situation de démonstration en géométrie, en particulier dans le cas où l'activité déductive se réduit à une simple inférence. Nous entreprenons, sur des situations réelles - réponses d'élèves à deux questionnaires -, les traitements statistiques des données recueillies, suivant les deux méthodes d'analyse présentées dans le § 1 : la classification hiérarchique (selon I.C. LERMAN) et la classification implicative. Nous dégageons ainsi quelques grandes classes de comportements erronés entre lesquelles nous tentons de nouveau d'établir des relations implicatives suivant le processus élaboré et développé dans le § 2.

#### **3.1. Présentation des questionnaires.**

Ces questionnaires concernent le début de l'apprentissage de la démonstration (classe de 5<sup>ème</sup>). Ils sont réalisés à partir d'un logiciel : le logiciel "PREMIER PAS", conçu au sein de l'équipe de didactique de Rennes pour aider l'enseignant à repérer et analyser les erreurs commises par l'élève dans les démonstrations à un pas. Ce logiciel propose à l'élève une liste de faits et une liste de théorèmes repérés par des numéros. Les questions concernent des inférences simples : hypothèse(s) - théorème - conclusion, présentant une ou plusieurs lacunes en fournissant les numéros des faits ou théorèmes appropriés.

Les questionnaires, appelés "6 questions" et "5 questions" que nous étudions ici, se rapportent à la symétrie centrale et, pour une question du 1<sup>er</sup> (Q5), à la transitivité du parallélisme. Les questions posées visent à étudier l'effet d'un certain nombre de variables liées à la logique de l'inférence et à la forme des énoncés proposés. Hypothèses et théorèmes sont donnés : l'élève doit compléter en choisissant un des faits de la liste à titre de conclusion.

## FAITS

- 1  $(EF)$  et  $(CD)$  sont symétriques par rapport au point  $I$
- 2  $[MN]$  est le symétrique de  $[PR]$  par rapport au point  $I$
- 3  $(AB)$  et  $(CD)$  sont symétriques par rapport au point  $O$
- 4  $(MN) \parallel (PR)$
- 5  $(CD) \parallel (EF)$
- 6  $(AB) \parallel (CD)$
- 7  $(AB) \parallel (EF)$
- 8  $MN = PR$
- 9  $CD = EF$
- 10  $AB = CD$
- 11  $AB = EF$

## THEOREMES

- 1 La symétrie centrale conserve les longueurs.
- 2 Si  $(D) \parallel (D')$  et  $(D') \parallel (D'')$  alors  $(D) \parallel (D'')$ .
- 3 Le symétrique d'une droite  $(D)$  par rapport à un point est une droite  $(D')$  parallèle à  $(D)$ .
- 4 Si deux droites sont symétriques par rapport à un point alors elles sont parallèles.
- 5 Deux segments symétriques par rapport à un point ont même longueur.
- 6 La symétrie centrale conserve les directions.

En fait, chaque question se présente schématiquement ainsi :

**Hypothèse :** fait n° p

**Théorème :** n° q

**Conclusion :** fait n° ?

Schématiquement, l'ensemble questions-réponses peut être présenté ainsi :

	HYPOTHESES	THEOREME	CONCLUSION à trouver
Q <sub>1</sub> {	Hypothèse : 1 Théorème : 3 → Concluion : 5	$(EF)$ et $(CD)$ symétriques par rapport à $I$	Le symétrique de $(D)$ par rapport à un point est $(D') \parallel (D)$  $(EF) \parallel (CD)$
Q <sub>2</sub> {	Hypothèse : 3 Théorème : 4 → Concluion : 6	$(AB)$ et $(CD)$ symétriques par rapport à $O$	Si 2 droites sont symétriques par rapport à un point alors elles sont parallèles  $(AB) \parallel (CD)$
Q <sub>3</sub> {	Hypothèse : 2 Théorème : 5 → Concluion : 8	$[MN]$ est symétrique de $[PR]$ par rapport à $I$	2 segments symétriques par rapport à un point ont même longueur  $MN = PR$
Q <sub>4</sub> {	Hypothèse : 3 Théorème : 6 → Concluion : 6	$(AB)$ et $(CD)$ symétriques par rapport à $O$	La symétrie centrale conserve les directions  $(AB) \parallel (CD)$
Q <sub>5</sub> {	Hypothèse : 6 et 5 Théorème : 2 → Concluion : 7	$(AB) \parallel (CD)$ et $(CD) \parallel (EF)$	Si $(D) \parallel (D')$ et $(D') \parallel (D'')$ alors $(D) \parallel (D'')$  $(AB) \parallel (EF)$
Q <sub>6</sub> {	Hypothèse : 2 Théorème : 1 → Concluion : 8	$[MN]$ est symétrique de $[PR]$ par rapport à $I$	La symétrie centrale conserve les longueurs  $MN = PR$

A la suite d'un 1<sup>er</sup> essai, l'élève est autorisé à faire un 2<sup>ème</sup> et dernier essai. Les questions sont indépendantes. Le questionnaire "5 questions" soumis aux élèves postérieurement au "6 questions" ne diffère de celui-ci que : - par la suppression de la question (Q<sub>5</sub>) non relative à la symétrie ponctuelle (il est en effet apparu que les réponses au "6 questions" étaient un peu biaisées par la présence de cette question relative à un autre concept) - et par l'ajout de 2 faits dont l'absence dans le "6 questions" semble avoir provoqué des confusions trop singulières :

$$12. (AB) = (CD)$$

$$13. [MN] = [PR].$$

Chaque modalité de réponse est codée par un triplet.

**Exemple.** 3-6-10 (Q<sub>4</sub>).

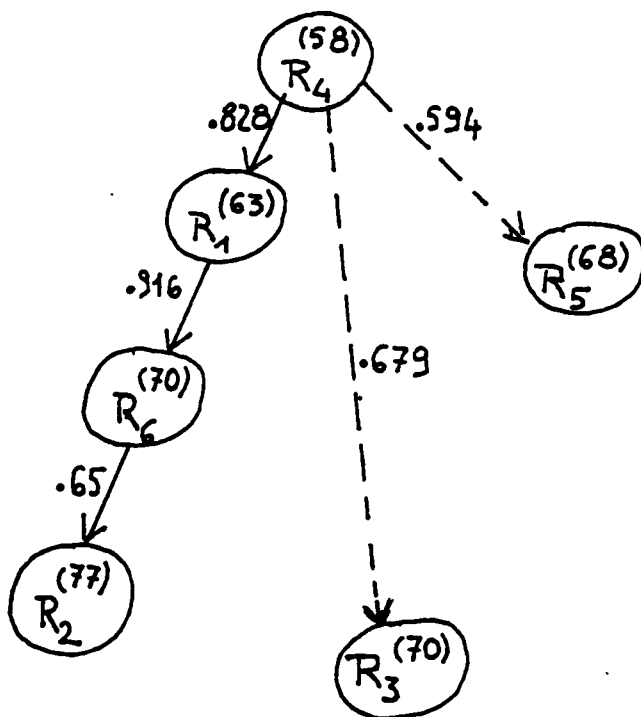
Hypothèse : (AB) et (CD) sont symétriques par rapport à 0.

Théorème : la symétrie centrale conserve les directions.

Conclusion donnée par l'élève :  $AB = CD$ .

### 3.2. Analyse hiérarchique et implicative des réponses au "6 questions".

Après calcul des indices de similarité et des intensités d'implication entre les 31 modalités de réponse prises 2 à 2 (80 élèves de 5<sup>ème</sup>), nous avons obtenu l'arbre hiérarchique suivant (cf. p. 23) :



Ci-contre, nous présentons le graphe implicatif entre les 6 modalités de réussite.



Dans l'arbre hiérarchique des similarités, tous les attributs s'agrègent en classes, les classes en sur-classes jusqu'à la classification complète. Nous voyons d'abord se regrouper les procédures erronées puis, à partir du niveau 20 seulement, toutes les réussites ; celles-ci ne se rejoignent entre elles qu'à l'avant-dernier niveau : dans ce questionnaire, la stabilité dans l'erreur semble donc plus importante que dans les réussites. 4 classes, également très stables, apparaissent nettement : les réussites précisément et 3 classes de procédures erronées que nous analyserons plus loin, mais remarquons dès à présent la grande proximité des réponses 1-3-1, 3-4-3, 2-5-2 et 3-6-3 qui sont répétition stricte de l'hypothèse en conclusion de l'exercice.

**Exemple. 3-6-3 ( $Q_4$ ) :** réponse obtenue 17 fois.

Hypothèse :  $(AB)$  et  $(CD)$  sont symétriques par rapport à 0.

Théorème : la symétrie centrale conserve les directions.

Conclusion :  $(AB)$  et  $(CD)$  sont symétriques par rapport à 0.

Ce type surprenant de réponse est sans doute lié à la forme du questionnaire, à une mauvaise compréhension de la consigne ou du logiciel "PREMIER PAS" ou des mots difficiles "directions" et "conserver" (avec le sens donné par l'enseignant dans sa classe) figurant dans le théorème.

Le graphe implicatif des réussites du "6 questions" traduit une relation de préordre partiel dans cet ensemble. Il est donc orienté et valué.  $R_4$ , réussite à la question ( $Q_4$ ) la plus complexe du questionnaire (et la moins bien réussie), en est la source.  $R_2$ ,  $R_3$ ,  $R_5$  en sont les puits :

$R_2$ , réussite à la question ( $Q_2$ ) dont le théorème est exprimé en "si .... alors", formulation facilitant, semble-t-il, le succès.

$R_3$ , relative à la conservation des longueurs dans la symétrie centrale. La longueur est une notion plus familière à l'élève que celle de direction.

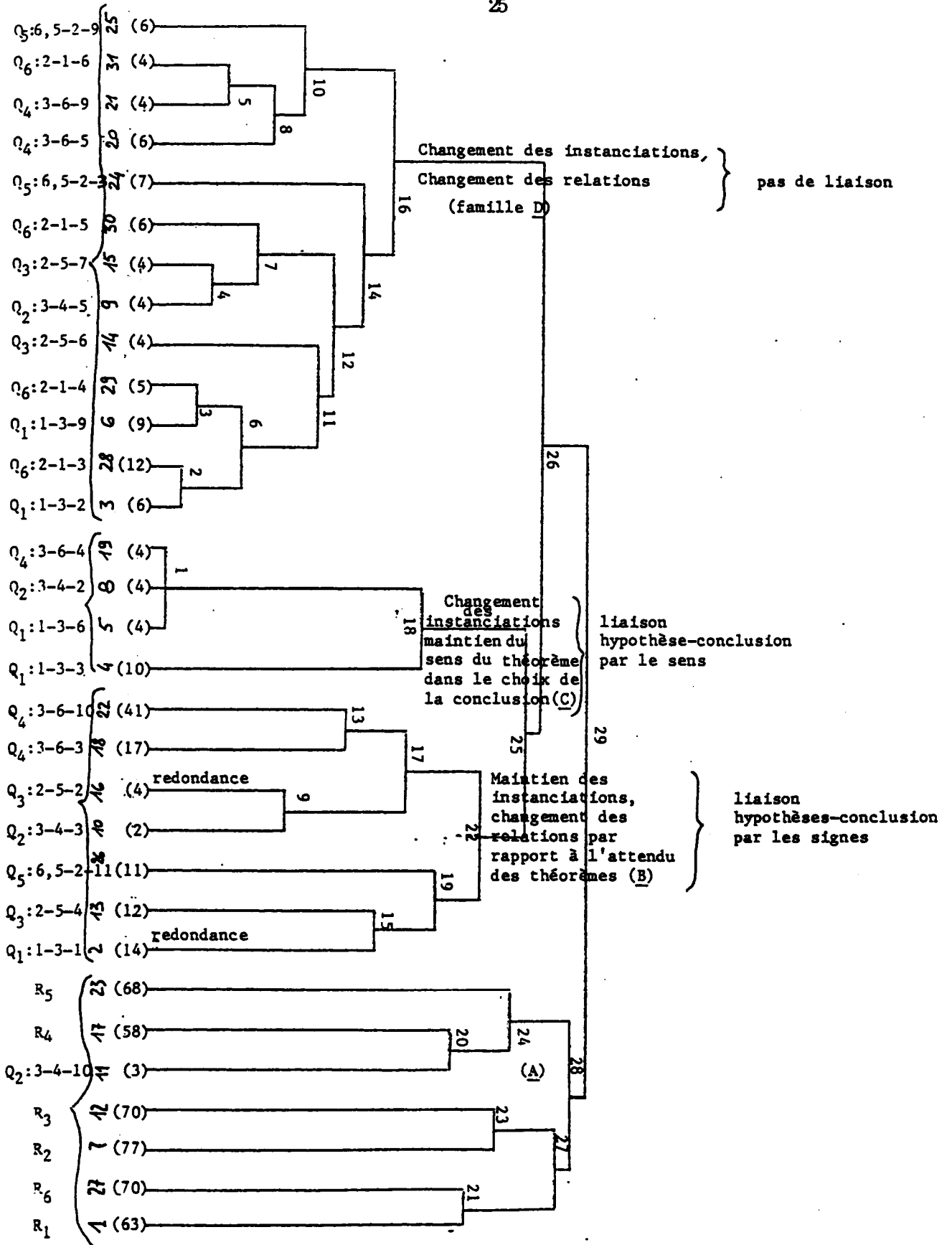
$R_5$ , seule question relative à la transitivité du parallélisme.

Les valeurs indiquées sont les intensités des implications.

Sur le même chemin, les réussites se placent dans l'ordre croissant de leurs effectifs.

Ainsi, dans le "6 questions",  $R_1 \Rightarrow R_2$  avec une intensité d'implication très forte : 0,938. Les questions  $Q_1$  et  $Q_2$  ne diffèrent que par l'expression de leur théorème. Celui de ( $Q_2$ ), en si....alors, facilite la réussite à cette question.

$R_4 \Rightarrow R_5$  très faiblement, avec une intensité d'implication de 0,594 ;  $R_5$  est la seule question non relative à la symétrie centrale et contenant une double hypothèse



Arbre des similarités 80 élèves

$R_3$  et  $R_6$ , réussites aux questions portant sur la conservation des longueurs dans une symétrie (propriété métrique), proches dans l'arbre des similarités ne sont pas liées implicativement, sans doute en raison de la présence du mot "conserver" dans le théorème de (Q6). Nous avons réfuté leur implication car la valeur obtenue pour l'intensité était inférieure à 0,5.

### 3.3. Analyse des modalités de réponses au questionnaire "6 questions" (52 élèves de 5<sup>ème</sup>).

Question par question, nous avons identifié et codé les procédures susceptibles d'apparaître. Nous les présentons ci-dessous :

A : bonne réponse

R : répétition

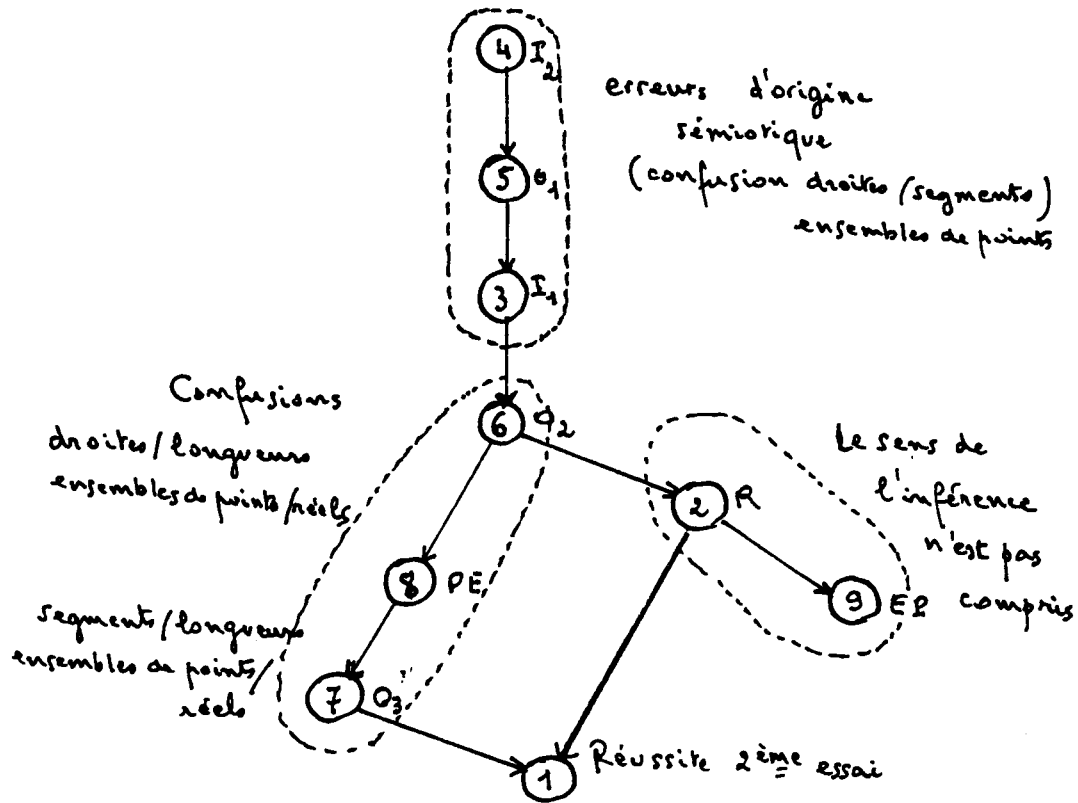
Changement d'instanciation	$\begin{cases} I_1 : \text{changement total d'instanciation} \\ I_2 : \text{changement partiel d'instanciation} \end{cases}$
Confusion entre objets	$\begin{cases} O_1 : \text{confusion entre droites et segments} \\ O_2 : \text{confusion entre droites et longueurs} \\ O_3 : \text{confusion entre segments et longueurs} \end{cases}$
Confusion entre relations	$\begin{cases} PE : // \text{remplacé par } = \\ EP : = \text{ remplacé par } // . \end{cases}$

Ces procédures sont spécifiques du questionnaire et permettent donc de définir les compteurs non-standard de notre analyse.

Variables	Nombre maximal d'apparitions chez un élève	Effectifs
1. Nombre de bonnes réponses au 2ème essai	5	66
2. Répétition	10	70
3. Changement total d'instanciation	10	29
4. Changement partiel d'instanciation	10	16
5. Confusion : droites/segments	6	10
6. Confusion : droites/longueurs	10	70
7. Confusion : segments/longueurs	4	49
8. // remplacé par = (PE)	6	64
9. = remplacé par // (EP)	4	31

Nous avons construit (cf. 1.3 p. 7) une méthode permettant de donner du sens et d'attacher une intensité à la quasi-implication d'une variable numérique ou fréquentielle sur une autre.

Ainsi, l'arbre implicatif obtenu après deux essais des élèves est le suivant :



Nous trouvons à l'origine de l'arbre les deux variables 4 et 5 correspondant, pensons-nous, à deux erreurs anodines :

- la variable 4 traduit un changement d'instanciation mais seulement partiel
- la variable 5 exprime la confusion entre droites et segments mais une droite et un segment sont tous deux des ensembles de points.

Le triplet {4,5,3} implique la variable 6 d'où partent deux ramifications qui se réunissent à la base de l'arbre en la réussite au 2<sup>ème</sup> essai (variable 1).

- à droite : seule la variable 2 (répétition) conduit à la réussite au 2<sup>ème</sup> essai ; ceci conforte les hypothèses avancées à différents endroits de la thèse d'A. LARHER sur ces procédures,
- à gauche, dans cet ordre : les variables 6 ( $O_2$  : confusion droites/longueurs), 8 (// remplacé par =) et 7 ( $O_3$  : confusion segments/longueurs).

. l'arbre précédent marque une différence entre les 3 confusions d'objets que nous avons codées  $O_1$ ,  $O_2$  et  $O_3$  :  $O_2$  et  $O_3$  sont deux confusions entre un ensemble de points (droite ou segment) et un nombre ; elles sont à distinguer en cela de  $O_1$ , confusion entre deux ensembles de points dont nous avons parlé plus haut ;  
 .  $O_2$  et  $O_3$  sont séparées dans l'arbre par la variable 8 et cette disjonction peut s'interpréter comme suit : à tout segment on peut attacher un réel, sa longueur, et le remplacement de l'un par l'autre (confusion  $O_3$  dans les égalités des

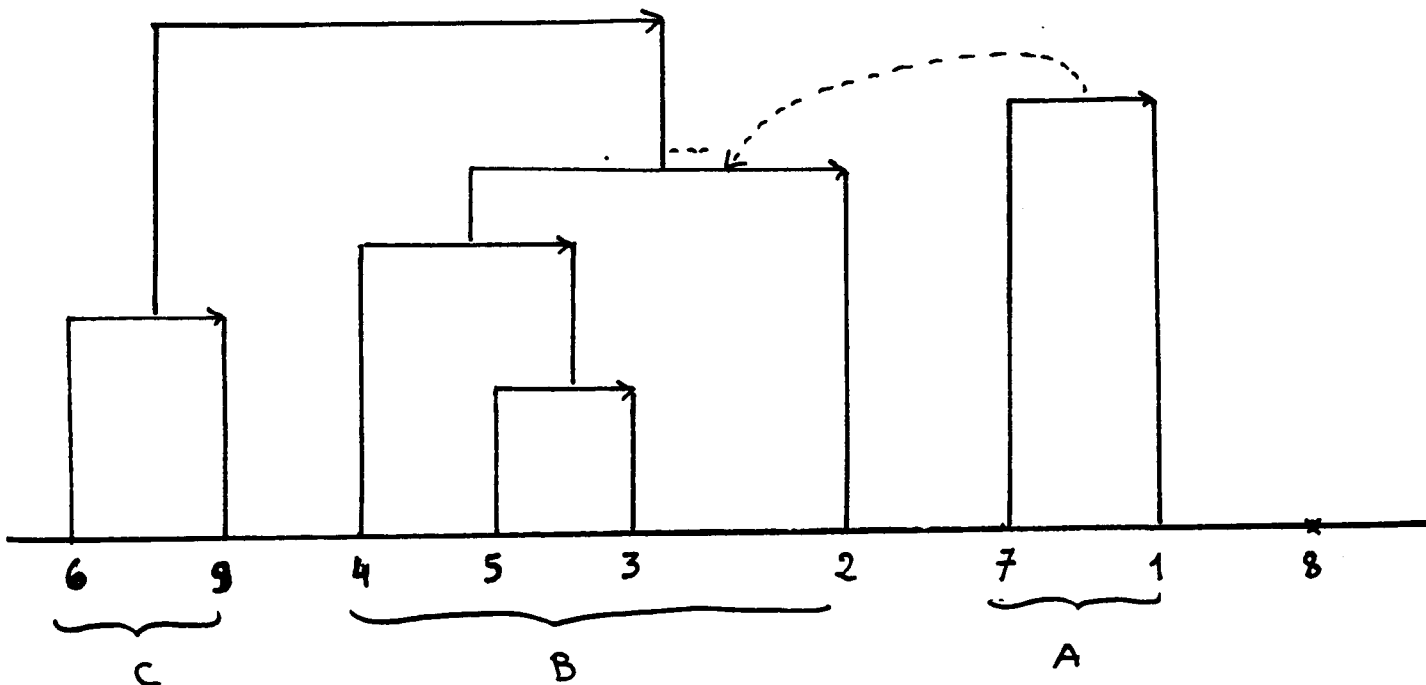
et le remplacement de l'un par l'autre (confusion  $\mathcal{O}_3$  dans les égalités des questions ( $Q_4$ ) et ( $Q_5$ )) relève certainement plus d'une erreur de notation que de sens ; ainsi la variable 7, correspondant à  $\mathcal{O}_3$ , implique le plus fortement la variable 1 (réussite au 2<sup>ème</sup> essai) ( $\varphi = 0,72$ ) alors que l'intensité d'implication de la variable 6 ( $\mathcal{O}_2$ ) vers la réussite est quasiment nulle ;

. c'est la variable 8 (// remplacé par =) qui sépare  $\mathcal{O}_2$  et  $\mathcal{O}_3$  ; nous avons la confirmation ici que la lecture de ce type de confusion doit se faire à deux niveaux :

- confusion de type sémiotique : le signe // se confond avec le signe =
- mais surtout la confusion de type sémantique : la relation = se substitue à la relation //.

Cette interprétation est elle-même confortée par la quasi-nullité de l'intensité d'implication de 9 vers 1 (réussite).

Les calculs de cohésion et l'application de l'algorithme de constitution de classe nous font obtenir la classification suivante :



Elle partitionne l'ensemble des 9 variables en 3 classes disjointes, ce qui n'est pas le cas de la classification hiérarchique des similarités où l'emboîtement se produit dès le 5<sup>ème</sup> niveau.

La classe centrale  $B$ , constituée des deux sous-classes (4,(5,3)) et (3) met en évidence l'origine majeure des erreurs au 1<sup>er</sup> essai : le changement d'instanciation, la confusion droites/segments et la répétition, erreur qui peut comme les deux premières ne pas être persistante. Est encore moins persistante l'erreur de notation représentée par la variable 7 puisque celle-ci se referme avec 1 (réussite au 2<sup>ème</sup> essai) en une classe  $A$  très cohérente.

Remarquons le regroupement  $C$  des 2 variables 6 et 9 qui, par contre, renforce notre hypothèse d'erreur persistante car conceptuelle. La résistance à l'association de 8 à (6,9), révélée par la très faible cohésion, souligne, à son tour, la disjonction des procédures représentées par 8 (// remplacé par =) et 9 (= remplacé par //). En cela, la hiérarchie implicative nous informe de façon différente de la hiérarchie des similarités, relativement linéaire ici dans ses emboîtements et, par suite, moins pertinente pour l'analyse cognitive que nous avons menée. Ce jugement n'affecte en rien les qualités informationnelles de la méthode dans d'autres cas comme nous l'avons constaté sur l'arbre des similarités du questionnaire "6 questions".

### 3.4. Synthèse.

Les analyses précédentes ont mis en évidence, dans les questionnaires "6 questions" et "5 questions" des types d'erreurs faites par de jeunes élèves en situation d'apprentissage de la démonstration mathématique ; à savoir 3 grandes familles de procédures erronées, très stables, mais non disjointes, se sont dégagées :

- changement d'instanciation : l'élève change les noms des objets
- répétition : l'élève répète en conclusion l'hypothèse de l'inférence
- changement de relation par rapport à l'attendu du théorème ou, plus spécifiquement, confusion entre le parallélisme et l'égalité.

L'examen des classifications hiérarchiques des similarités et des implications conduit à une épistémologie artificielle de l'inférence :

- 2 formes primitives de l'inférence :
  - . la répétition ou tautologie inféconde,
  - . le changement des noms des objets avec changement ou non de la relation attendue ;
- 1 forme plus évoluée où les relations attendues sont échangées avec des relations voisines, sans symétrie entre ces échanges ;
- 1 forme presque achevée où la réponse diffère de l'attendu par la seule écriture (confusion entre signifiants seuls) ;
- 1 forme achevée conduisant à la réussite.

#### § 4 - CONCLUSION.

Succédant à des mises à l'épreuve depuis plus de dix ans dans différentes recherches en didactique et en psychologie, les résultats obtenus à travers ces différents questionnaires, contrôlés localement par d'autres méthodes d'analyses de données (analyse factorielle et analyse hiérarchique), confortent l'approche implicative que nous avons adoptée. Par son caractère non symétrique, elle apporte une nouvelle synthèse dynamique de données, tout en complétant avantageusement les informations fournies par d'autres méthodes. L'extension à l'analyse implicative de classes présente un intérêt majeur par sa capacité à condenser les relations plus fines mais trop enserrées dans un réseau complexe. Des problèmes théoriques restent encore en suspens. Nul doute que l'examen des lois des intensités d'implication entre classes apportera une très intéressante information sur leur propre qualité. C'est, en particulier, l'objet de nouvelles autres recherches de notre équipe.

#### REFERENCES

- [ACID S., de CAMPOS L.M., GONZALEZ A., MOLINA R., PEREZ de la BLANCA N. 1991] - Learning with Castle - in Symbolic and quantitative approaches to uncertainty (R. KRUSE, P. SIEGEL), Springer-Verlag, 99-106.
- [AMARGE S., DUBOIS D., PRADE H. 1991], Imprecise quantifiers and conditional probabilities - in Symbolic and quantitative approaches to uncertainty (R. KRUSE, P. SIEGEL), Springer-Verlag, 33-37.
- [DIDAY E. 1991] - Towards a statistical theory of intentions for knowledge analysis, rapport de recherche 1494, INRIA Rocquencourt.
- [GAMMERMAN A., LUO Z. 1991] - Constructing Causal Trees from a medical database, Technical Report TR 91 002, Dep<sup>t</sup> of Computer Sci., Heriot-Watt Univ., Edimburgh.
- [GRAS R., 1979] - Contribution à l'étude expérimentale et à l'analyse de certaines acquisitions cognitives et de certains objectifs didactiques en mathématiques, Thèse d'Etat, Université de Rennes I, octobre 1979.
- [GRAS R. et LARHER A., 1989] - La quasi-implication : une méthode d'analyse de relations non symétriques entre attributs et entre classes d'attributs, Public. interne I.R.M.A.R., Rennes, 1989.
- [LARHER A., 1991] - Implication statistique et applications à l'analyse de démarches de preuve mathématique, Thèse de l'Université de Rennes I, février 1991.
- [LERMAN I.C., GRAS R., ROSTAM H., 1981] - Elaboration et évaluation d'un indice d'implication pour des données binaires, I et II, Mathématiques et Sciences Humaines n° 74, p 5-35 et n° 75, p 5-47, 1981.
- [LERMAN I.C., 1981] - Classification et analyse ordinaire des données, Dunod, 1981.
- [LOEVINGER J. 1947] - A systematic approach to the construction and evaluation of tests of ability, Psychological Monographs, 61, n° 4.
- [PEARL J. 1988] - Probabilistic Reasoning in intelligent systems, San Mateo, CA, Morgan Kaufmann.