Fundamental Concepts on Classification and Statistical Implicative Analysis for Modal Variables

Conceptos fundamentales de la clasificación y del análisis estadístico implicativo para variables modales

Larisa Zamora-Matamoros^a, Jorge Rey Díaz-Silvera^b, Lariza Portuondo-Mallet^c

Department of Mathematics, Faculty of Mathematics and Computer Science, University of Oriente, Santiago de Cuba, Cuba

Abstract

The present work offers some fundamental concepts of the Statistical Implicative Analysis for modal variables and proposes an index to establish the similarity between two modal variables. Expressions to calculate typicality and contribution of the individuals to the classes that are formed in the classification are also presented. This technique is illustrated by two examples, one with binary data, which allows the coincidence between the formulas presented and the existent for the binary case to be shown, and the other for modal data with more than two modalities.

Key words: Classificatory Analysis, Typicality, Contribution, Modal Variables.

Resumen

En el presente trabajo se ofrecen algunos conceptos fundamentales del análisis estadístico implicativo para el caso de variables modales y se propone un índice para establecer la similaridad entre dos variables modales, así como expresiones para el cálculo de la tipicalidad y contribución de los individuos a las clases que se forman en la clasificación. Con el objetivo de ilustrar la técnica presentada, es aplicada a dos juegos de datos, uno binario, el cual permite mostrar numéricamente la coincidencia de las fórmulas presentadas con las existentes para el caso de variables binarias, y otro modal con más de dos modalidades.

 ${\it Palabras}\ {\it clave:}$ análisis clasificatorio, tipicalidad, contribución, variables modales.

^aProfessor. E-mail: larisa@csd.uo.edu.cu

^bProfessor. E-mail: jdiaz@csd.uo.edu.cu

^cProfessor. E-mail: lportuondo@csd.uo.edu.cu

1. Introduction

Régis Gras and its collaborators have developed the Statistical Implicative Analysis (SIA), which allows association rules in a dataset to be established, crossing individuals and variables.

The initial objective of this method was to answer the question: if an object has a certain quality, could it have some other one? This is the case of binary variables (Gras 2000).

The research developed by Marc Bailleul from 1991 to 1994, (Bailleul & Gras 1994) showed the necessity to extend the concept of statistical implication to nonbinary variables, in particular to the modal variables.

Bailleul & Gras (1994) refer to the representation that the professors of mathematics make of their own teaching, for which, the professors have been asked to order a list of significant words depending on their importance. Their elections are no longer binary. The words selected by a professor are arranged in a scale where the top words are the most representative. The question of Bailleul is centered on questions of the type: "if I select a word x with the importance i_x then I select the word y with the importance $i_y \geq i_x$ ".

For the application of the SIA techniques, CHIC (Classification Hiérarchique Implicative et Cohésitive), (Couturier 2008) has been developed. This is a software of analysis of data that was initially conceived by Regis Gras and it has been successively improved for personal computers by Saddo Ag Almouloud, Harrisson Ratsimba-Rajohn and, in its latest versions, by Raphaël Couturier (Version 5.0), who incorporates the classification, by means of the similarity analysis.

Unlike the techniques of SIA, which are concerned with the cohesive and the implicative analysis, in the specialized literature, expressions for the determination of the similarity index, the significant nodes, the typicality and the contribution of individuals to the classes in the case of modal variables do not appear. For this reason, the present paper proposes a similarity index for the case of modal variables, as well as expressions for the determination of the typicality and contribution of the individuals to the classes formed during the construction of the similarity tree, extending the works of Bailleul & Gras (1994) and Lagrange (1998).

2. Statistical Implicative Analysis for Modal Variables

In this section, a brief summary of the theory developed by Bailleul & Gras (1994) and Lagrange (1998) for the case of modal variables will be presented. It will be the guideline to define a similarity index for these variables.

2.1. The Intensity of Propensity

Lagrange (1998) defines the "propensity" as a non-symmetric association of two variables whose values belong to the interval [0, 1] and introduces the "intensity of propensity" to decide whether or not an observed association is significant. He also proves that in the case of binary variables, the intensity of propensity is similar to the implication intensity.

Similar to the case of binary variables, let us consider a set I formed by n individuals and a set $A = \{a_1, a_2, \ldots, a_p\}$ formed by p modal variables, each one of which can take values in the set $M = \{k \in N : 1 \le k \le m\}$.

To transform the m modalities to values in the [0, 1] interval, a weight is assigned to each modality, varying from 1 (the strongest) to 0 by fractions of size $\frac{1}{m-1}$, by means of the function ψ defined by:

$$\begin{aligned}
\Psi : & \mathbf{M} \longrightarrow \mathbf{VAL} \\
& \mathbf{k} \longrightarrow \frac{k-1}{m-1}
\end{aligned} \tag{1}$$

where $VAL = \left\{0, \frac{1}{m-1}, \frac{2}{m-1}, \dots, \frac{m-2}{m-1}, 1\right\}$, conserving the notation used by Bailleul & Gras (1994).

Let a and b be two modal variables, Lagrange (1998) defines the non propensity index of a on b as:

$$S = \sum_{x \in I} \frac{\psi_a(x) (1 - \psi_b(x))}{n}$$
(2)

and to define the intensity of propensity, he considers a set E of size N sufficiently large, of which I is a random part. The random index is defined as:

$$Z = \sum_{x \in I} \frac{A_x \left(1 - B_x\right)}{n} \tag{3}$$

where $\{A_x\}_{x\in E}$ and $\{B_x\}_{x\in E}$ are sets of random variables with equal distributions, representative of the variables ψ_a and ψ_b respectively, for which the observed non propensity index s is one of its possible values. It is also proved that Z follows approximately a Normal distribution of mean $\mu = E[A_x(1-B_x)]$ and variance $\sigma^2 = \frac{E[A_x(1-B_x)]^2}{n}$, and therefore, for sufficiently large n, $\frac{Z-\mu}{\sigma}$ follows approximately a Normal distribution of mean zero and variance one.

Lagrange (1998) defines the propensity coefficient and the intensity of propensity of a on b as follows:

• Propensity coefficient:

$$\widetilde{q}(a,\overline{b}) = \frac{\sum_{x \in I} \frac{\psi_a(x) (1 - \psi_b(x))}{n} - m_a (1 - m_b)}{\sqrt{\frac{(s_a^2 + m_a^2) (s_b^2 + (1 - m_b)^2)}{n}}}$$
(4)

• Intensity of propensity:

$$\phi\left(-\widetilde{q}\left(a,\overline{b}\right)\right) = 1 - P\left(\frac{Z - m_a\left(1 - m_b\right)}{\sqrt{\frac{\left(s_a^2 + m_a^2\right)\left(s_b^2 + \left(1 - m_b\right)^2\right)}{n}}} < \widetilde{q}\left(a,\overline{b}\right)\right)$$
(5)

where:

- m_a (resp. m_b) is the empirical mean of ψ_a (resp. ψ_b),
- s_a^2 (resp. s_b^2) is the variance of ψ_a (resp. ψ_b), and
- ϕ is the standard normal distribution function.

For a risk level α (0 < α < 1), propensity of a on b will be assumed if $\phi\left(-\widetilde{q}\left(a,\overline{b}\right)\right) \geq 1-\alpha$.

2.2. Contribution of the Individuals

To define the typicality and the contribution of the individuals to a class or a path for the binary case, the approach defined by Ratsimba-Rajohn (1992) has been used, which presents the notion of respect of an implication among binary variables, in the following way:

$$\varphi_x \left(a, \overline{b} \right) = 1 \text{ if } [a \left(x \right) = 1, \ a \left(x \right) = 0] \text{ and } b \left(x \right) = 1,$$

$$\varphi_x \left(a, \overline{b} \right) = 0 \text{ if } a \left(x \right) = 1 \text{ and } b \left(x \right) = 0,$$

$$\varphi_x \left(a, \overline{b} \right) = p \text{ if } a \left(x \right) = 0 \text{ and } b \left(x \right) = 0,$$

where $\varphi_x(a, \overline{b})$ is the intensity of the implication of a on b for the individual x and p is a number between 0 and 1.

For modal variables, Bailleul & Gras (1994) only analyze the contribution of the individuals to a class or a path. They proceed in two steps:

First. The set of the possible values that the couples of variables can take is arranged according to the weights that have been assigned to each one of the modalities.

In this classification, they divide the couples of variables in two groups:

G1: The couples (a, b) that verify: $\psi_a(x) = \frac{m-j}{m-1}$ and $\psi_b(x) = \frac{i-j}{m-1}$, with $1 \le i \le (m-1)$ and $1 \le j \le i$, which contradict the implication $a \Rightarrow b$. The rank assigned to these couples of variables is given through the expression $rg(\psi_a(x), \psi_b(x)) = \frac{i(i-1)}{2} + j$.

G2: The couples (a, b) that respect the implication $a \Rightarrow b$ and verify $\psi_a(x) = \frac{m-i}{m-1}$ and $\psi_b(x) = \frac{m-i+j}{m-1}$, with $m \ge i \ge 1$, and decrement -1, and $0 \le j \le (i-1)$. The rank assigned to these couples of variables is given by the expression $rg(\psi_a(x), \psi_b(x)) = m^2 - \frac{i(i-1)}{2} + j + 1$.

Second. The value that measures the respect or the contradiction of the implication is related to the implication intensity.

Thus, the grade of adhesion of the individual x to the implication $a \Rightarrow b$ is defined as:

$$deg_{a,b}\left(x\right) = \frac{rg\left(\psi_{a}\left(x\right),\psi_{b}\left(x\right)\right) - 1}{m^{2} - 1}\varphi\left(a,\overline{b}\right)$$

$$\tag{6}$$

From this grade of adhesion, Bailleul & Gras (1994) define the distance d of an individual x to the class C formed by g sub-classes, for which the g generic couples (r_1, r_2, \ldots, r_g) and the g generic implications $(\varphi_1, \varphi_2, \ldots, \varphi_g)$ have been defined as (see Gras & Kuntz 2007):

$$d(x,C) = \left[\frac{1}{g}\sum_{i=1}^{g} \frac{\left(\varphi_i - deg_{r_i}(x)\right)^2}{1 - \varphi_i}\right]^{1/2}$$
(7)

and the contribution of an individual **x** to the class C as:

$$\gamma(x,C) = 1 - \frac{d(x,C)}{d(x_n,C)}$$
(8)

where x_n is the neutral individual, that is, the individual that assigns the value 0 to all the variables of the g generic couples r_i .

3. Classificatory Analysis for Modal Variables

In this section, an index to quantify the similarity among modal variables, the typicality and the contribution of the individuals to a class C is presented.

3.1. Similarity Index

Following the method to define the propensity coefficient, let us denote a set E of size N, with sufficiently large N, from which a random subset I of size n is selected, and a set $A = \{a_1, a_2, \ldots, a_p\}$ formed by p modal variables, each one of which can take values in the set $M = \{k \in N : 1 \leq k \leq m\}$. Specific weights are

assigned to each variable, from 1 (the strongest) to 0, in size fractions $\frac{1}{m-1}$ by means of the function ψ defined in (1).

Let $\{A_x\}_{x\in E}$ and $\{B_x\}_{x\in E}$ be two groups of random variables with equal distributions, representatives of the variables a and b with weights ψ_a and ψ_b respectively. Under the assumption of non a priori relationship between a and b, the 2N variables $\{A_x\}$, $\{B_x\}$ are independents and the N events $\{x \in I\}$ are independents with an occurrence probability equal to $\frac{n}{N}$.

Let us define the random variable:

$$\sum_{x \in I} \frac{A_x B_x}{n} \tag{9}$$

for which $\sum_{x \in I} \frac{\psi_a(x) \psi_b(x)}{n}$ is one of their possible values. This random variable follows approximately a Normal distribution of mean $\mu = E[A_x] E[B_x]$ and variance $\sigma^2 = \left(Var(A_x) + [E(A_x)]^2 \right) \left(Var(B_x) + [E(B_x)]^2 \right)$ and therefore, for n sufficiently large,

$$\frac{\sum_{x \in I} \frac{A_x B_x}{n} - \mu}{\sigma} \tag{10}$$

follows approximately a Normal distribution of mean zero and variance one.

The proof of this result is similar to the one presented by Lagrange (1998) for the random variable $\sum_{x \in I} \frac{A_x (1 - B_x)}{n}$, therefore it doesn't show up in this article.

The similarity index of a on b is defined as:

$$\widetilde{s}(a,b) = P\left(\frac{\sum_{x \in I} \frac{A_x B_x}{n} - m_a m_b}{\sqrt{\frac{(s_a^2 + m_a^2)(s_b^2 + m_b^2)}{n}}} < \frac{\sum_{x \in I} \frac{\psi_a(x)\psi_b(x)}{n} - m_a m_b}{\sqrt{\frac{(s_a^2 + m_a^2)(s_b^2 + m_b^2)}{n}}}\right)$$
(11)

where:

- m_a (m_b) is the empirical mean of ψ_a (ψ_b),
- $s_a^2(s_b^2)$ is the variance of $\psi_a(\psi_b)$,

which coincides with the similarity index defined for the binary case.

Let us prove this. For the binary case, there are just two modalities, this is, $M = \{1, 2\}$ and the unique assignable weights are 0 and 1, therefore:

•
$$\sum_{x \in I} \frac{A_x B_x}{n} = \frac{1}{n} Card \left(A \cap B\right)$$

Fundamental Concepts on Classification and Statistical Implicative Analysis

•
$$\sum_{x \in I} \frac{\psi_a(x) \psi_b(x)}{n} = \sum_{x \in I} \frac{a(x) b(x)}{n} = \frac{n_{a \wedge b}}{n}, \text{ where:}$$
$$n_{a \wedge b} = \{x \in I : \psi_a(x) = 1 \land \psi_b(x) = 1\} = \{x \in I : a(x) = 1 \land b(x) = 1\}$$

341

•
$$m_a = \frac{1}{n} \sum_{x \in I} \psi_a(x) = \frac{n_a}{n}$$
, where:
 $n_a = \{x \in I : \psi_a(x) = 1\} = \{x \in I : a(x) = 1\}$

•
$$s_a^2 = \hat{\sigma}^2(\psi_a) = \frac{1}{n} \left[\sum_{x \in I} \psi_a^2(x) - nm_a^2 \right] = m_a - m_a^2,$$

therefore, $s_a^2 + m_a^2 = m_a$, and

$$\begin{split} \widetilde{s}\left(a,b\right) &= P\left(\frac{\frac{Card\left(A \cap B\right)}{n} - \frac{n_{a}n_{b}}{n^{2}}}{\sqrt{\frac{n_{a}n_{b}}{n^{3}}}} < \frac{\frac{n_{a\wedge b}}{n} - \frac{n_{a}n_{b}}{n^{2}}}{\sqrt{\frac{n_{a}n_{b}}{n^{3}}}}\right) \\ &= P\left(\frac{Card\left(A \cap B\right) - \frac{n_{a}n_{b}}{n}}{\sqrt{\frac{n_{a}n_{b}}{n}}} < \frac{n_{a\wedge b} - \frac{n_{a}n_{b}}{n}}{\sqrt{\frac{n_{a}n_{b}}{n}}}\right) = s\left(a,b\right). \end{split}$$

3.2. Significant Nodes

For the determination of the significant nodes, a similar procedure to that of the binary case is used, but to define the initial and global pre-order on $A \times A$, the index $\tilde{s}(a, b)$ is used, this is, $G_{\tilde{s}}(\Omega) = \{((a, b), (c, d)) : \tilde{s}(a, b) < \tilde{s}(c, d)\}$. For more details, see (Zamora, Gregori & Orús 2009).

3.3. Typicality and Contribution to a Class C

The couple (a, b) so that:

$$\widetilde{s}(a,b) \ge \widetilde{s}(i,j), \quad \forall i \in A \text{ and } j \in B,$$

is called the generic couple of the class $C : A \sim B$. The number $\tilde{s}(a, b)$ is called generic similarity of C.

If the class C has g subclasses (C included), then there will be g generic similarities $(\tilde{s}_1, \tilde{s}_2, \ldots, \tilde{s}_g)$.

Based on the idea presented by Bailleul & Gras (1994), the grade of agreement between two modal variables a and b is defined. For that purpose:

First. The set of the possible values that the couples of variables can take is arranged, according to the weights that have been assigned to each modality.

The couples of this ordered set are divided into two groups:

G1: The couples (a, b) that verify the similarity $a \sim b$ or for which a(x) = b(x).

G2: The couples (a, b) that don't verify the similarity.

To establish a classification to the levels of similarity or non similarity, ranks are assigned to the couples of values that the couples of variables (a, b) can take for a given individual x, like it is shown below.

Ranks are placed on a $m \times m$ matrix (see Table 1 for the case m = 4), where the rank assigned to the cell (i, j), corresponding to the *i*-th modality of the variable a and the *j*-th modality of the variable *b*, is given by the expression:

$$Rg_{i,j} = \begin{cases} (m-j+1) + \frac{(i-1)(2m-i)}{2}, & \text{for } i \neq j \\ \\ \frac{m(m-1)}{2} + 1, & \text{for } i = j \end{cases}$$
(12)

with $1 \le i \le (m-1)$ and j decreasing from m to (i+1) with decrement (-1). It is also defined that $Rg_{i,j} = Rg_{j,i}$.

TA	BLE 1	: A	ssigr	ned 1	ang	es.
	i/j	1	2	3	4	
	1	7	3	2	1	
	2		$\overline{7}$	5	4	
	3			$\overline{7}$	6	
	4				7	

From this table we can appreciate that:

- The couple (1,4) corresponding to an individual for which $\psi_a(x) = 0$ and $\psi_b(x) = 1$ receives the lower rank, 1.
- For the rest of the cases, the cell receives a rank between 1 and 7, according to the level of discrepancy or disagreement between a(x) and b(x).

Second. The value that measures the respect or the contradiction to the full agreement is linked to the similarity index.

For that, the grade of agreement of the individual **x** to the relationship $a \sim b$ is defined as:

$$deg_{a,b}\left(x\right) = \frac{Rg_{i,j} - 1}{\frac{m\left(m-1\right)}{2}}\tilde{s}\left(a,b\right)$$
(13)

where a(x) = i and b(x) = j.

The typical subject is defined as the subject whose grade of agreement coincides with the value of the generic similarity for all the generic couples, this is, $deg_{r_i}(x) = \tilde{s}_i \ \forall r_i$.

The distance of an individual x to the class C formed by g sub-classes is then defined as:

$$d(x,C) = \left[\frac{1}{g} \sum_{i=1}^{g} \frac{(\tilde{s}_i - deg_{r_i}(x))^2}{1 - \tilde{s}_i}\right]^{1/2}$$
(14)

where (r_1, r_2, \ldots, r_g) are the g generic couples and $(\tilde{s}_1, \tilde{s}_2, \ldots, \tilde{s}_g)$ represents the g generic similarities.

The typicality is defined by the fact that certain individuals are representative of the population's behavior; this is, with a similarity index close to that of the formed class. The typicality of an individual x to the class C is calculated by the expression:

$$\gamma(x,C) = 1 - \frac{d(x,C)}{\max_{y \in I} d(y,C)}$$
(15)

Let us analyze the two extreme cases:

- If x is a typical subject, it should satisfy that $deg_{r_i} = \tilde{s}_i$ for all generic couples r_i , which implies that $Rg_{r_i} = \frac{m(m-1)}{2} + 1$, that is, a(x) = b(x) for all the generic couples $r_i = (a, b)$, and therefore, d(x, C) = 0 and $\gamma(x, C) = 1$.
- If x is in more disagreement with C, $d(x, C) = \max_{y \in I} d(y, C)$, which means that for all generic couples r_i , $deg_{r_i}(x) = 0$, that is, the rank that is assigned to all generic couples is 1, which occurs if for all generic couples $r_i = (a, b)$, a(x) = 1 and b(x) = m, or a(x) = m and b(x) = 1. However, if $d(x, C) = \max_{y \in I} d(y, C)$, then $\gamma(x, C) = 0$.

There are individuals more responsible than others in the formation of the class. The contribution is defined to evaluate this degree of responsibility. The contribution of an individual x to the class C is given by:

$$\widetilde{\gamma}(x,C) = 1 - \widetilde{d}(x,C) \tag{16}$$

where $\widetilde{d}(x,C) = \left[\frac{1}{g}\sum_{i=1}^{g} \left(1 - \deg_{r_i}^*(x)\right)^2\right]^{1/2}$ is the distance from the individual x to the class C and $\deg_{r_i}^* = \frac{Rg_{r_i} - 1}{\frac{m(m-1)}{2}}$.

Let us analyze the two extreme cases, the individual of higher contribution and the individual of lowest contribution to the formation of the class:

- If x is the individual of higher contribution to the formation of the class, then it satisfies that a(x) = b(x) and therefore $Rg_{r_i} = \frac{m(m-1)}{2} + 1$ for all generic couples $r_i = (a, b)$. Hence $deg_{r_i}^*(x) = 1$, $\forall r_i, \Rightarrow \tilde{d}(x, C) = 0 \Rightarrow \tilde{\gamma}(x, C) = 1$.
- If x is the individual of lowest contribution to the formation of the class, it satisfies that a(x) = 1 and b(x) = m or a(x) = m and b(x) = 1 and therefore $Rg_{r_i} = 1$ for all generic couples $r_i = (a, b)$. Hence $deg_{r_i}^*(x) = 0$, $\forall r_i$ and $\tilde{d}(x, C) = 1 \Rightarrow \tilde{\gamma}(x, C) = 0$.

4. Illustration of the Technique

The techniques presented have been programmed in language R (R Core Team 2014), in a file that will be denoted as ASIM.R. Two examples are presented, the first with binary data, to show the equivalence of the similarity index proposed with the existent one for the binary case. The second example deals with modal variables with more than two modalities, where the similarity index, the typicality and contribution of the individuals to the classes formed during the construction of the similarity tree are computed.

In the first example, seventeen binary variables are analyzed in a sample of 100 patients, from which 50 suffer from lung cancer and the other 50 don't. Before presenting the results of the application of both systems, CHIC and ASIM, it is necessary to point out that CHIC, for similarity analysis in the binary case, allows the user to choose the distribution (Poisson or Binomial); however, the results that CHIC shows, in both cases, correspond to the Poisson approximation to the Normal distribution. Given this difficulty, it is possible to compare the result shown by the CHIC for the binary case, with the result shown by ASIM for the modal case with two modalities.

The index proposed for the modal case with 2 modalities coincides with the index proposed for the binary case, in this case the random variable $\sum_{x \in I} \frac{A_x B_x}{n}$ follows approximately a Normal distribution of mean $\mu = E[A_x] E[B_x] = \frac{n_a n_b}{n}$ and variance $\sigma^2 = \left(Var(A_x) + [E(A_x)]^2 \right) \left(Var(B_x) + [E(B_x)]^2 \right) = \frac{n_a n_b}{n}$, distribution that coincides with the Poisson approximation to the Normal distribution.

The Figures 1 and 2 show the similarity indexes at level zero of the hierarchy. The coincidence of the results can be appreciated.

Console	Console E:/Larisa/Larisa_Personales/España/ASI_R/Modal_Agosto_19_2014/															
******	*******	*****	****	*****	*****	*****	*****	*****	****	*****	*****	*****	****	*****	*****	*****
		I	ndice	es de	simi	larida	ad al	nive	l cer	0						
******	*****	*****	****	****	*****	*****	****	****	****	****	****	****	****	****	****	*****
	CANCER EDAD	SEXO	PRB	PRM	PRN	PAG	PAP	POR	PCA	BEBE	FUMA	PVE	PRI	PS1	PS2	PS3
CANCER	1.00 0.99	0.93	0.45	0.41	0.66	0.96	0.58	0.99	0.94	0.93	0.98	0.28	0.95	0.94	0.53	0.02
EDAD	0.99 1.00	0.79	0.47	0.31	0.76	0.74	0.25	0.92	0.73	0.83	0.96	0.65	0.84	0.92	0.36	0.20
SEXO	0.93 0.79	1.00	0.36	0.44	0.73	0.78	0.10	1.00	0.53	1.00	0.95	0.01	0.95	0.82	0.55	0.08
PRB	0.45 0.47	0.36	1.00	0.00	0.00	0.68	0.53	0.28	0.55	0.20	0.50	0.40	0.81	0.51	0.28	0.91
PRM	0.41 0.31	0.44	0.00	1.00	0.00	0.28	0.42	0.17	0.26	0.65	0.27	0.30	0.36	0.80	0.33	0.54
PRN	0.66 0.76	0.73	0.00	0.00	1.00	0.57	0.56	0.97	0.74	0.70	0.78	0.83	0.28	0.14	0.89	0.04
PAG	0.96 0.74	0.78	0.68	0.28	0.57	1.00	0.93	0.89	0.59	0.81	0.83	0.86	0.90	0.48	0.56	0.38
PAP	0.58 0.25	0.10	0.53	0.42	0.56	0.93	1.00	0.12	0.79	0.08	0.56	0.91	0.39	0.80	0.26	0.72
POR	0.99 0.92	1.00	0.28	0.17	0.97	0.89	0.12	1.00	0.91	0.98	0.96	0.08	0.98	0.97	0.34	0.10
PCA	0.94 0.73	0.53	0.55	0.26	0.74	0.59	0.79	0.91	1.00	0.40	0.86	0.33	0.84	0.92	0.34	0.24
BEBE	0.93 0.83	1.00	0.20	0.65	0.70	0.81	0.08	0.98	0.40	1.00	0.95	0.05	0.83	0.82	0.49	0.15
FUMA	0.98 0.96	0.95	0.50	0.27	0.78	0.83	0.56	0.96	0.86	0.95	1.00	0.39	0.85	0.81	0.52	0.12
PVE	0.28 0.65	0.01	0.40	0.30	0.83	0.86	0.91	0.08	0.33	0.05	0.39	1.00	0.33	0.21	0.60	0.64
PRI	0.95 0.84	0.95	0.81	0.36	0.28	0.90	0.39	0.98	0.84	0.83	0.85	0.33	1.00	0.75	0.43	0.36
PS1	0.94 0.92	0.82	0.51	0.80	0.14	0.48	0.80	0.97	0.92	0.82	0.81	0.21	0.75	1.00	0.00	0.08
PS2	0.53 0.36	0.55	0.28	0.33	0.89	0.56	0.26	0.34	0.34	0.49	0.52	0.60	0.43	0.00	1.00	0.00
PS3	0.02 0.20	0.08	0.91	0.54	0.04	0.38	0.72	0.10	0.24	0.15	0.12	0.64	0.36	0.08	0.00	1.00

FIGURE 1: Similarity indexes at the level zero of the hierarchy with ASIM.R.

The Figures 3 and 4 show the classes formed during the construction of the similarity tree and its corresponding similarity indexes, obtained by ASIM.R.

Similarity E:\Larisa\España\ASI_R\Articulo_Modal\Datos\Cancer Pulmon.csv 1:2																				
Eile information	Г		-		-															
Frequency of c	L			18	/	/ /	/ /	/ /	/ /	/ /	/ /	/ /	/ /	/ /	/	/ /	/ /	/ /		7/7/
- Correlation coe	L		1	\$¥/2	\$//s	\$//\$	>//s	*//3	*//3	0//3	3//5	\$ // 5	s // 2	\$//:		~ //s	>//z	1/2	v // 8	2/
Index	Ы	CANCER	1.00	0.99	0.93	0.45	0.41	0.66	0.98	0.58	0.99	0.94	0.93	0.98	0.28	0.95	0.94	0.53	0.02	í
Value	Ľ	EDAD	0.99	1.00	0.79	0.47	0.31	0.76	0.74	0.25	0.92	0.73	0.83	0.96	0.65	0.84	0.92	0.38	0.20	1
	Ľ	SEXO	0.93	0.79	1.00	0.36	0.44	0.73	0.78	0.10	1.00	0.53	1.00	0.95	0.01	0.95	0.82	0.55	0.08	1
	Ľ	PRB	0.45	0.47	0.36	1.00	0.00	0.00	0.68	0.53	0.28	0.55	0.20	0.50	0.40	0.81	0.51	0.28	0.91	1
	Ľ	PRM	0.41	0.31	0.44	0.00	1.00	0.00	0.28	0.42	0.17	0.28	0.65	0.27	0.30	0.38	0.80	0.33	0.54	1
	Ľ	PRN	0.66	0.78	0.73	0.00	0.00	1.00	0.57	0.58	0.97	0.74	0.70	0.78	0.83	0.28	0.14	0.89	0.04	
	Ľ	PAG	0.96	0.74	0.78	0.68	0.28	0.57	1.00	0.93	0.89	0.59	0.81	0.83	0.86	0.90	0.48	0.58	0.38	1
	Ľ	PAP	0.59	0.25	0.10	0.52	0.42	0.58	0.92	1.00	0.12	0.00	0.01	0.56	0.91	0.39	0.90	0.00	0.72	
	Ľ	POP	0.00	0.92	1.00	0.00	0.17	0.97	0.00	0.12	1.00	0.91	0.00	0.98	0.09	0.99	0.97	0.24	0.10	
	Ľ	PCA	0.00	0.72	0.52	0.55	0.28	0.74	0.59	0.72	0.91	1.00	0.40	0.00	0.00	0.00	0.97	0.24	0.10	
	Ľ	PCA	0.54	0.73	1.00	0.00	0.20	0.74	0.05	0.75	0.91	0.40	1.00	0.00	0.35	0.04	0.92	0.34	0.24	
	Ľ	DEDE	0.93	0.65	1.00	0.20	0.05	0.70	0.01	0.08	0.96	0.40	1.00	0.95	0.05	0.65	0.02	0.49	0.15	
	Ľ	FUMA	0.98	0.90	0.95	0.50	0.27	0.78	0.65	0.50	0.30	0.80	0.95	1.00	0.39	0.85	0.01	0.52	0.12	
	Ľ	PVE	0.28	0.00	0.01	0.40	0.30	0.85	0.80	0.91	0.08	0.33	0.05	0.39	1.00	0.33	0.21	0.00	0.04	
	Ľ	PRI	0.95	0.84	0.95	0.81	0.38	0.28	0.90	0.39	0.98	0.84	0.83	0.85	0.33	1.00	0.75	0.43	0.35	
	Ľ	PS1	0.94	0.92	0.82	0.51	0.80	0.14	0.48	0.80	0.97	0.92	0.82	0.81	0.21	0.75	1.00	0.00	0.08	
		PS2	0.53	0.38	0.55	0.28	0.33	0.89	0.58	0.28	0.34	0.34	0.49	0.52	0.60	0.43	0.00	1.00	0.00	
		PS3	0.02	0.20	0.08	0.91	0.54	0.04	0.38	0.72	0.10	0.24	0.15	0.12	0.64	0.38	0.08	0.00	1.00	i i

FIGURE 2: Similarity indexes at the level zero of the hierarchy, obtained by CHIC.

Console	E:/Lari	sa/Larisa_Personales/España/ASI_R/Modal_Agosto_19_2014/ 🔅
*****	*****	**************************************
Nivel Nivel Nivel Nivel Nivel Nivel Nivel Nivel	1 : 2 : 3 : 4 : 5 : 6 : 7 : 8 :	"(SEX0,BEBE)" "((SEX0,BEBE),POR)" "((CANCER,EDAD)" "((CANCER,EDAD),FUMA)" "(PAG,PAP)" "(((SEX0,BEBE),POR),PRI)" "(PCA,PS1)"
Nivel Nivel Nivel Nivel Nivel Nivel Nivel	9 : 10 : 11 : 12 : 13 : 14 : 15 :	"(((CANCER,EDAD),FUMA),(((SEXO,BEBE),POR),PRI))" "(PRN,PS2)" "((PAG,PAP),PVE)" "((((CANCER,EDAD),FUMA),(((SEXO,BEBE),POR),PRI)),(PCA,PS1))" "((((CANCER,EDAD),FUMA),(((SEXO,BEBE),POR),PRI)),(PCA,PS1)),(PRN,PS2))" "(((((CANCER,EDAD),FUMA),(((SEXO,BEBE),POR),PRI)),(PCA,PS1)),(PRN,PS2)),((PAG,PAP),PVE))" "(((PR,PS3),PRM)"

FIGURE 3: Classes obtained by running ASIM.R.

Figure 5 shows the same results, by executing the software CHIC. The coincidence of both results can be appreciated.

Figures 6a and 6b show the similarity trees obtained by means of each program. Again the results are coincident.

Through the running of the above example and from the execution of other datasets, we can say that the execution times between CHIC and ASIM.R are similar.

Next, the results of a modal variable case with five modalities are presented. To accomplish this, the file "DatosMusicaModalII.txt" is used. This file corresponds to the degree of preference of twenty students to fifteen different kinds of music. The degree of preference is assessed by a scale from 1 to 5, 5 being the highest value of preference.

Console	E:/Laris	a/Larisa_Personales/España/ASI_R/Modal_Agosto_19_2014/ 🔗
*****	*****	******
		Clasificacion por niveles
*****	****	******
		indice
Nivel	1 :	0.999074
Nivel	2 :	0.997657
Nivel	3 :	0.988194
Nivel	4 :	0.955888
Nivel	5:	0.931400
Nivel	6 :	0.931062
Nivel	7 :	0.921997
Nivel	8 :	0.914071
Nivel	9 :	0.907252
Nivel	10 :	0.893640
Nivel	11 :	0.819734
Nivel	12 :	0.691713
Nivel	13 :	0.584519
Nivel	14 :	0.304926
Nivel	15 :	0.287690

FIGURE 4: Similarity indexes obtained by running ASIM.R.

¹⁸ /wr Similiarity E\Larisa\España\ASL_R\Articulo_Moda\Datos\Cancer Pulmon.csv 1:1	- • ×
Classification at level : 1 : (SEXO BEBE) similarity : 0.999074	^
Classification at level : 2 : ((SEXO BEBE) POR) similarity : 0.997657	
Classification at level : 3 : (CANCER EDAD) similarity : 0.988194	
Classification at level : 4 : ((CANCER EDAD) FUMA) similarity : 0.955888	
Classification at level : 5 : (PAG PAP) similarity : 0.9314	
Classification at level : 6 : (((SEXO BEBE) POR) PRI) similarity : 0.931062	
Classification at level : 7 : (PCA PS1) similarity : 0.921997	
Classification at level : 8 : (PRB PS3) similarity : 0.914071	E
Classification at level : 9 : (((CANCER EDAD) FUMA) (((SEXO BEBE) POR) PRI)) similarity : 0.907252	
Classification at level : 10 : (PRN PS2) similarity : 0.89364	
Classification at level : 11 : ((PAG PAP) PUE) similarity : 0.819734	
Classification at level : 12 : ((((CANCER EDAD) FUNA) (((SEXO BEBE) POR) PRI)) (PCA PS1)) similarity : 0.691714	
Classification at level : 13 : (((((CANCER EDAD) FUMA) (((SEXO BEBE) POR) PRI)) (PCA PS1)) (PRN PS2)) similarity : 0.58452	
Classification at level : 14 : (((((CANCER EDAD) FUMA) (((SEXO BEBE) POR) PRI)) (PCA PS1)) (PRN PS2)) ((PAG PAP) PUE)) similarit	y : 0.304926
Classification at level : 15 : ((PRB PS3) PRM) similarity : 0.28769	-
m m	►

FIGURE 5: Classes and similarity indexes obtained by executing CHIC.

Figures 7, 8 and 9 show the similarity indexes obtained at the level zero of the hierarchy, the grouping of variables in each level of this hierarchy, and the indexes to which they were joined.

The similarity index of the class 1 formed by variables HIP and PUN (see Figures 8 and 9) is 0.934143; it means that HIP and PUN are the most similar music types in the population of which the sample under study has been selected.

The second class is formed by the variables JAZ and HEA. This class also shows a big similarity index (0.922925), but smaller than the previous one (see Figures 8 and 9). Then we can say that the preference for these music types is high, but smaller than for the previous group (HIP and PUN). The other classes are interpreted in a similar way.



FIGURE 6: Similarity trees.



FIGURE 7: Similarity indexes at the level zero of the hierarchy.



FIGURE 8: Classes obtained.

*****	***		***********	*******
				Clasificacion por niveles
*****	***	***	*******	*********************************
			indice	
Nivel	1	:	0.934143	
Nivel	2	:	0.922925	
Nivel	3	:	0.900668	
Nivel	4	:	0.889664	
Nivel	5	:	0.851791	
Nivel	6	:	0.814453	
Nivel	7	:	0.644702	
Nivel	8	:	0.613585	
Nivel	9	:	0.497993	
Nivel	10) :	0.327310	
Nivel	11	:	0.285067	
Nivel	12	:	0.042301	

FIGURE 9: Similarity indexes obtained.

Figure 10 presents the similarity tree for the music data classification shown in Figures 8 and 9.



The typicality of the individuals and the optimal group are shown for each one of the formed classes. In Figure 11 we can see that individuals for which the value of the typicality in a class is 0 exist, for example, the individual 7 in the class 3. These individuals are in more disagreement, this is, the values taken for the variables in each generic couple for these individuals are respectively the minimum (modality 1) and maximum (modality 5) value or vice versa.

On the contrary, there are individuals whose typicality value in a class is 1, for example, the individual 1 of the class 3. The values taken for the variables in each generic couple for these individuals are the same. The degree of agreement among these variables for this individual coincides with the similarity index, being hence the most typical individual.

Figures 13 and 14 show the contributions of the individuals and the optimal group. The results related to the contribution can be interpreted as follows. When the contribution of an individual to a class is 1, for example, individual 1 of class

Console	E:/Larisa/Lari	isa Personale	es/España/A	SI_R/Modal	Agosto_19_20)14/ 🖘				
*******	*******	*******	********	*******	******	********	********	********	********	*******
		Tipica	lidades de	e los ind	ividuos a	las clase	5			
*******	********	********	********	********	*********	********	********	*********	*********	********
	Clase[1]	Clase[2]	Clase[3]	Clase[4]	Clase[5]	Clase[6]	Clase[7]	Clase[8]	Clase[9]	Clase[10]
Indiv_1	1.00000	1.00000	1.00000	0	1.0000000	1	1	0	0.3110467	0.5715400
Indiv_2	1.00000	1.00000	1.00000	1	1.0000000	1	1	1	1.0000000	0.9330531
Indiv_3	1.00000	1.00000	1.00000	0	1.0000000	1	1	1	0.2309112	0.5715400
Indiv_4	1.00000	1.00000	1.00000	1	1.0000000	1	0	0	1.0000000	0.5715400
Indiv_5	1.00000	1.00000	0.25000	1	1.0000000	1	1	1	0.4527748	0.9330531
Indiv_6	1.00000	1.00000	1.00000	1	1.0000000	1	1	0	0.4637086	0.8928850
Indiv_7	1.00000	1.00000	0.00000	0	1.0000000	1	1	0	0.0000000	1.0000000
Indiv_8	1.00000	0.25000	0.34375	1	0.4696699	1	1	1	0.5211780	0.4696699
Indiv_9	1.00000	1.00000	1.00000	1	0.4696699	1	1	1	0.4637086	0.3799092
Indiv 10	1,00000	1.00000	1.00000	1	1,0000000	1	1	1	0.8659272	1,0000000

FIGURE 11: Typicalities of the first 10 individuals to the first 10 classes.

```
Console E:/Larisa/Larisa Personales/España/ASI_R/Modal_Agosto_19_2014/

Grupos optimales de individuos para las tipicalidades:

Optimal Group to Clase[1]: (HIP,PUN)

17 1 2 3 4 5 6 7 8 9 10 12 13 14 15 16 18 19 20

Card :19

Optimal Group to Clase[2]: (JAZ,HEA)

20 1 2 3 4 5 6 7 9 10 11 12 14 16 17 18 19

Card :17

Optimal Group to Clase[3]: (OPE,CLA)

1 2 3 4 6 9 10 11 12 13 14 16 17 18 19 20

Card :16

Optimal Group to Clase[4]: (FLA,RAP)

2 4 5 6 8 9 10 11 12 13 14 15 16 17 18 19

Card :16

Optimal Group to Clase[5]: ((JAZ,HEA),REG)

1 2 3 4 5 6 7 10 11 12 14 16 18 19

Card :14

Optimal Group to Clase[6]: (MAQ,POP)

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20

Card :20
```

FIGURE 12: Optimal group for the first 6 classes.

1, the values of the variables for these individuals, in each generic couple, are equal. This individual evaluates in a similar way his preference for the music types that form the vector of generic similarities for the class and therefore it is more responsible in the creation of the class.

If the contribution of an individual to a class is 0, for example, individual 1 of class 4, then its distance to the class is the biggest possible (1). This is the individual that contributes the least to the formation of the class (FLA, RAP), doesn't like the music FLA and assigns the biggest value possible to the preference for the music RAP. It is the individual of lowest contribution for declaring FLA and RAP as similar music types.

For space reasons, only a part of the typicalities, contributions and optimal group are shown.

Console	E:/Larisa/Lari	isa Personale	es/España/A	SI_R/Modal	Agosto_19_2	014/ 🖘 **********		********	*********
*******	********	********	********	********	********	********	********	********	*********
	Clase[1]	Clase[2]	Clase[3]	Clase[4]	Clase[5]	Clase[6]	Clase[7]	Clase[8]	Clase[9]
Indiv_1	1.00000	1.00000	1.00000	0	1.0000000	1	1	0	0.4156445
Indiv_2	1.00000	1.00000	1.00000	1	1.0000000	1	1	1	1.0000000
Indiv_3	1.00000	1.00000	1.00000	0	1.0000000	1	1	1	0.3094292
Indiv_4	1.00000	1.00000	1.00000	1	1.0000000	1	0	0	1.0000000
Indiv_5	1.00000	1.00000	0.25000	1	1.0000000	1	1	1	0.5669873
Indiv_6	1.00000	1.00000	1.00000	1	1.0000000	1	1	0	0.4226497
Indiv_7	1.00000	1.00000	0.00000	0	1.0000000	1	1	0	0.1835034
Indiv_8	1.00000	0.25000	0.34375	1	0.4696699	1	1	1	0.6211139
Indiv_9	1.00000	1.00000	1.00000	1	0.4696699	1	1	1	0.4226497
Indiv 10	1.00000	1.00000	1.00000	1	1.0000000	1	1	1	0.8556624

FIGURE 13: Contribution of the first 10 individuals to the first 9 classes.

```
Console E:/Larisa/Larisa Personales/España/ASI_R/Modal_Agosto_19_2014/ >

Grupos optimales de individuos para las contribuciones:

Optimal Group to Clase[1]: (HIP,PUN)

17 1 2 3 4 5 6 7 8 9 10 12 13 14 15 16 18 19 20

Card :19

Optimal Group to Clase[2]: (JAZ,HEA)

20 1 2 3 4 5 6 7 9 10 11 12 14 16 17 18 19

Card :17

Optimal Group to Clase[3]: (OPE,CLA)

1 2 3 4 6 9 10 11 12 13 14 16 17 18 19 20

Card :16

Optimal Group to Clase[4]: (FLA,RAP)

2 4 5 6 8 9 10 11 12 13 14 15 16 17 18 19

Card :16

Optimal Group to Clase[5]: ((JAZ,HEA),REG)

1 2 3 4 5 6 7 10 11 12 14 16 18 19

Card :14

Optimal Group to Clase[6]: (MAQ,POP)

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20

Card :20
```

FIGURE 14: Optimal group for the first 6 classes.

5. Conclusions

In this work we proposed a similarity index among modal variables, and expressions for the calculation of the typicalities and contributions of the individuals to a class, for these types of variables. It has been demonstrated, theoretically and numerically, that the expression for the calculation of the similarity index for modal variables with two modalities coincides with the expression for the binary case, thus the index proposed is a generalization of the binary case. Following the idea of Bailleul & Gras (1994), we propose an expression to calculate the grade of agreement between two modal variables and from this, expressions are given to calculate the typicalities and contributions of the individuals to the formed classes.

[Received: June 2014 — Accepted: December 2014]

References

- Bailleul, M. & Gras, R. (1994), 'L'implication statistique entre variables modales', Mathématiques et Sciences Humaines 128, 47–51.
- Couturier, R. (2008), 'CHIC: Cohesive hierarchical implicative classification', Studies in Computational Intelligence (SCI) 127, 41–53.
- Gras, R. (2000), 'Les fondements de l'analyse statistique implicative', *Quaderni di Ricerca in Didattica* 9, 187–208.
- Gras, R. & Kuntz, K. (2007), Analyse statistique implicative (ASI), en réponse à des problèmes fondateurs, in 'Apports Théoriques à l'Analyse Statistique Implicative et Applications', Université Jaume I, Castellón, pp. 15–40.
- Lagrange, J. B. (1998), 'Analyse implicative d'un ensemble de variables numériques; application au traitement d'un questionnaire à résponses modales ordonnées', *Revue de Statistique Appliquée* **46**(1), 71–93.
- R Core Team (2014), R: A Language and Environment for Statistical Computing, R Foundation for Statistical Computing. *http://CRAN.R-project.org
- Ratsimba-Rajohn, H. (1992), Contribution à l'étude de la hiérarchie implicative, application à l'analyse de la gestión didactique des phénomènes d'ostension et de contradiction, Master Thesis, Université de Rennes, France.
- Zamora, L., Gregori, P. & Orús, P. (2009), Conceptos Fundamentales del Análisis Estadístico Implicativo (ASI) y su Soporte Computacional CHIC, Departamento de Matemáticas, Universitat Jaume I de Castellón, España.