



HAL
open science

Analyse de la forme d'un coefficient d'association entre variables qualitatives

Mohamed Ouali Allah

► **To cite this version:**

Mohamed Ouali Allah. Analyse de la forme d'un coefficient d'association entre variables qualitatives. [Rapport de recherche] RR-1366, INRIA. 1991. inria-00075194

HAL Id: inria-00075194

<https://inria.hal.science/inria-00075194>

Submitted on 24 May 2006

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

INRIA

UNITÉ DE RECHERCHE
INRIA-RENNES

Institut National
de Recherche
en Informatique
et en Automatique

Domaine de Voluceau
Rocquencourt
B.P.105
78153 Le Chesnay Cedex
France
Tél.: (1) 39 63 55 11

Rapports de Recherche

N° 1366

Programme 6
Calcul scientifique, Modélisation et
Logiciels numériques

ANALYSE DE LA FORME D'UN COEFFICIENT D'ASSOCIATION ENTRE VARIABLES QUALITATIVES

Mohamed OUALI - ALLAH

Janvier 1991



ANALYSE DE LA FORME D'UN COEFFICIENT D'ASSOCIATION ENTRE VARIABLES QUALITATIVES

Mohamed OUALI ALLAH

Publication Interne n° 554 - Octobre 1990 - 26 Pages

- Programme 6 -

Résumé

Dans le but de comparer deux variables qualitatives, nous les représentons au moyen d'un codage sur l'ensemble des couples d'objets. Nous considérons un coefficient général d'association basé sur une normalisation de nature combinatoire et statistique (centrage et réduction), par rapport à une hypothèse d'absence de lien.

L'objet de ce travail est de fournir des nouvelles formes de cet indice ayant un caractère très synthétique. Nous distinguerons notamment les cas où le codage est symétrique ou antisymétrique. Ces expressions permettront d'appréhender aisément le comportement asymptotique du coefficient.

Mots-clés : - Coefficient d'association - Variable qualitative - Codage - Normalisation - Hypothèse d'absence de lien - Forme limite

ANALYSIS OF THE FORM OF AN ASSOCIATION COEFFICIENT BETWEEN QUALITATIVE VARIABLES

Abstract

For comparing two qualitative variables, we represent them with a coding on object-couples. We consider a general association coefficient based on a combinatory and statistical normalization with respect to the independence hypothesis.

The aim of this work is to provide new simplified forms for this coefficient. In particular we distinguish the symmetrical and the skew-symmetrical coding cases. These new expressions give an insight into the asymptotic behaviour of the coefficient.

Key words : - Association coefficient - Qualitative variable - Coding - Normalization - Non-linkage hypothesis - Limit form

Table des matières

| | | |
|----------|---|-----------|
| 1 | Introduction | 4 |
| 2 | Un coefficient général d'association | 5 |
| 2.1 | Centrage-réduction | 6 |
| 2.1.1 | Calcul de l'espérance | 6 |
| 2.1.2 | Calcul de la variance | 6 |
| 2.2 | Symétrie et antisymétrie | 7 |
| 3 | Expressions du coefficient centré-réduit | 10 |
| 3.1 | Expression ϕ | 11 |
| 3.1.1 | Cas symétrique | 11 |
| 3.1.2 | Cas antisymétrique | 12 |
| 3.2 | Expression ψ | 13 |
| 3.2.1 | Cas symétrique | 13 |
| 3.2.2 | Cas antisymétrique | 15 |
| 4 | Comportement asymptotique | 16 |
| 4.1 | Signes de ϕ et ψ | 17 |
| 4.2 | Forme limite | 18 |
| 4.3 | Réduction géométrique | 19 |
| 5 | Conclusion | 21 |
| | Bibliographie | 22 |

1 Introduction

Nous nous situons ici dans un contexte classique de la comparaison deux à deux d'un ensemble de variables qualitatives observées sur un ensemble \mathcal{O} d'objets.

Depuis les travaux de M.G. KENDALL [6] de nombreux chercheurs s'accordent sur la représentation d'une variable qualitative au moyen d'une relation binaire sur \mathcal{O} qui peut être évaluée : L.J. HUBERT [3] [5] , I.C. LERMAN [9] [11] , F. MARCOTORCHINO & MICHAUD [14] , B. MONJARDET & J.P. BARTHELEMY [15] ...

Nous considérons donc un coefficient général d'association entre deux relations "valuées", dont la normalisation statistique suppose le calcul de l'espérance et surtout de la variance. N. MANTEL dans [13], et I.C. LERMAN dans [12], avec des approches différentes, aboutissent à une expression générale du coefficient centré-réduit.

Nous reprenons d'une façon plus analytique les calculs de normalisation, et nous développons des expressions plus synthétiques de la variance; chacune de ces expressions ayant son intérêt propre. Nous distinguons notamment le traitement du cas d'un codage symétrique de celui d'un codage antisymétrique.

La relative simplicité des expressions obtenues du coefficient centré-réduit, permet d'étudier aisément son comportement asymptotique, et d'aboutir à une forme limite de même type que dans [12], c'est à dire un rapport pur à un facteur en n près (n étant le nombre d'objets).

2 Un coefficient général d'association

Notre étude porte sur la situation générale envisagée dans [7], de la comparaison de deux codages ou "valuations" sur l'ensemble $\mathcal{O} \times \mathcal{O}$ des couples d'objets. Ainsi chaque variable v définit un codage $\{c_{ij}^v / (i, j) \in I \times I\}$ sur $\mathcal{O} \times \mathcal{O}$. Où $I = \{1, 2, \dots, i, \dots, n\}$ est l'ensemble d'indéxation de \mathcal{O} .

Le point de vue que nous considérons ici est celui introduit par I.C. LERMAN où l'hypothèse d'absence de liaison permet la découverte de l'expression formelle de coefficients d'association entre variables qualitatives .

Relativement à la comparaison de deux variables qualitatives v et w , on introduit l'indice appelé "brut" :

$$(1) \quad s(v, w) = \sum_{i \neq j} c_{ij}^v c_{ij}^w$$

L'hypothèse d'absence de lien, on dit aussi d'indépendance, à caractère permutatif, consiste à associer à l'indice brut (1) deux variables aléatoires duales :

$$(2) \quad S = s(v, w^*) = \sum_{i \neq j} c_{ij}^v c_{i^*j^*}^w$$

$$S' = s(v^*, w) = \sum_{i \neq j} c_{i^*j^*}^v c_{ij}^w$$

où $i^* = \tau(i)$ et $j^* = \tau(j)$, τ étant une permutation aléatoire dans l'ensemble - muni d'une probabilité uniforme - des permutations sur l'ensemble I .

L'étude de la variable aléatoire S (resp. S') a intéressé de nombreux statisticiens et ce, dans différents contextes. Le plus classique est celui des tests d'indépendance :

P. ARABIE [1], L.J. HUBERT [2] [4], G. LE CALVE [7], N. MANTEL [13], P.W. MIELKE [16] ...

Les variables aléatoires S et S' ont la même distribution (cf. [7]). Celle-ci, dans la presque totalité des cas, suit asymptotiquement une loi normale (cf. [10],[11]). En notant $\mathcal{E}(S)$ l'espérance et $\sigma^2(S)$ la variance de cette variable aléatoire, nous pouvons définir le coefficient centré-réduit :

$$(3) \quad Q(v, w) = \frac{s(v, w) - \mathcal{E}(S)}{\sigma(S)}$$

2.1 Centrage-réduction

2.1.1 Calcul de l'espérance

La variable aléatoire S définie dans (2) a pour espérance :

$$\mathcal{E}(S) = \sum_{i \neq j} c_{ij}^v \mathcal{E}[c_{i^*j^*}^w]$$

or :

$$\mathcal{E}[c_{i^*j^*}^w] = \sum_{i' \neq j'} c_{i'j'}^w P_{i'j'}$$

où : $P_{i'j'}$ est la probabilité pour que $i^* = i'$ et $j^* = j'$

En désignant par $n^{[p]} = n(n-1)\dots(n-p+1)$ la $p^{\text{ème}}$ puissance factorielle de n , on a :

$$P_{i'j'} = \frac{1}{n^{[2]}}$$

D'où l'expression de l'espérance :

$$(4) \quad \mathcal{E}(S) = \frac{1}{n^{[2]}} \sum_{i \neq j} c_{ij}^v \sum_{i' \neq j'} c_{i'j'}^w$$

2.1.2 Calcul de la variance

Pour le calcul de $\sigma^2(S) = \mathcal{E}(S^2) - \mathcal{E}^2(S)$, développons d'abord S^2 :

$$\begin{aligned} S^2 = & \sum_{i \neq j} c_{ij}^v{}^2 c_{i^*j^*}^w{}^2 + \sum_{i \neq j} c_{ij}^v c_{ji}^v c_{i^*j^*}^w c_{j^*i^*}^w + \sum_{i \neq j \neq k} c_{ij}^v c_{ik}^v c_{i^*j^*}^w c_{i^*k^*}^w \\ & + \sum_{i \neq j \neq k} c_{ji}^v c_{ki}^v c_{j^*i^*}^w c_{k^*i^*}^w + 2 \sum_{i \neq j \neq k} c_{ij}^v c_{ki}^v c_{i^*j^*}^w c_{k^*i^*}^w + \sum_{i \neq j \neq k \neq l} c_{ij}^v c_{kl}^v c_{i^*j^*}^w c_{k^*l^*}^w \end{aligned}$$

où les indices de sommation sont mutuellement distincts. Nous calculons alors de la même façon qu'au paragraphe précédent l'espérance de S^2 :

$$\begin{aligned} \mathcal{E}(S^2) = & \sum_{i \neq j} c_{ij}^v{}^2 \sum_{i' \neq j'} c_{i'j'}^w{}^2 P_{i'j'} + \sum_{i \neq j} c_{ij}^v c_{ji}^v \sum_{i' \neq j'} c_{i'j'}^w c_{j'i'}^w P_{i'j'} \\ & + \sum_{i \neq j \neq k} c_{ij}^v c_{ik}^v \sum_{i' \neq j' \neq k'} c_{i'j'}^w c_{i'k'}^w P_{i'j'k'} + \sum_{i \neq j \neq k} c_{ji}^v c_{ki}^v \sum_{i' \neq j' \neq k'} c_{j'i'}^w c_{k'i'}^w P_{i'j'k'} \\ & + 2 \sum_{i \neq j \neq k} c_{ij}^v c_{ki}^v \sum_{i' \neq j' \neq k'} c_{i'j'}^w c_{k'i'}^w P_{i'j'k'} + \sum_{i \neq j \neq k \neq l} c_{ij}^v c_{kl}^v \sum_{i' \neq j' \neq k' \neq l'} c_{i'j'}^w c_{k'l'}^w P_{i'j'k'l'} \end{aligned}$$

$$\begin{aligned}
(5) \quad \mathcal{E}(S^2) = & \frac{1}{n[2]} \sum_{i \neq j} c_{ij}^v{}^2 \sum_{i' \neq j'} c_{i'j'}^w{}^2 + \frac{1}{n[2]} \sum_{i \neq j} c_{ij}^v c_{ji}^v \sum_{i' \neq j'} c_{i'j'}^w c_{j'i'}^w \\
& + \frac{1}{n[3]} \sum_{i \neq j \neq k} c_{ij}^v c_{ik}^v \sum_{i' \neq j' \neq k'} c_{i'j'}^w c_{i'k'}^w + \frac{1}{n[3]} \sum_{i \neq j \neq k} c_{ji}^v c_{ki}^v \sum_{i' \neq j' \neq k'} c_{j'i'}^w c_{k'i'}^w \\
& + \frac{2}{n[3]} \sum_{i \neq j \neq k} c_{ij}^v c_{ki}^v \sum_{i' \neq j' \neq k'} c_{i'j'}^w c_{k'i'}^w + \frac{1}{n[4]} \sum_{i \neq j \neq k \neq l} c_{ij}^v c_{kl}^v \sum_{i' \neq j' \neq k' \neq l'} c_{i'j'}^w c_{k'l'}^w
\end{aligned}$$

Introduisons les notations suivantes :

$$\begin{aligned}
\alpha &= \sum_{i \neq j} c_{ij} & \beta &= \sum_{i \neq j} c_{ij}^2 & \beta' &= \sum_{i \neq j} c_{ij} c_{ji} & \Gamma &= \sum_{i \neq j \neq k} c_{ij} c_{ik} \\
\Gamma' &= \sum_{i \neq j \neq k} c_{ij} c_{ki} & \Gamma'' &= \sum_{i \neq j \neq k} c_{ji} c_{ki} & \delta &= \sum_{i \neq j \neq k \neq l} c_{ij} c_{kl}
\end{aligned}$$

Nous appliquons ces notations à chacune des variables v et w :

$$\alpha_v = \sum_{i \neq j} c_{ij}^v \quad \alpha_w = \sum_{i \neq j} c_{ij}^w \quad \dots$$

Nous obtenons alors, à partir de (4) et (5), l'expression de l'espérance et de la variance dans le cas général (voir aussi [12] et [13]) :

$$\begin{aligned}
(I) \quad & \mathcal{E}(S) = \frac{1}{n[2]} \alpha_v \alpha_w \\
& \sigma^2(S) = \frac{1}{n[2]} (\beta_v \beta_w + \beta'_v \beta'_w) + \frac{1}{n[3]} (\Gamma_v \Gamma_w + \Gamma''_v \Gamma''_w + 2\Gamma'_v \Gamma'_w) + \frac{1}{n[4]} \delta_v \delta_w - \left(\frac{1}{n[2]} \alpha_v \alpha_w \right)^2
\end{aligned}$$

2.2 Symétrie et antisymétrie

Les relations valuées à comparer sont généralement, toutes les deux soit :

- symétriques : $c_{ij} = c_{ji}$ pour tout couple (i, j) de $I \times I$
- antisymétriques : $c_{ij} = -c_{ji}$ pour tout couple (i, j) de $I \times I$

Dans ce cas on a :

$$\Gamma'' = \Gamma \quad \beta' = \pm \beta \quad \Gamma' = \pm \Gamma$$

D'où :

$$\beta'_v \beta'_w = \beta_v \beta_w \quad \Gamma'_v \Gamma'_w = \Gamma_v \Gamma_w$$

Nous pouvons remarquer par ailleurs qu'on a dans tous les cas :

$$\delta = \alpha^2 - (\beta + \beta' + \Gamma + \Gamma'' + 2\Gamma')$$

D'où :

$$\begin{cases} \delta = \alpha^2 - 2\beta - 4\Gamma & \text{si les deux codages sont symétriques} \\ \delta = \alpha = 0 & \text{si les deux codages sont antisymétriques} \end{cases}$$

Dans ces conditions l'expression de la variance devient :

$$(6) \quad \sigma^2(S) = \frac{2}{n^{[2]}} \beta_v \beta_w + \frac{4}{n^{[3]}} \Gamma_v \Gamma_w - \left(\frac{1}{n^{[2]}} \alpha_v \alpha_w \right)^2 + \frac{1}{n^{[4]}} (\alpha_v^2 - 2\beta_v - 4\Gamma_v) (\alpha_w^2 - 2\beta_w - 4\Gamma_w)$$

D'autre part, les sommes triples Γ peuvent être remplacées par des sommes doubles notées γ :

$$\Gamma = \sum_{i \neq j \neq k} c_{ij} c_{ik} = \sum_i \left(\sum_{j, j \neq i} c_{ij} \right)^2 - \sum_{i \neq j} c_{ij}^2 = \gamma - \beta$$

L'expression (6) devient alors :

$$(7) \quad \sigma^2(S) = \frac{2}{n^{[2]}} \beta_v \beta_w + \frac{4}{n^{[3]}} (\gamma_v - \beta_v) (\gamma_w - \beta_w) - \left(\frac{1}{n^{[2]}} \alpha_v \alpha_w \right)^2 + \frac{1}{n^{[4]}} (\alpha_v^2 + 2\beta_v - 4\gamma_v) (\alpha_w^2 + 2\beta_w - 4\gamma_w)$$

Remarque :

Les deux derniers termes de (6) et (7) sont nuls dans le cas antisymétrique.

Si comme l'indique I.C. LERMAN dans [12], l'expression (6) est plus claire d'un point de vue formellement conceptuel, elle présente par contre l'handicap de contenir des sommes triples (Γ). A l'inverse l'expression (7) de N. MANTEL [13] se prête mieux aux calculs puisqu'elle n'utilise que des sommes doubles.

De notre côté, nous démontrons dans la proposition ci-dessous, une nouvelle expression de la variance dans le cas (le plus complexe) où le codage est symétrique :

Proposition 1 :

si les codages $\{c_{ij}\}$ sont symétriques alors :

$$(8) \quad \sigma^2(S) = \frac{2}{n^{[3]}} \left[\frac{\alpha_v^2 \alpha_w^2}{n^{[2]}} + (n-1)\beta_v \beta_w - 2\gamma_v \gamma_w \right. \\ \left. + \frac{n-1}{n-3} \left(\frac{\alpha_v^2}{n-1} + \beta_v - 2\gamma_v \right) \left(\frac{\alpha_w^2}{n-1} + \beta_w - 2\gamma_w \right) \right]$$

Démonstration :

En développant $\sigma^2(S)$ dans le cas symétrique, nous obtenons :

$$\sigma^2(S) = \frac{2(2n-3)}{n^{[2]n^{[4]}}} \alpha_v^2 \alpha_w^2 + \frac{2}{n(n-3)} \beta_v \beta_w + \frac{4(n+1)}{n^{[4]}} \gamma_v \gamma_w \\ + \left[\frac{2}{n^{[4]}} (\alpha_v^2 \beta_w + \alpha_w^2 \beta_v) - \frac{4}{n^{[4]}} (\alpha_v^2 \gamma_w + \alpha_w^2 \gamma_v) - \frac{4(n-1)}{n^{[4]}} (\beta_v \gamma_w + \beta_w \gamma_v) \right]$$

Il suffit alors de remarquer que la partie entre crochets est égale à :

$$\frac{2(n-1)}{n^{[4]}} \left[\left(\frac{\alpha_v^2}{n-1} + \beta_v - 2\gamma_v \right) \left(\frac{\alpha_w^2}{n-1} + \beta_w - 2\gamma_w \right) - \frac{\alpha_v^2 \alpha_w^2}{(n-1)^2} - \beta_v \beta_w - 4\gamma_v \gamma_w \right]$$

D'où le résultat final.

En comparant notre expression (8) aux expressions (6) et (7), nous pouvons remarquer qu'elle présente un double intérêt :

- Comme dans (6), chacun des trois premiers termes ne dépend que d'une seule somme α , β ou γ .
- Comme dans (7), elle ne contient que des sommes doubles.

3 Expressions du coefficient centré-réduit

L'expression du coefficient centré-réduit (3), dans le cas général, est fournie par les formules (I). Ce cas présente peu d'intérêt pour nous car les codages qu'on utilise pour tous les types de variables sont soit symétriques, soit antisymétriques.

Nous chercherons donc ici à exprimer ce coefficient dans le cas d'un codage symétrique (resp. antisymétrique), le plus clairement possible, en utilisant des notations classiques. Pour cela nous introduisons les "moments centrés" ϕ et ψ définis par :

$$\phi_{vw} = \frac{1}{n^{[2]}} \sum_{i \neq j} c_{ij}^v c_{ij}^w - \left(\frac{1}{n^{[2]}} \sum_{i \neq j} c_{ij}^v \right) \left(\frac{1}{n^{[2]}} \sum_{i \neq j} c_{ij}^w \right)$$

$$\psi_v = \frac{1}{n(n-1)^2} \sum_i \left(\sum_{j, j \neq i} c_{ij}^v \right)^2 - \left(\frac{1}{n^{[2]}} \sum_{i \neq j} c_{ij}^v \right)^2 \quad (\text{idem pour } \psi_w)$$

Remarques :

- ϕ_{vw} est une covariance entre les deux tables $\{c_{ij}^v \mid 1 \leq i \neq j \leq n\}$ et $\{c_{ij}^w \mid 1 \leq i \neq j \leq n\}$ induites par les variables v et w
- ϕ_{vv} qui sera notée ϕ_v est donc la variance de la table $\{c_{ij}^v \mid 1 \leq i \neq j \leq n\}$
- ψ_v est la variance des marges de la table $\{c_{ij}^v \mid 1 \leq i \neq j \leq n\}$
- la deuxième partie de ces moments est nulle dans le cas antisymétrique.

Avec ces notations, le numérateur de (3) s'écrit :

$$(9) \quad s(v, w) - \mathcal{E}(S) = n^{[2]} \phi_{vw}$$

Reste à exprimer la variance de S , pour cela nous traiterons séparément les deux cas symétrique et antisymétrique. Dans chacun de ces deux cas nous démontrons deux nouvelles expressions de $\sigma^2(S)$. Chacune se présente sous la forme d'une somme de deux termes : l'un ne dépend que d'un seul moment ϕ ou ψ , l'autre étant une combinaison linéaire des deux. Nous en déduisons deux formes distinctes du coefficient centré-réduit, intitulées :

- Expression ϕ : un des deux termes de la variance ne dépend que de ϕ
- Expression ψ : un des deux termes de la variance ne dépend que de ψ

3.1 Expression ϕ

3.1.1 Cas symétrique

Proposition 2 :

si les codages $\{c_{ij}\}$ sont symétriques alors :

$$\sigma^2(S) = \frac{(2n^{[2]})^2}{n-2} \left[\frac{n^2-1}{n(n-3)} \left(\psi_v - \frac{\phi_v}{n+1} \right) \left(\psi_w - \frac{\phi_w}{n+1} \right) + \frac{1}{2(n+1)} \phi_v \phi_w \right]$$

Démonstration :

Nous développons comme dans la proposition 1, l'expression de $\sigma^2(S)$:

$$\begin{aligned} \sigma^2(S) &= \frac{2(2n-3)}{n^{[2]}n^{[4]}} \alpha_v^2 \alpha_w^2 + \frac{2}{n(n-3)} \beta_v \beta_w + \frac{2}{n^{[4]}} (\alpha_v^2 \beta_w + \alpha_w^2 \beta_v) \\ &+ \left[\frac{4(n+1)}{n^{[4]}} \gamma_v \gamma_w - \frac{4}{n^{[4]}} (\alpha_v^2 \gamma_w + \alpha_w^2 \gamma_v) - \frac{4(n-1)}{n^{[4]}} (\beta_v \gamma_w + \beta_w \gamma_v) \right] \end{aligned}$$

La partie entre crochets est égale à :

$$\begin{aligned} \frac{4}{n^{[4]}(n+1)} &\{ [(n+1)\gamma_v - (n-1)\beta_v + \alpha_v^2] [(n+1)\gamma_w - (n-1)\beta_w + \alpha_w^2] \\ &- \alpha_v^2 \alpha_w^2 - (n-1)^2 \beta_v \beta_w - (n-1)(\alpha_v^2 \beta_w + \alpha_w^2 \beta_v) \} \end{aligned}$$

D'où :

$$\begin{aligned} \sigma^2(S) &= \frac{4}{n^{[4]}(n+1)} [(n+1)\gamma_v - (n-1)\beta_v + \alpha_v^2] [(n+1)\gamma_w - (n-1)\beta_w + \alpha_w^2] \\ &+ \frac{2\alpha_v^2 \alpha_w^2}{n^{[2]}n^{[3]}(n+1)} + \frac{2\beta_v \beta_w}{(n-2)(n+1)} - \frac{2(\alpha_v^2 \beta_w + \alpha_w^2 \beta_v)}{n^{[3]}(n+1)} \\ &= \frac{4}{n^{[4]}(n+1)} [(n+1)\gamma_v - (n-1)\beta_v + \alpha_v^2] [(n+1)\gamma_w - (n-1)\beta_w + \alpha_w^2] \\ &+ \frac{2}{(n+1)(n-2)} \left(\beta_v - \frac{\alpha_v^2}{n^{[2]}} \right) \left(\beta_w - \frac{\alpha_w^2}{n^{[2]}} \right) \end{aligned}$$

Il suffit ensuite de remarquer que :

$$\begin{aligned} (n+1)\gamma_v - (n-1)\beta_v + \alpha_v^2 &= n(n-1)^2 \left[(n+1) \frac{\gamma_v}{n(n-1)^2} - \frac{\beta_v}{n^{[2]}} + n \left(\frac{\alpha_v}{n^{[2]}} \right)^2 \right] \\ &= n(n-1)^2 \left[(n+1) \left(\frac{\gamma_v}{n(n-1)^2} - \left(\frac{\alpha_v}{n^{[2]}} \right)^2 \right) - \left(\frac{\beta_v}{n^{[2]}} - \left(\frac{\alpha_v}{n^{[2]}} \right)^2 \right) \right] \\ &= n(n-1)^2 [(n+1)\psi_v - \phi_v] \end{aligned}$$

et que :

$$\beta_v - \frac{\alpha_v^2}{n^{[2]}} = n^{[2]} \left[\frac{\beta_v}{n^{[2]}} - \left(\frac{\alpha_v}{n^{[2]}} \right)^2 \right] = n^{[2]} \phi_v$$

D'où le résultat final.

Nous déduisons de ce résultat en tenant compte de (9), l'expression ϕ du coefficient centré-réduit dans le cas symétrique :

$$(II) \quad Q(v, w) = \frac{\sqrt{n-2}}{2} \frac{\phi_{vw}}{\sqrt{\frac{n^2-1}{n(n-3)} \left(\psi_v - \frac{\phi_v}{n+1} \right) \left(\psi_w - \frac{\phi_w}{n+1} \right) + \frac{\phi_v \phi_w}{2(n+1)}}$$

3.1.2 Cas antisymétrique

Proposition 3 :

si les codages $\{c_{ij}\}$ sont antisymétriques alors :

$$\sigma^2(S) = 4n(n-1)^2 \left[\frac{n-1}{n-2} \left(\psi_v - \frac{\phi_v}{n-1} \right) \left(\psi_w - \frac{\phi_w}{n-1} \right) + \frac{1}{2(n-1)} \phi_v \phi_w \right]$$

Démonstration :

Comme nous l'avons vu précédemment, si les codages sont antisymétriques alors $\alpha = \delta = 0$, d'où :

$$\begin{aligned} \sigma^2(S) &= \frac{4}{n^{[3]}} (\gamma_v - \beta_v)(\gamma_w - \beta_w) + \frac{2}{n^{[2]}} \beta_v \beta_w \\ &= \frac{4n(n-1)^3}{n-2} \left(\frac{\gamma_v}{n(n-1)^2} - \frac{1}{n-1} \frac{\beta_v}{n^{[2]}} \right) \left(\frac{\gamma_w}{n(n-1)^2} - \frac{1}{n-1} \frac{\beta_w}{n^{[2]}} \right) \\ &\quad + 2n^{[2]} \frac{\beta_v}{n^{[2]}} \frac{\beta_w}{n^{[2]}} \end{aligned}$$

Or dans le cas antisymétrique on a :

$$\frac{1}{n(n-1)^2} \gamma = \psi \quad \text{et} \quad \frac{1}{n^{[2]}} \beta = \phi$$

D'où le résultat final.

Nous déduisons de ce résultat en tenant compte de (9), l'expression ϕ du coefficient centré-réduit dans le cas antisymétrique :

$$(II') \quad Q(v, w) = \frac{\sqrt{n}}{2} \frac{\phi_{vw}}{\sqrt{\frac{n-1}{n-2} \left(\psi_v - \frac{\phi_v}{n-1} \right) \left(\psi_w - \frac{\phi_w}{n-1} \right) + \frac{\phi_v \phi_w}{2(n-1)}}$$

Les deux expressions (II) et (II') sont équivalentes, et le terme dominant de la variance dépend des deux moments ϕ et ψ .

3.2 Expression ψ

3.2.1 Cas symétrique

Proposition 4 :

si les codages $\{c_{ij}\}$ sont symétriques alors :

$$\sigma^2(S) = \frac{4n^2(n-1)^3}{(n-2)^2} \left[\psi_v \psi_w + \frac{n-1}{2n(n-3)} \left(\frac{n-2}{n-1} \phi_v - 2\psi_v \right) \left(\frac{n-2}{n-1} \phi_w - 2\psi_w \right) \right]$$

Démonstration :

Nous développons comme dans la proposition 2 l'expression de $\sigma^2(S)$:

$$\begin{aligned} \sigma^2(S) = & \frac{2(2n-3)}{n[2]n[4]} \alpha_v^2 \alpha_w^2 + \frac{4(n+1)}{n[4]} \gamma_v \gamma_w - \frac{4}{n[4]} (\alpha_v^2 \gamma_w + \alpha_w^2 \gamma_v) \\ & + \left[\frac{2}{n(n-3)} \beta_v \beta_w + \frac{2}{n[4]} (\alpha_v^2 \beta_w + \alpha_w^2 \beta_v) - \frac{4(n-1)}{n[4]} (\beta_v \gamma_w + \beta_w \gamma_v) \right] \end{aligned}$$

La partie entre crochets est égale à :

$$\begin{aligned} \frac{2(n-1)}{n[4](n-2)} & \left[\left(\frac{\alpha_v^2}{n-1} + (n-2)\beta_v - 2\gamma_v \right) \left(\frac{\alpha_w^2}{n-1} + (n-2)\beta_w - 2\gamma_w \right) \right. \\ & \left. - 4\gamma_v \gamma_w - \left(\frac{\alpha_v \alpha_w}{n-1} \right)^2 + \frac{2}{n-1} (\alpha_v^2 \gamma_w + \alpha_w^2 \gamma_v) \right] \end{aligned}$$

D'où :

$$\begin{aligned}
\sigma^2(S) &= \frac{2(n-1)}{n^{[4]}(n-2)} \left(\frac{\alpha_v^2}{n-1} + (n-2)\beta_v - 2\gamma_v \right) \left(\frac{\alpha_w^2}{n-1} + (n-2)\beta_w - 2\gamma_w \right) \\
&\quad + \frac{4n\beta_v\beta_w}{n^{[2]}(n-2)} + \frac{4(n-1)\alpha_v^2\alpha_w^2}{(n^{[3]})^2} - \frac{4(\alpha_v^2\gamma_v + \alpha_w^2\gamma_w)}{n^{[3]}(n-2)} \\
&= \frac{2(n-1)}{n^{[4]}(n-2)} \left(\frac{\alpha_v^2}{n-1} + (n-2)\beta_v - 2\gamma_v \right) \left(\frac{\alpha_w^2}{n-1} + (n-2)\beta_w - 2\gamma_w \right) \\
&\quad + \frac{4}{(n-1)(n-2)^2} \left(\gamma_v - \frac{\alpha_v^2}{n} \right) \left(\gamma_w - \frac{\alpha_w^2}{n} \right)
\end{aligned}$$

Il suffit ensuite de remarquer que :

$$\begin{aligned}
\frac{\alpha_v^2}{n-1} + (n-2)\beta_v - 2\gamma_v &= n(n-1)^2 \left(\frac{n}{n-1} \left(\frac{\alpha_v}{n^{[2]}} \right)^2 + \frac{n-2}{n-1} \frac{\beta_v}{n^{[2]}} - 2 \frac{\gamma_v}{n(n-1)^2} \right) \\
&= n(n-1)^2 \left[\frac{n-2}{n-1} \left(\frac{\beta_v}{n^{[2]}} - \left(\frac{\alpha_v}{n^{[2]}} \right)^2 \right) - 2 \left(\frac{\gamma_v}{n(n-1)^2} - \left(\frac{\alpha_v}{n^{[2]}} \right)^2 \right) \right] \\
&= n(n-1)^2 \left(\frac{n-2}{n-1} \phi_v - 2\psi_v \right)
\end{aligned}$$

et que :

$$\begin{aligned}
\gamma_v - \frac{\alpha_v^2}{n} &= n(n-1)^2 \left(\frac{\gamma_v}{n(n-1)^2} - \left(\frac{\alpha_v}{n^{[2]}} \right)^2 \right) \\
&= n(n-1)^2 \psi_v
\end{aligned}$$

D'où le résultat final.

Nous déduisons de ce résultat en tenant compte de (9), l'expression ψ du coefficient centré-réduit dans le cas symétrique :

$$\text{(III)} \quad Q(v, w) = \frac{n-2}{2\sqrt{n-1}} \frac{\phi_{vw}}{\sqrt{\psi_v\psi_w + \frac{n-1}{2n(n-3)} \left(\frac{n-2}{n-1} \phi_v - 2\psi_v \right) \left(\frac{n-2}{n-1} \phi_w - 2\psi_w \right)}}$$

3.2.2 Cas antisymétrique

Proposition 5 :

si les codages $\{c_{ij}\}$ sont antisymétriques alors :

$$\sigma^2(S) = 4(n-1)^3 \left[\psi_v \psi_w + \frac{1}{2(n-2)} \left(\frac{n}{n-1} \phi_v - 2\psi_v \right) \left(\frac{n}{n-1} \phi_w - 2\psi_w \right) \right]$$

Démonstration :

Comme dans la proposition 3, nous développons l'expression de $\sigma^2(S)$:

$$\begin{aligned} \sigma^2(S) &= \frac{4}{n^{[3]}} \gamma_v \gamma_w + \frac{2n}{n^{[3]}} \beta_v \beta_w - \frac{4}{n^{[3]}} (\beta_v \gamma_w + \beta_w \gamma_v) \\ &= \left(\frac{4}{n^{[3]}} - \frac{8}{nn^{[3]}} \right) \gamma_v \gamma_w + \frac{2}{nn^{[3]}} \left[n^2 \beta_v \beta_w - 2n(\beta_v \gamma_w + \beta_w \gamma_v) + 4\gamma_v \gamma_w \right] \\ &= \frac{4}{nn^{[2]}} \gamma_v \gamma_w + \frac{2}{nn^{[3]}} (n\beta_v - 2\gamma_v)(n\beta_w - 2\gamma_w) \\ &= 4(n-1)^3 \frac{\gamma_v \gamma_w}{n^2(n-1)^4} \\ &\quad + \frac{2(n-1)^3}{n-2} \left(\frac{n}{n-1} \frac{\beta_v}{n^{[2]}} - 2 \frac{\gamma_v}{n(n-1)^2} \right) \left(\frac{n}{n-1} \frac{\beta_w}{n^{[2]}} - 2 \frac{\gamma_w}{n(n-1)^2} \right) \end{aligned}$$

Or dans le cas antisymétrique on a :

$$\frac{1}{n(n-1)^2} \gamma = \psi \quad \text{et} \quad \frac{1}{n^{[2]}} \beta = \phi$$

D'où le résultat final.

Nous déduisons de ce résultat en tenant compte de (9), l'expression ψ du coefficient centré-réduit dans le cas antisymétrique :

$$(III') \quad Q(v, w) = \frac{n}{2\sqrt{n-1}} \frac{\phi_{vw}}{\sqrt{\psi_v \psi_w + \frac{1}{2(n-2)} \left(\frac{n}{n-1} \phi_v - 2\psi_v \right) \left(\frac{n}{n-1} \phi_w - 2\psi_w \right)}}$$

Les deux expressions (III) et (III') sont équivalentes, et le terme dominant de la variance ne dépend que de ψ .

4 Comportement asymptotique

Pour n suffisamment grand ($n \simeq n-1 \simeq n-2 \simeq n-3$), les expressions ϕ du coefficient centré-réduit (II) et (II') sont identiques. Il en est de même pour les expressions ψ (III) et (III'). Nous obtenons donc la même expression pour $Q(v, w)$ dans le cas symétrique que dans le cas antisymétrique :

Expression ϕ :

$$(IV) \quad Q(v, w) = \frac{\sqrt{n}}{2} \frac{\phi_{vw}}{\sqrt{(\psi_v - \frac{\phi_v}{n})(\psi_w - \frac{\phi_w}{n}) + \frac{\phi_v \phi_w}{2n}}}$$

Expression ψ :

$$(V) \quad Q(v, w) = \frac{\sqrt{n}}{2} \frac{\phi_{vw}}{\sqrt{\psi_v \psi_w + \frac{1}{2n}(\phi_v - 2\psi_v)(\phi_w - 2\psi_w)}}$$

Suivant le codage utilisé et le type des variables étudiées, l'une ou l'autre des deux formes (IV) ou (V) peut paraître plus adéquate.

L'expression ψ semble plus adaptée pour étudier la forme limite du coefficient centré-réduit. En effet, le terme dominant de la variance ne dépend que du seul moment ψ , et le deuxième terme est une combinaison linéaire de ψ et ϕ indépendante de n .

Pour extraire une forme limite de l'indice, il faudrait s'assurer que le terme dominant est bien positif, et préciser éventuellement dans quelles conditions il pourrait s'annuler. Pour cela nous allons étudier les signes de ψ et ϕ :

4.1 Signes de ϕ et ψ

Comme nous l'avons déjà remarqué ϕ et ψ sont des variances, donc positives. Nous pouvons le vérifier à l'aide des deux propositions suivantes :

Proposition 6 :

$$\begin{cases} \psi \geq 0 \\ \psi = 0 \Leftrightarrow c_i \text{ est constant } \forall i \end{cases}$$

où :

$$c_i = \sum_{j, j \neq i} c_{ij}$$

Démonstration :

$$\begin{aligned} \psi &= \frac{1}{n(n-1)} \sum_i \left(\sum_{j, j \neq i} c_{ij} \right)^2 - \left(\frac{1}{n^{[2]}} \sum_i \sum_{j \neq i} c_{ij} \right)^2 \\ &= \frac{1}{(n^{[2]})^2} \left[n \sum_i c_i^2 - \left(\sum_i c_i \right)^2 \right] \\ &= \frac{1}{(n^{[2]})^2} \left[n \sum_i c_i^2 - \sum_i c_i^2 - \sum_{i \neq j} c_i c_j \right] \\ &= \frac{1}{(n^{[2]})^2} \left[(n-1) \sum_i c_i^2 + \frac{1}{2} \sum_{i \neq j} \left((c_i - c_j)^2 - c_i^2 - c_j^2 \right) \right] \\ &= \frac{1}{(n^{[2]})^2} \left[(n-1) \sum_i c_i^2 + \frac{1}{2} \sum_{i, j} (c_i - c_j)^2 - (n-1) \sum_i c_i^2 \right] \\ &= \frac{1}{2(n^{[2]})^2} \sum_{i, j} (c_i - c_j)^2 \end{aligned}$$

D'où le résultat final.

Remarque :

Si le codage est antisymétrique alors : $\psi = 0 \Leftrightarrow c_i = 0 \forall i$

Proposition 7 :

$$\begin{cases} \phi \geq 0 \\ \phi = 0 \Leftrightarrow c_{ij} \text{ est constant } \forall i \neq j \end{cases}$$

Démonstration :

$$\begin{aligned} \phi &= \frac{1}{n^{[2]}} \sum_{i \neq j} c_{ij}^2 - \left(\frac{1}{n^{[2]}} \sum_{i \neq j} c_{ij} \right)^2 \\ &= \frac{1}{(n^{[2]})^2} \left[n^{[2]} \sum_{i \neq j} c_{ij}^2 - \left(\sum_{i \neq j} c_{ij} \right)^2 \right] \\ &= \frac{1}{(n^{[2]})^2} \left[n^{[2]} \sum_{i \neq j} c_{ij}^2 - \sum_{i \neq j} c_{ij}^2 - \sum_{\substack{(i,j) \neq (i',j') \\ i \neq j, i' \neq j'}} c_{ij} c_{i'j'} \right] \\ &= \frac{1}{(n^{[2]})^2} \left[(n^{[2]} - 1) \sum_{i \neq j} c_{ij}^2 + \frac{1}{2} \sum_{\substack{(i,j) \neq (i',j') \\ i \neq j, i' \neq j'}} [(c_{ij} - c_{i'j'})^2 - c_{ij}^2 - c_{i'j'}^2] \right] \\ &= \frac{1}{(n^{[2]})^2} \left[(n^{[2]} - 1) \sum_{i \neq j} c_{ij}^2 + \frac{1}{2} \sum_{\substack{(i,j) \neq (i',j') \\ i \neq j, i' \neq j'}} (c_{ij} - c_{i'j'})^2 - (n^{[2]} - 1) \sum_{i \neq j} c_{ij}^2 \right] \\ &= \frac{1}{2(n^{[2]})^2} \sum_{\substack{i \neq j, i' \neq j'}} (c_{ij} - c_{i'j'})^2 \end{aligned}$$

D'où le résultat final.

Remarque :

Si le codage est antisymétrique alors : $\phi = 0 \Leftrightarrow c_{ij} = 0 \forall i \neq j$

4.2 Forme limite

Le codage $\{c_{ij}^v / (i, j) \in I \times I\}$ ne dépend pas de la taille n de l'ensemble des objets, mais uniquement du nombre de modalités et du type de la variable v . Les quantités ϕ et ψ sont donc asymptotiquement indépendantes de n .

Nous pouvons donc supposer que ϕ et ψ ont des limites finies. D'où quand n tend vers l'infini, le deuxième terme de la variance dans (V) : $\frac{1}{2n} (\phi_v - 2\psi_v) (\phi_w - 2\psi_w)$ devient négligeable devant le premier terme : $\psi_v \psi_w$.

Dans ces conditions la forme limite du coefficient centré-réduit s'écrit :

$$(VI) \quad Q(v, w) \simeq \frac{\sqrt{n}}{2} \frac{\phi_{vw}}{\sqrt{\psi_v \psi_w}}$$

Nous avons montré, dans la proposition 6, que la quantité ψ peut s'annuler si c_i est constant. L'interprétation de cette condition dépend de la nature de la variable et de son codage. Nous montrerons dans un autre rapport (à paraître) que ψ ne s'annule jamais pour les variables ordinales. En revanche, si les modalités sont non ordonnées, ψ peut s'annuler dans un seul cas, celui par exemple des variables nominales à modalités uniformément distribuées.

Dans ces conditions, si $\psi_v = \psi_w = 0$, la forme limite du coefficient centré-réduit s'écrit :

$$(10) \quad Q(v, w) \simeq \frac{n}{\sqrt{2}} \frac{\phi_{vw}}{\sqrt{\phi_v \phi_w}}$$

Remarque :

Le dénominateur de (10) est non nul. En effet, comme nous l'avons montré dans la proposition 7, ϕ_v ne s'annule que dans le cas trivial où la variable v induit le même codage sur tous les couples d'objets.

4.3 Réduction géométrique

Nous avons donc montré que le coefficient $Q(v, w)$ entre deux relations valuées est, à un facteur en n près, un rapport "pur" asymptotiquement indépendant de n .

Si le coefficient (VI) correspond à une notion de ressemblance statistique, la "réduction géométrique" suivante : $Q(v, w) = \frac{Q(v, w)}{\sqrt{Q(v, v)Q(w, w)}}$ permettra d'éliminer le facteur en n , et de mettre davantage l'accent sur la "similarité des formes". Nous obtenons ainsi un coefficient de nature corrélative.

Nous obtenons donc à partir de l'expression (VI) :

$$(VII) \quad Q(v, w) \simeq \frac{\phi_{vw}}{\sqrt{\phi_v \phi_w}}$$

Remarque :

Si $\psi_v = \psi_w = 0$ la "réduction géométrique" de $Q(v, w)$ dans (10) donne le même résultat. L'expression (VII) reste donc valable même si $\psi_v = \psi_w = 0$.

5 Conclusion

Dans ce travail, nous avons synthétisé et condensé les expressions d'un coefficient d'association qui se présentait, à priori, de façon complexe. Il s'en est notamment dégagé, dans son expression limite, une forme corrélative.

Un autre volet de notre travail, consiste à élaborer un codage approprié à chaque type de variable qualitative, et à expliciter l'expression de l'indice qui en découle, en utilisant le support des tables de contingences.

Ceci fera l'objet d'un futur article, déjà mentionné, où nous considérons toutes les variables qualitatives comme des variables préordonnances, de la manière suivante : nous attribuons à chaque couple (k, l) de modalités d'une même variable v , un rang noté $rg_v(k, l)$ qui dépendra de la nature de cette variable. Nous en déduisons un codage sur l'ensemble des couples d'individus :

$c_{ij}^v = rg_v(k, l)$ si l'individu i (resp. j) possède la modalité k (resp. l) de la variable v .

Cette structure descriptive des variables qualitatives, tout en demeurant très générale, nous semble la plus riche et la moins arbitraire, d'autant plus qu'elle peut être directement fournie par l'expert. D'autre part, ce codage permet le calcul d'un coefficient d'association entre deux variables qualitatives de natures différentes.

Bibliographie

- [1] P. ARABIE & L.J. HUBERT *Comparing partitions* Journal of Classification vol. 2, 1985.
- [2] H.E. DANIELS *The relation between measures of correlation in the universe of sample permutations* Biometrika vol. 33, 1944.
- [3] L.J. HUBERT *A Relationship between the assignment problem and some statistical techniques* Quality and Quantity vol. 10, 1976.
- [4] L.J. HUBERT *Combinatorial Data Analysis : Association and Partial Association for Nominal Data* Psychometrika vol. 50, 1985.
- [5] L.J. HUBERT *Assignment Methods in Combinatorial Data Analysis* Decker, New York, 1987.
- [6] M.G. KENDALL *Rank Correlations Methods* Griffin, Londres, 1962.
- [7] G. LE CALVE *Un indice de similarité pour les variables de type quelconques* Statistiques et analyse des données, 1976.
- [8] I.C. LERMAN *Les bases de la Classification Automatique* Villars, Paris, 1970.
- [9] I.C. LERMAN *Etude distributionnelle de statistiques de proximité entre structures finies de même type* Cahiers du B.U.R.O. N° 19, 1973.
- [10] I.C. LERMAN *Formal analysis of a general notion of proximity between variables* North Holland, 1977.
- [11] I.C. LERMAN *Classification et analyse ordinale des données* Dunod, Paris, 1981.
- [12] I.C. LERMAN *Analyse de la forme limite de coefficients statistiques d'association entre variables relationnelles* P.I. N° 367 I.R.I.S.A., 1987.
- [13] N. MANTEL *Detection of disease clustering and a generalized regression approach* Cancer Research vol. 27, 1967.

-
- [14] F. MARCOTORCHINO & MICHAUD *Optimisation en Analyse Ordinale des Données* Masson, Paris, 1979.
- [15] B. MONJARDET & J.P. BARTHELEMY *Ajustement et résumé de Données relationnelles : les relations centrales* North Holland, 1980.
- [16] P.W. MIELKE *On asymptotic non normality of null distributions of MRPP Statistics* Theory and Methods, 1979.

LISTE DES DERNIERES PUBLICATIONS INTERNES IRISA

- PI 544 **IMPLEMENTATION AND EVALUATION OF DISTRIBUTED SYNCHRONIZATION ON A DISTRIBUTED MEMORY PARALLEL MACHINE**
André COUVERT, René PEDRONO, Michel RAYNAL
Juillet 1990, 14 Pages.
- PI 545 **ESTIMATION OF NETWORK RELIABILITY ON A PARALLEL MACHINE BY MEANS OF A MONTE CARLO TECHNIQUE**
Mohamed EL KHADIRI, Raymond MARIE, Gerardo RUBINO
Août 1990, 20 Pages.
- PI 546 **LIMIT THEOREMS FOR MIXING PROCESSES**
Bernard DELYON
Septembre 1990, 22 Pages.
- PI 547 **PERFORMANCES DES COMMUNICATIONS SUR LE T-NODE**
Frédéric GUIDEC
Septembre 1990, 38 Pages.
- PI 548 **LES PREDICATS COLLECTIFS : UN MOYEN D'EXPRESSION DU CONTROLE DU PARALLELISME OU EN PROLOG**
René QUINIOU, Laurent TRILLING
Septembre 1990, 34 Pages.
- PI 549 **NORMALISATION SOUS HYPOTHESE D'ABSENCE DE LIEN APPLICATION AU CAS NOMINAL**
François DAUDE
Septembre 1990, 42 Pages.
- PI 550 **MULTISCALE SIGNAL PROCESSING : FROM QMF TO WAVELETS**
Albert BENVENISTE
Septembre 1990, 28 Pages.
- PI 551 **ON THE TRANSITION GRAPHS OF AUTOMATA AND GRAMMARS**
Didier CAUCAL, Roland MONFORT
Septembre 1990, 46 Pages.
- PI 552 **ERREURS DE CALCUL DES ORDINATEURS**
Jocelyne ERHEL
Septembre 1990, 58 Pages.
- PI 553 **SEQUENTIAL FUNCTIONS**
Boubakar GAMATIE
Octobre 1990, 16 Pages.
- PI 554 **ANALYSE DE LA FORME D'UN COEFFICIENT D'ASSOCIATION ENTRE VARIABLES QUALITATIVES**
Mohamed OUALI ALLAH
Octobre 1990, 26 Pages.
- PI 555 **APPROXIMATION BY NONLINEAR WAVELET NETWORKS**
Qinghua ZHANG, Albert BENVENISTE
Octobre 1990, 16 Pages.

ISSN 0249 - 6399