# A Topological Multiple Correspondence Analysis

Rafik Abdesselam

COACTIS-ISH Management Sciences Laboratory - Human Sciences Institute,
University of Lyon, Lumière Lyon 2
Campus Berges du Rhône, 69635 Lyon Cedex 07, France
(E-mail: `rafik.abdesselam@univ-lyon2.fr`)
(http://perso.univ-lyon2.fr/~rabdesse/fr/)

**Abstract**

Topological Multiple Correspondence Analysis (TMCA) studies a group of categorical variables defined on the same set of individuals. Its a topological method of data analysis that consists of exploring, analyzing and representing the associations between several qualitative variables in the context of multiple correspondence anal-ysis. It compares and classifies proximity measures to select the best one according to the data under consideration, then analyzes, interprets and visualizes with graphic representations, the possible associations between several categorical variables relat-ing to, the known problem of Multiple Correspondence Analysis (MCA). Based on the notion of neighborhood graphs, some of these proximity measures are more-or-less equivalent. A topological equivalence index between two measures is defined and statistically tested according to the degree of description of the associations between the modalities of these qualitative variables.

We compare proximity measures and propose a topological criterion for choosing the best association measure, adapted to the data considered, from among some of the most widely used proximity measures for categorical data. The principle of the proposed approach is illustrated using a real data set with conventional proximity measures for binary variables from the literature. The first step is to find the proximity measure that can best adapted to the data; the second step is to use this measure to perform the TMCA.

*Keywords:* Burt table; proximity measure; neighborhood graph; adjacency ma-trix; topological equivalence; graphical representations.

## 1 Introduction

Similarity measures play an important role in many areas of data analysis. The results of any operation involving structuring, clustering or classifying objects are strongly dependent on the proximity measure chosen. The user has to select one measure among many existing ones. Yet, according to the notion of topological equivalence chosen, some measures are more-or-less equivalent. The concept of topological equivalence uses the basic notion of local neighbor-hood. We define the topological equivalence between two proximity measures, in the context of association between several categorical variables, through the topological structure induced by each measure.

Multiple correspondence analysis (MCA) is an important methodology among factorial techniques due to the extent of its field of application. It allows us, among others things, to describe large binary tables, such as socio-economic surveys, and usually answers questions on modalities.

This method is a generalization of correspondence analysis (CA); it concerns the relations between or within a set of p ( $p > 2$ ) qualitative variables simultaneously observed on n individuals. Generally the variables are homogeneous in the sense that they revolve around a particular theme.

In order to understand and act on situations that are represented by a set of objects, very often we are required to compare them. Humans perform this comparison subconsciously using the brain. In the context of artificial intelligence, however, we should be able to describe how the machine might perform this comparison. In this context, one of the basic elements that must be specified is the proximity measure between objects.

Certainly, application context, prior knowledge, data type and many other factors can help in identifying the appropriate measure. For instance, if the objects to be compared are described by Boolean vectors, we can restrict our comparisons to a class of measures specifically devoted to this type of data. However, the number of candidate measures may still remain quite large. Can we consider that all those measures remaining are equivalent and just pick one of them at random? Or are there some that are equivalent and, if so, to what extent? This information might interest a user when seeking a specific measure. For instance, in information retrieval, choosing a given proximity measure is an important issue. We effectively know that the result of a query depends on the measure used. For this reason, users may wonder which one is more useful? Very often, users try many of them, randomly or sequentially, seeking a "suitable" measure. If we could provide a framework that allows the user to compare proximity measures in order to identify those that are similar, they would no longer need to try out all measures.

The present study proposes a new framework for comparing proximity measures in order to choose the best one in the context of association between a set of qualitative variables. The aim is to establish a TMCA.

We deliberately ignore the issue of the appropriateness of the proximity measure, as it is still an open and challenging question currently being studied. The comparison of proximity measures can be analyzed from various angles.

The comparison of objects, situations or ideas is an essential task in order to assess a situation, to rank preferences, to structure a set of tangible or abstract elements, and so on. In a word, to understand and act, we have to compare. These comparisons that the brain naturally performs, however, must be clarified if we want them to be done by a machine. For this purpose, we use proximity measures. A proximity measure is a function which measures the similarity or dissimilarity between two objects within a set. These proximity measures have mathematical properties and specific axioms. But are such measures equivalent? Can they be used in practice in an undifferentiated way? Do they produce the same learning database that will serve to find the membership class of a new object? If we know that the answer is negative, then how do we decide which one to use? Of course, the context of the study and the type of data being considered can help in selecting a few possible proximity measures, but which one should we choose from this selection as the best measure for summarizing the association?

We find this problematic also in the context of TMCA. The eventual links or associations between all the qualitative variables partly depends on the learning database being used. The results of multiple correspondence analysis can change according to the selected proximity measure.

Several studies on the topological equivalence of proximity measures have been proposed, [3] [15] [4] [12] [21], also in discrimination context [2], but none of these propositions has an association objective between several categorical variables. An approach in the case of association between two qualitative variables has been proposed in [1].

Therefore, this article focuses on how to construct the best adjacency matrix induced by a proximity measure, taking into account the association between all the modalities of the qualitative variables.

This paper is organized as follows. In section 2, after recalling the basic notions of structure, graph and topological equivalence, we present the proposed method, how to build an adjacency matrix associated with a proximity measure in the context of association between several qualitative variables, how to compare and statistically test the degree of topological equivalence between proximity measures and how to select the best measure to describe multiple associations. Section 3 presents an illustrative example using real data. The conclusion of this work is given in section 4.

Table 7, shown in the appendix, summarizes some classic proximity measures used for binary data [20], we give on $\{0, 1\}^n$ the definition of 22 of them.

We assume that we have at our disposal $\{x^k; k = 1, .., p\}$ a set of $p > 2$ qualitative variables, partitions of $n = \sum_{k=1}^{p} n_k$ individuals-objects into $m_k$ modalities-subgroups. The interest lies in whether there is a topological association between all these variables. Let us denote:
- $X_k = X_{(n,m_k)}$ the disjonctif table, data matrix associated to the $m_k$ dummy variables of the qualitative variable $x^k$ with $n$ rows-objects and $m_k$ columns-modalities, we check that $\Sigma_{k=1}^{m_k} x_i^k = 1, \ \forall_i$ and $\Sigma_{i=1}^{n} x_i^k = n_k$
- $X_{(n,m)} = [X_1|X_2| \cdots |X_p]$ the indicator matrix, juxtaposition of the $p$ binary tables $X_k$, with $n$ rows-objects and $m = \sum_{k=1}^{p} m_k$ columns-modalities, we check that $\Sigma_{k=1}^{m_k} x_i^k = p, \ \forall_i$ and $\Sigma_{i=1}^{n} \Sigma_{k=1}^{m_k} x_i^k = np$.

An alternative coding of such data is as a Burt matrix, a square symmetric modalities-by-modalities matrix formed from all two-way contingency tables of pairs of variables, including on the block diagonal the cross-tabulations of each variable with itself.
- $\mathcal{B}_{(m,m)} = {}^t X \, X$ the symmetric Burt matrix of the two-way cross-tabulations of the $p$ variables,
- $W_{(m,m)} = diag[\mathcal{B}]$ is the diagonal marginal frequency matrices.
- $U = \mathbb{1}_m \, {}^t \mathbb{1}_m$ is the $m \times m$ matrix of 1s, $I_m$ the $m \times m$ identity matrix where $\mathbb{1}_m$ denotes the m indicator vector of 1s and $\mathbb{1}_n$ the n indicator vector of 1s.

The dissimilarity matrices associated with proximity measures are computed from data given by the Burt table $\mathcal{B}$. The attributes of any two points' modalities' $x^k$ and $x^l$ in $\{0, 1\}^n$ of the proximity measures can be easily written

and calculated from the following matrices. Computational complexity is thus considerably reduced.

- $A_{(m,m)} = (a_{kl}) = \mathcal{B}$, the Burt matrix

whose element, $a_{kl} = |x^k \cap x^l| = \sum_{i=1}^{n} x_i^k x_i^l$ is the number of attributes common to both points $x^k$ and $x^l$,

- $B_{(m,m)} = (b_{kl}) = {}^t X \, (\mathbb{1}_n \, {}^t\mathbb{1}_m - X) = {}^t X \, \mathbb{1}_n \, {}^t\mathbb{1}_m - {}^t X \, X$
$$= W \, \mathbb{1}_m \, {}^t\mathbb{1}_m - A = W \, U - A$$

whose element, $b_{kl} = |X^k - X^l| = |X^k \cap \overline{X^l}| = \sum_{i=1}^{n} x_i^k(1 - x_i^l)$ is the number of attributes present in $x^k$ but not in $x^l$,

- $C_{(m,m)} = (c_{kl}) = {}^t(\mathbb{1}_n \, {}^t\mathbb{1}_m \, - \, X) \, X = {}^t(\mathbb{1}_n \, {}^t\mathbb{1}_m) \, X - {}^t X \, X$
$$= \mathbb{1}_m \, {}^t\mathbb{1}_n \, X \, - \, {}^t X \, X = UW - A$$

whose element, $c_{kl} = |X^l - X^k| = |X^l \cap \overline{X^k}| = \sum_{i=1}^{n} x_i^l(1 - x_i^k)$ is the number of attributes present in $x^l$ but not in $x^k$.

- $D_{(m,m)} = (d_{kl}) = {}^t(\mathbb{1}_n \, {}^t\mathbb{1}_m \, - \, X) \, (\mathbb{1}_n \, {}^t\mathbb{1}_m \, - \, X)$
$$= \mathbb{1}_m \, {}^t\mathbb{1}_n \, \mathbb{1}_n \, {}^t\mathbb{1}_m \, - \, \mathbb{1}_m \, {}^t\mathbb{1}_n \, X \, - \, {}^t X \, \mathbb{1}_n \, {}^t\mathbb{1}_m \, + \, {}^t X \, X$$
$$= n\mathbb{1}_m \, {}^t\mathbb{1}_m \, - \, UW - WU + A = nU - UW - WU + A$$
$$= nU - (A + B + C)$$

whose element, $d_{kl} = |\overline{X^k} \cap \overline{X^l}| = \sum_{i=1}^{n}(1 - x_i^k)(1 - x_i^l)$ is the number of attributes in neither $x^k$ or $x^l$.

$X^k = \{i/x_i^k = 1\}$ and $X^l = \{i/x_i^l = 1\}$ are the sets of attributes present in data point-modality $x^k$ and $x^l$ respectively, and $|.|$ the cardinality of a set.

The attributes are linked by the relation:
$$\forall k = 1, p \, ; \, \forall l = 1, p \quad a_{kl} + b_{kl} + c_{kl} + d_{kl} = n.$$

Together, the four dependent quantities $a_{kl}, b_{kl}, c_{kl}$ and $d_{kl}$ can be used to construct the $2 \times 2$ contingency table presented in Table 1, where the information can be summarized by an index of similarity (affinity, resemblance, association, coexistence). As a general symbol for a similarity coefficient the capital letter S will be used. A list of 22 similarity coefficients is given in Table 7 in Appendix.

**Table 1.** The four depedent quantities between two binary modalities $x^{kr}$ and $x^{ls}$

|  | $x_i^{ls} = 1$ | $x_i^{ls} = 0$ | Total |
|---|---|---|---|
| $x_i^{kr} = 1$ | $a_{kl}$ | $b_{kl}$ | $a_{kl} + b_{kl}$ |
| $x_i^{kr} = 0$ | $c_{kl}$ | $d_{kl}$ | $c_{kl} + d_{kl}$ |
| Total | $a_{kl} + c_{kl}$ | $b_{kl} + d_{kl}$ | n |

## 2  Topological Correspondence

Topological equivalence is based on the concept of the topological graph also referred to as the neighborhood graph. The basic idea is actually quite simple: two proximity measures are equivalent if the corresponding topological

graphs induced on the set of objects remain identical. Measuring the similarity between proximity measures involves comparing the neighborhood graphs and measuring their similarity. We will first define more precisely what a topological graph is and how to build it. Then, we propose a measure of proximity between topological graphs that will subsequently be used to compare the proximity measures.

Consider a set $E = \{x^{11}, \ldots, x^{1m_1}, \ldots, x^{p1}, \ldots, x^{pm_p}\}$ of $m = \sum_{j=1}^{p} m_j$ modalities in $\{0,1\}^n$, associated with the $p$ qualitative variables $x^j$ with $m_j$ modalities. We can, by means of a proximity measure $u$, define a neighborhood relationship $V_u$ to be a binary relationship on $E \times E$. There are many possibilities for building this neighborhood binary relationship.
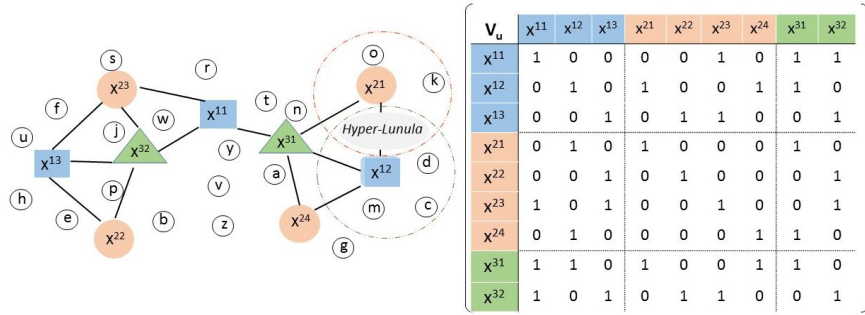
Thus, for a given proximity measure $u$, we can build a neighborhood graph on a set of objects-modalities, where the vertices are the modalities and the edges are defined by a property of the neighborhood relationship.

Many definitions are possible to build this binary neighborhood relationship. One can choose the Minimal Spanning Tree (MST) [10], the Gabriel Graph (GG) [14] or, as is the case here, the Relative Neighborhood Graph (RNG) [18].

For any given proximity measure $u$, we construct the associated adjacency binary symmetric matrix $V_u$ of order $m = \sum_{j=1}^{p} m_j$, where, all pairs of neighboring modalities $(x^{kr}, x^{ls})$, where $k, l = 1, p$ ; $r = 1, m_k$ and $s = 1, m_l$, satisfy the following RNG property.

*Property 1.* Relative Neighborhood Graph (RNG)

$$
\begin{cases}
V_u(x^{kr}, x^{ls}) = 1 & if \ \ u(x^{kr}, x^{ls}) \ \leq \ \max[u(x^{kr}, x^{qt}), u(x^{qt}, x^{ls})]; \\
& \forall x^{kr}, \ x^{ls}, \ x^{qt} \in E, \ x^{qt} \neq x^{kr} \ \ and \ \ x^{qt} \neq x^{ls} \\
V_u(x^{kr}, x^{ls}) = 0 & otherwise
\end{cases}
$$



| $V_u$ | $x^{11}$ | $x^{12}$ | $x^{13}$ | $x^{21}$ | $x^{22}$ | $x^{23}$ | $x^{24}$ | $x^{31}$ | $x^{32}$ |
|---|---|---|---|---|---|---|---|---|---|
| $x^{11}$ | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 1 |
| $x^{12}$ | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 1 | 0 |
| $x^{13}$ | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 1 |
| $x^{21}$ | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 |
| $x^{22}$ | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 |
| $x^{23}$ | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 |
| $x^{24}$ | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 |
| $x^{31}$ | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 1 | 0 |
| $x^{32}$ | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 1 |

**Fig. 1.** RNG example with nine groups-modalities - Associated adjacency matrix

This means that if two modalities $x^{kr}$ and $x^{ls}$ which verify the RNG property are connected by an edge, the vertices $x^{kr}$ and $x^{ls}$ are neighbors.

Thus, for any proximity measure given, $u$, we can associate an adjacency matrix $V_u$, of binary and symmetrical order $m$. Figure 1 illustrates an example

of RNG in $\mathbb{R}^2$ of a set of n objects-individuals around nine modalities associated with three qualitative variables $x^1$, $x^2$ and $x^3$ with three, four and two modalities respectively.

For example, for the second modality of the first variable and the first modality of the second variable, $V_u(x^{12}, x^{21}) = 1$, it means that on the geometrical plane, the hyper-Lunula (intersection between the two hyperspheres centered on the two modalities $x^{12}$ and $x^{21}$) is empty.

For a given neighborhood property (MST, GG or RNG), each measure $u$ generates a topological structure on the objects in $E$ which are totally described by the adjacency binary matrix $V_u$. In this paper, we chose to use the Relative Neighbors Graph (GNR).

## 2.1 Comparison and selection of proximity measures

First we compare different proximity measures according to their topological similarity in order to regroup them and to better visualize their resemblances.

To measure the topological equivalence between two proximity measures $u_i$ and $u_j$, we propose to test if the associated adjacency matrices $V_{u_i}$ and $V_{u_j}$ are different or not. The degree of topological equivalence between two proximity measures is measured by the following property of concordance.

*Property 2.* Topological equivalence index between two adjacency matrices

$$S(V_{u_i}, V_{u_j}) = \frac{1}{m^2} \sum_{k=1}^{p} \sum_{r=1}^{m_k} \sum_{l=1}^{p} \sum_{s=1}^{m_l} \delta_{kr\,ls}(x^{kr}, x^{ls})$$

with $\quad \delta_{kr\,ls}(x^{kr}, x^{ls}) = \left\{ \begin{array}{ll} 1 & \text{if } V_{u_i}(x^{kr}, x^{ls}) = V_{u_j}(x^{kr}, x^{ls}) \\ 0 & \text{otherwise.} \end{array} \right.$

Then, in our case, we want to compare these different proximity measures according to their topological equivalence in a context of association. So we define a criterion for measuring the spacing from the independence or no association position.

A contingency table is one of the most common ways to summarize categorical data. Generally, interest lies in whether there is an association between the row variable and the column variable that produce the table; sometimes there is further interest in describing the strength of that association. The data can arise from several different sampling frameworks, and the interpretation of the hypothesis of no association depends on the framework. The question of interest is whether there is an association between the two variables.

We construct the adjacency matrix denoted by $V_{u_*}$, which corresponds best to the Burt table. Thus, to examine similarities between the modalities we examine the gap between each profile-modality and its average profile, that is, the gap to independence. This best adjacency matrix can be written as follows:

*Property 3.* Reference adjacency matrix

$$\left\{ \begin{array}{ll} V_{u_*}(x^{kr}, x^{ls}) = 1 & if \ \frac{\mathcal{B}_{kr\,ls}}{\mathcal{B}_{kr\,..}} \geq \frac{\mathcal{B}_{kr\,..}}{np^2}; \quad \forall k, l = 1, p \ ; \ r = 1, m_k \ and \ s = 1, m_l \\ V_{u_*}(x^{kr}, x^{ls}) = 0 & otherwise \end{array} \right.$$

$\mathcal{B}_{kr\,ls} = \Sigma_{i=1}^n x_i^{kr} x_i^{ls}$, element of the Burt matrix that corresponds to the number of individuals who have the modality r of the variable k and the modality s of the variable l,

$\mathcal{B}_{kr\,..} = \Sigma_{l=1}^p \Sigma_{s=1}^{m_s} b_{kr\,ls}$ is the row margin of the modality r of the variable k,

$\frac{\mathcal{B}_{kr\,ls}}{\mathcal{B}_{kr\,..}}$ is the row profile of the modality r of the variable k,

$\frac{\mathcal{B}_{kr_{..}}}{np^2}$ is the average profile of the modality r of the variable k, $np^2$ being the total number.

The binary and symmetric adjacency matrix $V_{u_*}$ is associated with an unknown proximity measure denoted $u_*$ and called a reference measure.

Thus, with this reference proximity measure we can establish $S(V_{u_i}, V_{u_*})$ the topological equivalence of association between the modalities of the p variables by measuring the percentage of similarity between the adjacency matrix $V_{u_i}$ and the reference adjacency matrix $V_{u_*}$.

In order to graphically describe the similarities between proximity measures, we can for example apply the notion of themascope, [11], which is a methodological sequence of a clustering method on the results of a factorial method. In this case, a Principal Component Analysis (PCA) followed by a Hierarchical Ascendant Classification (HAC) were performed upon the 22 component dissimilarity matrix defined by $[D]_{ij} = D(V_{u_i}, V_{u_j}) = 1 - S(V_{u_i}, V_{u_j})$ to partition them into homogeneous groups and to view their similarities in order to see which measures are close to one another.

We can use any classic visualization techniques to achieve this. For example, we can build a dendrogram of hierarchical clustering of the proximity measures. We can also use multidimensional scaling or any other technique, such as Laplacian projection, to map the 22 proximity measures into a two dimensional space.

Finally, in order to evaluate and determine the closest class of proximity measures to the reference measure $u_*$, we project the latter as a supplementary element into the two data analysis methods, positioned by the dissimilarity vector with 22 components $[D]_{*i} = 1 - S(V_{u_*}, V_{u_i})$.

## 2.2   Statistical comparisons between two proximity measures

In this section, we use Cohen's kappa coefficient [7], to test statistically the degree of topological equivalence between two proximity measures. This non parametric test compares these measures based on their associated adjacency matrices.

The comparison between indices of proximity measures has also been studied by [16], [17] and [8] from a statistical perspective. The authors proposed an approach that compares similarity matrices obtained by each proximity measure, using Mantel's test [13], in a pairwise manner.

Cohen's nonparametric Kappa test is the statistical test best suited to compare matched binary data. The Kendall or Spearman coefficient compares matched continuous data. It makes it possible in this context to measure the agreement or the concordance of the binary values of two adjacency matrices associated with two proximity measures.

Let $V_{u_i}$ and $V_{u_j}$ be adjacency matrices associated with two proximity measures $u_i$ and $u_j$. To compare the degree of topological equivalence between these two measures, we propose to test if the associated adjacency matrices are statistically different or not, using a non-parametric test of paired data. These binary and symmetric matrices of order $m$, are unfolded in two vector-matched components, consisting of $\frac{m(m+1)}{2}$ values: the $m$ diagonal values and the $\frac{m(m-1)}{2}$ values above or below the diagonal.

The degree of topological equivalence between two proximity measures is estimated from the Kappa coefficient, computed on the $2 \times 2$ contingency table formed by the two binary vectors, using the following property:

*Property 4.* Kappa coefficient

$$\widehat{\kappa} = \widehat{\kappa}(V_{u_i}, V_{u_j}) = \frac{P_o - P_e}{1 - P_e}$$

where,

$P_o = \frac{2}{m(m+1)} \sum_{k=0}^{1} n_{kk}$    is the observed proportion of concordance, and

$P_e = \frac{4}{m^2(m+1)^2} \sum_{k=0}^{1} n_{k.} n_{.k}$   represents the expected proportion of concordance under the assumption of independence.

The Kappa coefficient is a real number, without dimension, between $-1$ and $+1$. The concordance is higher the closer the value of Kappa is to 1 and the maximum concordance is reached ($\widehat{\kappa} = 1$) when $P_o = 1$ and $P_e = 0.5$. When there is perfect independence, $\widehat{\kappa} = 0$ with $P_o = P_e$, and in the case of total mismatch, $\widehat{\kappa} = -1$ with $P_o = 0$ and $P_e = 0.5$.

The true value of the Kappa coefficient in the population is a random variable that approximately follows a Gaussian law of mean $E(\kappa)$ and variance $Var(\kappa)$. The null hypothesis $H_0$ is $\kappa = 0$ against the alternative hypothesis $H_1 : \kappa > 0$. We formulate the null hypothesis $H_0 : \kappa = 0$, independence of agreement or concordance. The concordance becomes higher as $\kappa$ tends towards 1, and is a perfect maximum if $\kappa = 1$. It is equal to $-1$ in the case of a perfect discordance.

We also test the topological equivalence between each proximity measure $u_i$ and the perfect measure $u_*$ by comparing the adjacency matrices $V_{u_i}$ and $V_{u_*}$.

## 2.3   Graphical representation of the topological associations

In order to represent graphically the possible topological links between the $m$ modalities of the $p$ qualitative variables, we use Multidimensional Scaling (MDS). It allows to visualize a proximity matrix (similarity or dissimilarity) and makes it possible to pass from a proximity matrix between a set of n objects to the coordinates of these same objects in a p-dimensional space. We propose to carry out the classical MDS, namely factorial analysis on similarity $V_{u*}$ or dissimilarity $D_{u*} = U - V_{u*}$ table [6]. The topological Correspondence Analysis (TMCA) returns to perform the following PCA:

*Property 5.* TMCA consist to perform the PCA of the triple $\{V_{u*} \; ; \; M \; ; \; D_m\}$,

where, $V_{u*}$ is the adjacency matrix associated with the proximity measure $u*$, the most appropriate measure for the considered data, $M = I_m$ is the identity matrix of order $m$ and $D_m = \frac{W}{np}$ is the weighted diagonal matrix of modality weights.

One can also opt for a normalized PCA if one wishes to give the same weight to all the variables in the calculation of the distance between two modalities.

This topological analysis leads to the spectral decomposition of the M-symmetric and positive matrix ${}^t V_{u*} D_m V_{u*} M$, whose explained inertia is equal to $\frac{1}{np} trace({}^t V_{u*} W V_{u*})$, with the first $m - p - 1$ non-zero eigenvalues.
We can thus establish the topological correspondence analysis of each of the 22 proximity measures $u_i$ considered, by carrying out the PCA of the $V_{u_i}$ adjacency data table.

Aid for the interpretation of TMCA results are those of PCA. Graphical representations on factorial plans allow to visualize and identify the topological links between the modalities of the variables. As in weighted PCA, we consider the most significant modalities on the axes, that is the modalities having both a strong relative contribution and a good quality of representation, measured by the square cosine of the angle formed by the point-modality and its projection on the factorial plane considered.

## 3    Application to real data

To illustrate the TMCA, we considered the data displayed in Table 2 of a study on female entrepreneurship conducted in Dakar Senegal in 2014. These data were collected from 153 female entrepreneurs of the Dakar region, and their objective here is to give a topological description of the sample's demographic features: age, marital status, number of children and level of studies.

In a metric and classical context, we simply have to apply an MCA on the homogeneous set of the four characteristics of the female entrepreneurs. The main numerical and graphical results of this MCA, given in Table 9 in the Appendix and in Figure 4, will be compared to those of the proposed TMCA.

In a topological context, the main results of the proposed method are presented in the following tables and graphs, which allow us to visualize proximity measures close to each other and to select the one that best describes the associations between the modalities of the four characteristics of the sample population.

An HAC algorithm based on the Ward criterion [19] was used in order to characterize classes of proximity measure relative to their similarities.
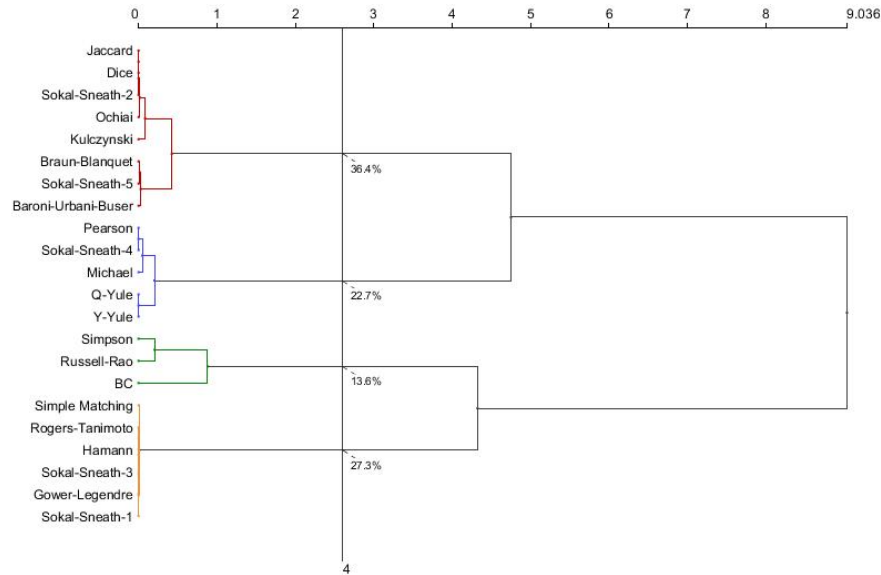
The reference measure $u_*$ is projected as a supplementary element. The dendrogram of Figure 2 represents the hierarchical tree of the 22 proximity measures considered.

---

Aggregation based on the criterion of the loss of minimal inertia.

**Table 2.** Burt table - Female Entrepreneurship in Dakar - Senegal

| | Variables | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Modalities | Age | | | Marital status | | | | Number of children | | | Level of studies | | |
| Under 25 | 22 | 0 | 0 | 18 | 2 | 1 | 1 | 13 | 3 | 6 | 3 | 1 | 18 |
| 25 to 50 years | 0 | 80 | 0 | 16 | 9 | 21 | 34 | 14 | 11 | 55 | 58 | 5 | 17 |
| Over 50 | 0 | 0 | 51 | 3 | 8 | 24 | 16 | 8 | 35 | 8 | 30 | 10 | 11 |
| Single | 18 | 16 | 3 | 37 | 0 | 0 | 0 | 20 | 3 | 14 | 9 | 1 | 27 |
| Divorcee | 2 | 9 | 8 | 0 | 19 | 0 | 0 | 3 | 10 | 6 | 13 | 5 | 1 |
| Monogamous bride | 1 | 21 | 24 | 0 | 0 | 46 | 0 | 7 | 21 | 18 | 26 | 5 | 15 |
| Polygamous bride | 1 | 34 | 16 | 0 | 0 | 0 | 51 | 5 | 15 | 31 | 43 | 5 | 3 |
| No children | 13 | 14 | 8 | 20 | 3 | 7 | 5 | 35 | 0 | 0 | 11 | 5 | 19 |
| From 1 to 3 children | 3 | 11 | 35 | 3 | 10 | 21 | 15 | 0 | 49 | 0 | 27 | 9 | 13 |
| More than 3 children | 6 | 55 | 8 | 14 | 6 | 18 | 31 | 0 | 0 | 69 | 53 | 2 | 14 |
| Illiterate-Primary | 3 | 58 | 30 | 9 | 13 | 26 | 43 | 11 | 27 | 53 | 91 | 0 | 0 |
| Secondary | 1 | 5 | 10 | 1 | 5 | 5 | 5 | 5 | 9 | 2 | 0 | 16 | 0 |
| Superior | 18 | 17 | 11 | 27 | 1 | 15 | 3 | 19 | 13 | 14 | 0 | 0 | 46 |



**Fig. 2.** Hierarchical tree of the proximity measures

Table 3 summarizes the main results of the chosen partition into four homogeneous classes of proximity measure, obtained from the cut of the hierarchical tree of Figure 2.

Moreover, in view of the results in Table 3, the reference measure $u_*$ is closer to the third class consisting of Russell-Rao, Simpson and BC measures for which there is a strong topological association between the modalities of the variables among the 22 proximity measures considered. We will have a stronger

**Table 3.** PCA & HAC results - Assignment of the reference measure

| Class number | Class 1 | Class 2 | Class 3 | Class 4 |
|---|---|---|---|---|
| Frequency | 8 | 5 | 3 | 6 |
| Proximity measure | $u_{Jaccard}$ $u_{Dice}$ $u_{Sokal-Sneath-2}$ $u_{Ochiai}$ $u_{Kulczynski}$ $u_{Baroni-Urbani-Buser}$ $u_{Sokal-Sneath-5}$ $u_{Braun-Blanquet}$ | $u_{Pearson}$ $u_{Sokal-Sneath-4}$ $u_{Q-Yule}$ $u_{Y-Yule}$ $u_{Michael}$ | $u_{Russell-Rao}$ $u_{Simpson}$ $u_{BC}$ | $u_{Simple-Matching}$ $u_{Rogers-Tanimoto}$ $u_{Hamann}$ $u_{Sokal-Sneath-3}$ $u_{Gower-Legendre}$ $u_{Sokal-Sneath-1}$ |
| Reference measure | | | $u_*$ | |

association between the variables of the typical profile of the entrepreneur in Dakar Senegal.

It was shown in [21], by means of a series of experiments, that the choice of proximity measure has an impact on the results of a supervised or unsupervised classification.

For any pair of proximity measures $(u_i; u_j)$ given in Table 7 in the Appendix, we will show how to build and apply the Kappa test in order to compare two adjacency matrices to measure and test their topological equivalence $S(V_{u_i} ; V_{u_j})$.

Let $V_{u_*}$ and $V_{RR}$ the reference and Russell-Rao adjacency matrices, the topological equivalence between the reference and Russell-Rao proximity measures equal $S(V_{u_*}, V_{RR}) = 79.88\%$. These matrices are unfolded to two vectors comprising the $\frac{m(m+1)}{2} = 91$ diagonal and upper-diagonal values. These two binary vectors are two dummy variables represented in the same sample size of 91 pairs of objects. We then formulated the null hypothesis, $H_0 : \kappa = 0$, that the topological equivalence between reference and Russell-Rao proximity measures is not significant according to the considered data.

**Table 4.** Kappa statistic - Reference and Russell-Rao measures

| | $V_{u_{RR}} = 0$ | $V_{u_{RR}} = 1$ | Total |
|---|---|---|---|
| $V_{u*} = 0$ | 51 | 3 | 54 |
| $V_{u*} = 1$ | 14 | 23 | 37 |
| Total | 65 | 26 | 91 |

Table 4 shows the $2 \times 2$ contingency table observed between the two binary vectors associated to the reference and Russell-Rao proximity measures. Thus, for this example, the calculated Kappa value $\hat{\kappa} = 0.5939$ corresponds to a p-value less than 0.01%. Since this probability is lower than a pre-specified significance level of 5%, the null hypothesis that $\kappa = 0$ for these data (no agreement) is rejected.

We can therefore conclude that the topological equivalence between the two proximity measures measured by $S(V_{RR} ; V_{u_*}) = 79.88\%$, is significant.

Table 8 given in the Appendix, summarizes the similarities and Kappa statistic values between all pairs of proximity measures formed with the 22 measures considered and the unknown reference measure $u_*$, in a topological framework. The values below the diagonal correspond to the similarities $S(V_{u_i}, V_{u_j})$ and the values above the diagonal are the Kappa coefficients $\widehat{\kappa}(V_{u_i}, V_{u_j})$. All Kappa statistical tests are significant with $\alpha \leq 5\%$ level of significance.

**Table 5.** Measures with perfect topological equivalence

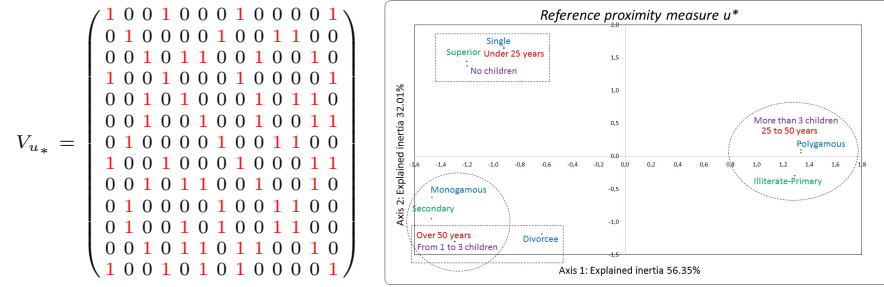| Group 1 | Group 2 | Group 3 |
|---|---|---|
| $u_{Jaccard}$ | $u_{Pearson}$ | $u_{Simple-Matching}$ |
| $u_{Dice}$ | $u_{Sokal-Sneath-4}$ | $u_{Rogers-Tanimoto}$ |
| $u_{Sokal-Sneath-2}$ | | $u_{Hamann}$ |
| | | $u_{Sokal-Sneath-3}$ |
| | | $u_{Gower-Legendre}$ |
| | | $u_{Sokal-Sneath-1}$ |

The similarities in pairs between the 22 proximity measures differ somewhat: some are closer than others. Some measures are in perfect topological equivalence $S(V_{u_i}, V_{u_j}) = 1$ with a perfect concordance $\widehat{\kappa}(V_{u_i}, V_{u_j}) = 1$; these are therefore identical for the data considered, as is the case with the measure groups presented in Table 5.

**Table 6.** Row and Average profiles

| Row-Profiles | Age | | | Marital status | | | | Number of child | | | Level of studies | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Under 25 years | 0.25 | 0 | 0 | 0.205 | 0.023 | 0.011 | 0.011 | 0.148 | 0.034 | 0.068 | 0.034 | 0.011 | 0.205 |
| 25 to 50 years | 0 | 0.25 | 0 | 0.050 | 0.028 | 0.066 | 0.106 | 0.044 | 0.034 | 0.172 | 0.181 | 0.016 | 0.053 |
| Over 50 years | 0 | 0 | 0.25 | 0.015 | 0.039 | 0.118 | 0.078 | 0.039 | 0.172 | 0.039 | 0.147 | 0.049 | 0.054 |
| Single | 0.122 | 0.108 | 0.020 | 0.25 | 0 | 0 | 0 | 0.135 | 0.020 | 0.095 | 0.061 | 0.007 | 0.182 |
| Divorcee | 0.026 | 0.118 | 0.105 | 0 | 0.25 | 0 | 0 | 0.040 | 0.132 | 0.079 | 0.171 | 0.066 | 0.013 |
| Monogamous | 0.005 | 0.114 | 0.130 | 0 | 0 | 0.25 | 0 | 0.038 | 0.114 | 0.098 | 0.141 | 0.027 | 0.082 |
| Polygamous | 0.005 | 0.167 | 0.078 | 0 | 0 | 0 | 0.25 | 0.025 | 0.074 | 0.152 | 0.211 | 0.025 | 0.015 |
| No children | 0.093 | 0.100 | 0.057 | 0.143 | 0.021 | 0.050 | 0.036 | 0.25 | 0 | 0 | 0.079 | 0.036 | 0.136 |
| From 1 to 3 child | 0.015 | 0.056 | 0.179 | 0.015 | 0.051 | 0.107 | 0.077 | 0 | 0.25 | 0 | 0.138 | 0.046 | 0.066 |
| More than 3 child | 0.022 | 0.199 | 0.029 | 0.051 | 0.022 | 0.065 | 0.112 | 0 | 0 | 0.25 | 0.192 | 0.007 | 0.051 |
| Illiterate-Primary | 0.008 | 0.159 | 0.082 | 0.025 | 0.036 | 0.071 | 0.118 | 0.030 | 0.074 | 0.146 | 0.25 | 0 | 0 |
| Secondary | 0.016 | 0.078 | 0.156 | 0.016 | 0.078 | 0.078 | 0.078 | 0.078 | 0.141 | 0.031 | 0 | 0.25 | 0 |
| Superior | 0.098 | 0.092 | 0.060 | 0.147 | 0.005 | 0.082 | 0.016 | 0.103 | 0.071 | 0.076 | 0 | 0 | 0.25 |
| Average profile | 0.036 | 0.131 | 0.083 | 0.061 | 0.031 | 0.075 | 0.083 | 0.057 | 0.080 | 0.113 | 0.149 | 0.026 | 0.075 |

The adjacency matrix $V_{u*}$ associated to the best adapted proximity measure $u*$ to the considered data is established from the profile table 6. Figure 3 shows on the main first TMCA plan, the significant links between the modalities of the signage of female entrepreneurship. The links are materialized by geometric shapes.

Figure 4 presents, for comparison, on the first factorial plan, the graphical representations of the multiple correspondence analyses, the topological (TMCA) and the classical (MCA) [9], [5].

$$V_{u_*} = \begin{pmatrix} 1 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 1 & 0 & 0 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 0 & 1 & 1 & 0 & 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 & 1 & 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 1 & 0 & 1 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$



**Fig. 3.** TMCA - Adjacency matrix and Graphical representation

Unlike the MCA method which describes only three strong links, the TMCA highlights four: two opposing on the first factorial axis (56.35%) and the other two on the second factorial axis (32.01%).



**Fig. 4.** Comparison TMCA & MCA

Considering percentages of inertia presented in Appendix Table 9 which represent the associations between all modalities, we restrict the comparison of the graphical representations to the two first factorial axes.

We can represent Figure 5, the different TMCA of 4 of the 22 proximity measures considered. One can moreover give Figure 6, the graphical representation associated with a perfect topological independence.

**Fig. 5.** Jaccard (Class 1), Pearson (Class 2), Russel & Rao (Class 3) and Simple Matching (Class 4) proximity measures



**Fig. 6.** Adjacency identity matrix and Perfect independence representation

## 4 Conclusion

This paper proposes a TMCA which is a new topological method of multiple correspondence analysis that enriches the classical methods of qualitative data analysis. This work compares existing proximity measures to perform a topological multiple correspondence analysis based on the notion of neighborhood graphs according to considered data. Future work involves extending this topological approach to other factorial methods of data analysis, especially to analyze the correlation structure of a set of continuous variables - Topological Principal Component Analysis (TPCA), or to synthesize the relations existing between two groups of continuous variables - Topological Canonical Analysis (TCA).

# 5   Appendix

**Table 7.** Some proximity measures.

| Measures | Similarity | Dissimilarity |
|---|---|---|
| Jaccard | $s_1 = \frac{a}{a+b+c}$ | $u_1 = 1 - s_1$ |
| Dice, Czekanowski | $s_2 = \frac{2a}{2a+b+c}$ | $u_2 = 1 - s_2$ |
| Kulczynski | $s_3 = \frac{1}{2}\left(\frac{a}{a+b} + \frac{a}{a+c}\right)$ | $u_3 = 1 - s_3$ |
| Driver, Kroeber and Ochiai | $s_4 = \frac{a}{\sqrt{(a+b)(a+c)}}$ | $u_4 = 1 - s_4$ |
| Sokal and Sneath 2 | $s_5 = \frac{a}{a+2(b+c)}$ | $u_5 = 1 - s_5$ |
| Braun-Blanquet | $s_6 = \frac{a}{max(a+b,a+c)}$ | $u_6 = 1 - s_6$ |
| Simpson | $s_7 = \frac{a}{min(a+b,a+c)}$ | $u_7 = 1 - s_7$ |
| Kendall, Sokal-Michener | $s_8 = \frac{a+d}{a+b+c+d}$ | $u_8 = 1 - s_8$ |
| Russell and Rao | $s_9 = \frac{a}{a+b+c+d}$ | $u_9 = 1 - s_9$ |
| Rogers and Tanimoto | $s_{10} = \frac{a+d}{a+2(b+c)+d}$ | $u_{10} = 1 - s_{10}$ |
| Pearson $\phi$ | $s_{11} = \frac{ad-bc}{\sqrt{(a+b)(a+c)(d+b)(d+c)}}$ | $u_{11} = \frac{1-s_{11}}{2}$ |
| Hamann | $s_{12} = \frac{a+d-b-c}{a+b+c+d}$ | $u_{12} = \frac{1-s_{12}}{2}$ |
| bc | | $u_{13} = \frac{4bc}{(a+b+c+d)^2}$ |
| Sokal and Sneath 5 | $s_{14} = \frac{ad}{\sqrt{(a+b)(a+c)(d+b)(d+c)}}$ | $u_{14} = 1 - s_{14}$ |
| Michael | $s_{15} = \frac{4(ad-bc)}{(a+d)^2+(b+c)^2}$ | $u_{15} = \frac{1-s_{15}}{2}$ |
| Baroni, Urbani and Buser | $s_{16} = \frac{a+\sqrt{ad}}{a+b+c+\sqrt{ad}}$ | $u_{16} = 1 - s_{16}$ |
| Yule Q | $s_{17} = \frac{ad-bc}{ad+bc}$ | $u_{17} = \frac{1-s_{17}}{2}$ |
| Yule Y | $s_{18} = \frac{\sqrt{ad}-\sqrt{bc}}{\sqrt{ad}+\sqrt{bc}}$ | $u_{18} = \frac{1-s_{18}}{2}$ |
| Sokal and Sneath 4 | $s_{19} = \frac{1}{4}\left(\frac{a}{a+b} + \frac{a}{a+c} + \frac{d}{d+b} + \frac{d}{d+c}\right)$ | $u_{19} = 1 - s_{19}$ |
| Sokal and Sneath 3 | | $u_{20} = \frac{b+c}{a+d}$ |
| Gower and Legendre | $s_{21} = \frac{a+d}{a+\frac{(b+c)}{2}+d}$ | $u_{21} = 1 - s_{21}$ |
| Sokal and Sneath 1 | $s_{22} = \frac{2(a+d)}{2(a+d)+b+c}$ | $u_{22} = 1 - s_{22}$ |

**Table 8.** Similarities $S(V_{u_i}, V_{u_j})$ & Kappa coefficient $\widehat{\kappa}(V_{u_i}, V_{u_j})$

| Measure | Sokal-Sneath-1 | Gower-Legendre | Sokal-Sneath-3 | Sokal-Sneath-4 | Y-Yule | Q-Yule | Baroni-Urbani-Buser | Michael | Sokal-Sneath-5 | BC | Hamann | Pearson | Rogers-Tanimoto | Russell-Rao | Simple Matching | Simpson | Braun-Blanquet | Sokal-Sneath-2 | Ochiai | Kulczynski | Dice | Jaccard |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Reference measure $u_*$ | 0.67 | 0.67 | 0.67 | 0.67 | 0.67 | 0.67 | 0.67 | 0.67 | 0.67 | 0.62 | 0.67 | 0.67 | 0.67 | 0.59 | 0.67 | 0.57 | 0.64 | 0.69 | 0.67 | 0.67 | 0.69 | 0.69 |
| Jaccard | 0.71 | 0.71 | 0.71 | 0.77 | 0.72 | 0.72 | 0.92 | 0.77 | 0.87 | 0.56 | 0.71 | 0.77 | 0.71 | 0.74 | 0.72 | 0.66 | 0.90 | 1 | 0.97 | 0.92 | 0.96 | 1 |
| Dice | 0.71 | 0.71 | 0.71 | 0.77 | 0.72 | 0.72 | 0.92 | 0.77 | 0.87 | 0.56 | 0.71 | 0.77 | 0.71 | 0.74 | 0.72 | 0.66 | 0.90 | 1 | 0.97 | 0.92 | 1 | 0.96 |
| Kulczynski | 0.68 | 0.68 | 0.68 | 0.74 | 0.74 | 0.74 | 0.85 | 0.77 | 0.84 | 0.58 | 0.68 | 0.79 | 0.68 | 0.81 | 0.68 | 0.74 | 0.82 | 0.92 | 0.99 | 1 | 0.96 | 0.96 |
| Ochiai | 0.71 | 0.71 | 0.71 | 0.79 | 0.72 | 0.72 | 0.90 | 0.79 | 0.89 | 0.53 | 0.68 | 0.77 | 0.68 | 0.76 | 0.68 | 0.68 | 0.87 | 0.97 | 1 | 0.99 | 0.99 | 0.99 |
| Sokal-Sneath-2 | 0.71 | 0.71 | 0.71 | 0.82 | 0.72 | 0.72 | 0.97 | 0.82 | 0.87 | 0.56 | 0.71 | 0.82 | 0.71 | 0.74 | 0.71 | 0.56 | 0.90 | 1 | 0.94 | 0.92 | 0.95 | 0.95 |
| Braun-Blanquet | 0.63 | 0.63 | 0.63 | 0.64 | 0.69 | 0.69 | 0.59 | 0.64 | 0.68 | 0.74 | 0.63 | 0.64 | 0.63 | 0.76 | 0.63 | 0.56 | 1 | 0.90 | 0.90 | 0.90 | 0.90 | 0.90 |
| Simpson | 0.63 | 0.63 | 0.63 | 0.66 | 0.69 | 0.69 | 0.59 | 0.64 | 0.70 | 0.65 | 0.60 | 0.66 | 0.60 | 0.60 | 0.71 | 1 | 0.80 | 0.56 | 0.66 | 0.66 | 0.66 | 0.66 |
| Simple Matching | 0.71 | 0.71 | 0.71 | 0.79 | 0.74 | 0.74 | 0.85 | 0.74 | 0.74 | 0.79 | 1 | 0.79 | 1 | 0.82 | 1 | 0.83 | 0.87 | 0.71 | 0.71 | 0.68 | 0.71 | 0.72 |
| Russell-Rao | 0.79 | 0.79 | 0.79 | 0.69 | 0.74 | 0.74 | 0.69 | 0.74 | 0.74 | 0.79 | 0.79 | 0.79 | 0.91 | 1 | 0.82 | 0.83 | 0.87 | 0.74 | 0.76 | 0.86 | 0.80 | 0.74 |
| Rogers-Tanimoto | 0.91 | 0.91 | 0.91 | 0.85 | 0.74 | 0.74 | 0.69 | 0.74 | 0.74 | 0.69 | 0.79 | 1 | 1 | 0.60 | 0.71 | 0.63 | 0.63 | 0.71 | 0.68 | 0.66 | 0.71 | 0.71 |
| Pearson | 0.95 | 0.95 | 0.79 | 0.85 | 0.74 | 0.74 | 0.85 | 0.95 | 0.85 | 0.69 | 0.79 | 1 | 0.79 | 0.79 | 0.71 | 0.66 | 0.64 | 0.74 | 0.77 | 0.77 | 0.77 | 0.77 |
| Hamann | 0.79 | 0.79 | 0.79 | 0.79 | 0.74 | 0.74 | 0.69 | 0.74 | 0.74 | 0.79 | 1 | 0.91 | 1 | 0.85 | 0.88 | 0.83 | 0.80 | 0.71 | 0.66 | 0.60 | 0.60 | 0.63 |
| BC | 0.79 | 0.69 | 0.69 | 0.54 | 0.74 | 0.74 | 0.90 | 0.74 | 0.58 | 1 | 0.91 | 0.93 | 0.88 | 0.87 | 0.86 | 0.86 | 0.80 | 0.80 | 0.79 | 0.81 | 0.80 | 0.80 |
| Sokal-Sneath-5 | 0.85 | 0.85 | 0.85 | 0.85 | 0.74 | 0.79 | 0.85 | 0.95 | 1 | 0.81 | 0.88 | 0.98 | 0.88 | 0.85 | 0.88 | 0.86 | 0.94 | 0.87 | 0.95 | 0.93 | 0.94 | 0.94 |
| Michael | 0.79 | 0.79 | 0.79 | 0.69 | 0.74 | 0.79 | 0.90 | 1 | 0.93 | 0.83 | 0.88 | 0.98 | 0.88 | 0.87 | 0.88 | 0.83 | 0.92 | 0.96 | 0.91 | 0.91 | 0.89 | 0.89 |
| Baroni-Urbani-Buser | 0.92 | 0.92 | 0.85 | 0.90 | 0.79 | 0.79 | 1 | 0.93 | 0.95 | 0.79 | 0.86 | 0.93 | 0.86 | 0.87 | 0.86 | 0.81 | 0.99 | 0.96 | 0.95 | 0.93 | 0.96 | 0.96 |
| Q-Yule | 0.72 | 0.79 | 0.74 | 0.72 | 0.95 | 1 | 0.91 | 0.95 | 0.91 | 0.88 | 0.88 | 0.98 | 0.88 | 0.85 | 0.88 | 0.86 | 0.89 | 0.87 | 0.95 | 0.88 | 0.87 | 0.87 |
| Y-Yule | 0.72 | 0.79 | 0.74 | 0.72 | 1 | 0.98 | 0.91 | 0.95 | 0.91 | 0.88 | 0.88 | 0.98 | 0.88 | 0.85 | 0.88 | 0.86 | 0.89 | 0.87 | 0.95 | 0.88 | 0.87 | 0.87 |
| Sokal-Sneath-4 | 0.79 | 0.79 | 0.79 | 1 | 0.98 | 0.98 | 0.93 | 0.98 | 0.93 | 0.86 | 0.91 | 1 | 0.91 | 0.85 | 0.91 | 0.83 | 0.92 | 0.89 | 0.88 | 0.91 | 0.89 | 0.89 |
| Sokal-Sneath-3 | 1 | 1 | 1 | 0.74 | 0.88 | 0.88 | 0.86 | 0.88 | 0.88 | 0.91 | 1 | 0.91 | 1 | 0.82 | 0.87 | 0.83 | 0.87 | 0.87 | 0.86 | 0.86 | 0.87 | 0.87 |
| Gower-Legendre | 1 | 1 | 1 | 0.79 | 0.88 | 0.88 | 0.86 | 0.88 | 0.88 | 0.91 | 1 | 0.91 | 1 | 0.82 | 0.87 | 0.83 | 0.87 | 0.87 | 0.86 | 0.86 | 0.87 | 0.87 |
| Sokal-Sneath-1 | 1 | 1 | 1 | 0.79 | 0.88 | 0.88 | 0.86 | 0.88 | 0.88 | 0.91 | 1 | 0.91 | 1 | 0.82 | 0.87 | 0.83 | 0.87 | 0.87 | 0.86 | 0.86 | 0.87 | 0.87 |
| Reference measure $u_*$ | 0.83 | 0.83 | 0.83 | 0.83 | 0.83 | 0.83 | 0.83 | 0.83 | 0.83 | 0.81 | 0.83 | 0.83 | 0.83 | 0.80 | 0.83 | 0.79 | 0.82 | 0.85 | 0.83 | 0.83 | 0.85 | 0.85 |
| Measure | Sokal-Sneath-1 | Gower-Legendre | Sokal-Sneath-3 | Sokal-Sneath-4 | Y-Yule | Q-Yule | Baroni-Urbani-Buser | Michael | Sokal-Sneath-5 | BC | Hamann | Pearson | Rogers-Tanimoto | Russell-Rao | Simple Matching | Simpson | Braun-Blanquet | Sokal-Sneath-2 | Ochiai | Kulczynski | Dice | Jaccard |

Examples:

$$S(u_{Kulczynski}\ ,\ u_{Jaccard}) = 0.96$$
$$\widehat{\kappa}(u_{Jaccard}\ ,\ u_{Kulczynski}) = 0.92\ ;\ \ p-value < 0.01\%$$

All Kappa statistical tests are significant with $\alpha \leq 5\%$ level of Significance.

**Table 9.** Eigenvalues associated with the topological and classical multiple correspondence analyses

| TMCA | Axis | Eigenvalue | Proportion | Cumulative |
|---|---|---|---|---|
| | 1 | 1.609 | 56.35% | 56.35% |
| | 2 | 0.914 | 32.01% | 88.36% |
| | 3 | 0.159 | 5.56% | 93.91% |
| | 4 | 0.087 | 3.06% | 96.97% |
| | 5 | 0.032 | 1.12% | 98.10% |
| | 6 | 0.027 | 0.95% | 99.05% |
| | 7 | 0.015 | 0.53% | 99.59% |
| $m - p - 1 \rightarrow$ | 8 | 0.012 | 0.41% | 100.00% |
| | Total | 2.855 | 100.00% | 100.00% |

| MCA | Axis | Eigenvalue | Proportion | Cumulative |
|---|---|---|---|---|
| | 1 | 0.585 | 26.01% | 26.01% |
| | 2 | 0.462 | 20.52% | 46.53% |
| | 3 | 0.285 | 12.67% | 59.20% |
| | 4 | 0.222 | 9.85% | 69.05% |
| | 5 | 0.212 | 9.40% | 78.45% |
| | 6 | 0.166 | 7.39% | 85.84% |
| | 7 | 0.126 | 5.60% | 91.44% |
| | 8 | 0.101 | 4.48% | 95.92% |
| $m - p \rightarrow$ | 9 | 0.092 | 4.08% | 100.00% |
| | Total | 2.250 | 100.00% | 100.00% |

# References

1. Abdesselam, R.: Selection of proximity measures for a Topological Correspondence Analysis. *In a Book Series*, 5th Stochastic Modeling Techniques and Data Analysis, International Conference, Chania, Greece, C.H. Skiadas (Ed), 11–24 (2018).

2. Abdesselam, R.: A Topological Discriminant Analysis. *In book Chapter, Volume 3, Data Analysis and Applications 2: Utilization of Results in Europe and 0ther Topics*, J. Bozeman and C. Skiadas Editors, ISTE Science Publishing, Wiley, 167–178 (2018).

3. Batagelj, V., Bren, M.: Comparing resemblance measures. In Proc. International Meeting on Distance Analysis (Distancia'92) (1992).

4. Batagelj, V., Bren, M.: Comparing resemblance measures. *In Journal of classification*, 12, 73–90 (1995).

5. Benzcri, J.P.: L'Analyse des Données. Tome 1 : La Taxinomie. Tome 2 : L'analyse des correspondances, *"2ème édition Dunod, Paris* (1976).

6. Caillez, F. and Pagès, J.P.: Introduction à l'Analyse des données. *S.M.A.S.H., Paris*, 1976.

7. Cohen, J.: A coefficient of agreement for nominal scales. *Educ Psychol Meas*, Vol 20, 27–46 (1960).

8. Demsar, J.: Statistical comparisons of classifiers over multiple data sets. *The journal of Machine Learning Research*, Vol. 7, 1–30 (2006).

9. Escofier, B.: Une représentation des variables dans l'analyse des correspondances multiples. *Revue de statistique Appliquées*, 27, 37–47 (1979).

10. Kim, J.H. and Lee, S.: Tail bound for the minimal spanning tree of a complete graph. *In Statistics & Probability Letters*, 4, 64, 425–430 (2003).

11. Lebart, L.: Stratégies du traitement des données d'enquêtes. *La Revue de MODULAD, 3, 21–29 (1989).*

12. *Lesot, M. J., Rifqi, M. and Benhadda,H.: Similarity measures for binary and numerical data: a survey.* In IJKESDP, 1, 1, 63-84 (2009).

13. *Mantel, N.: A technique of disease clustering and a generalized regression approach.* In Cancer Research, 27, 209–220 (1967).

14. *Park, J. C.,Shin, H. and Choi, B. K.: Elliptic Gabriel graph for finding neighbors in a point set and its application to normal vector estimation.* In Computer-Aided Design Elsevier, 38, 6, 619–626 (2006).

15. *Rifqi, M., Detyniecki, M. and Bouchon-Meunier, B.: Discrimination power of measures of resemblance.* IFSA'03 Citeseer (2003).

16. *Schneider, J. W. and Borlund, P.: Matrix comparison, Part 1: Motivation and important issues for measuring the resemblance between proximity measures or ordination results.* In Journal of the American Society for Information Science and Technology, 58, 11, 1586–1595 (2007).

17. *Schneider, J. W. and Borlund, P.: Matrix comparison, Part 2: Measuring the resemblance between proximity measures or ordination results by use of the Mantel and Procrustes statistics.* In Journal of the American Society for Information Science and Technology, 11, 58, 1596–1609 (2007).

18. *Toussaint, G. T.: The relative neighbourhood graph of a finite planar set.* In Pattern recognition, 12, 4, 261–268 (1980).

19. *Ward, J. R.: Hierarchical grouping to optimize an objective function.* In Journal of the American statistical association JSTOR, 58, 301, 236–244 (1963).

20. *Warrens, M. J.: Bounds of resemblance measures for binary (presence/absence) variables. In Journal of Classification,* Springer, 25, 2, 195–208 (2008).

21. *Zighed, D., Abdesselam, R., and Hadgu, A.: Topological comparisons of proximity measures.* In the 16th PAKDD 2012 Conference. *In P.-N. Tan et al., Eds. Part I, LNAI 7301, Springer-Verlag Berlin Heidelberg, 379–391 (2012).*