

CORRESPONDENCE ANALYSIS

Hervé Abdi & Michel Béra

Title: Correspondence Analysis

Name: Hervé Abdi¹, Michel Béra²

Affil./Addr. 1: School of Behavioral and Brain Sciences
The University of Texas at Dallas
Richardson, TX 75080, USA
herve@utdallas.edu

Affil./Addr. 2: Centre d'Étude et de Recherche en Informatique et Communica-
tions
Conservatoire National des Arts et Métiers
F-75141 Paris Cdex 03, France
michel.bera@cnam.fr

Correspondence Analysis

Synonyms

Dual Scaling, optimal scaling, homogeneity analysis.

Glossary

CA: Correspondence analysis

component: A linear combination of the variables of a data table

dimension: see component

factor: see component

GSVD: Generalized singular value decomposition

PCA: Principal component analysis

SVD: Singular value decomposition

Abdi, H., & Béra, M. (to appear 2014). Correspondence Analysis. In R. Alhajj and J. Rokne (Eds.), "Encyclopedia of Social Networks and Mining." New York: Springer Verlag.

Introduction

Correspondence analysis [CA; see 11; 22; 21; 13; 16; 17; 7; 20; 1; 9] is an extension of principal component analysis (PCA, for details, see [8]) tailored to handle nominal variables. Originally, CA was developed to analyze contingency tables in which a sample of observations is described by two nominal variables, but it was rapidly extended to the analysis of any data matrices with non-negative entries. The origin of CA can be traced to early work of Pearson ([24]) or Fisher, but the modern version of correspondence analysis and its geometric interpretation comes from the 1960s in France and is associated with the French school of “data analysis” (*analyse des donn ées*) and was developed under the leadership of Jean-Paul Benzécri. As a technique, it was often discovered (and re-discovered) and so variations of CA can be found under several different names such as “dual-scaling,” “optimal scaling,” “homogeneity analysis,” or “reciprocal averaging.” The multiple identities of correspondence analysis are a consequence of its large number of properties: Correspondence analysis can be defined as an optimal solution for a lot of apparently different problems.

Key Points

CA transforms a data table into two sets of new variables called *factor scores* (obtained as linear combinations of, respectively the rows and columns): One set for the rows and one set for the columns. These factor scores give the best representation of the similarity structure of, respectively, the rows and the columns of the table. In addition, the factors scores can be plotted as maps, that optimally display the information in the original table. In these maps, rows and columns are represented as points whose coordinates are the factor scores and where the dimensions are also called *factors*, *components* (by analogy with PCA), or simply *dimensions*. Interestingly, the factor scores of the

rows and the columns have the same variance and, therefore, rows and columns can be conveniently represented in one single map.

In correspondence analysis, the total variance (often called inertia) of the factor scores is proportional to the independence chi-square statistic of this table and therefore the factor scores in CA decompose this χ^2 into orthogonal components.

Correspondence Analysis: Theory and Practice

Notations

Matrices are denoted in upper case bold letters, vectors are denoted in lower case bold, and their elements are denoted in lower case italic. Matrices, vectors, and elements from the same matrix all use the same letter (*e.g.*, \mathbf{A} , \mathbf{a} , a). The transpose operation is denoted by the superscript \top , the inverse operation is denoted by $^{-1}$. The identity matrix is denoted \mathbf{I} , vectors or matrices of ones are denoted $\mathbf{1}$, matrices or vectors of zeros are denoted $\mathbf{0}$. When provided with a square matrix, the **diag** operator gives a vector with the diagonal elements of this matrix. When provided with a vector, the **diag** operator gives a diagonal matrix with the elements of the vector as the diagonal elements of this matrix. When provided with a square matrix, the **trace** operator gives the sum of the diagonal elements of this matrix.

The data table to be analyzed by CA is a contingency table (or at least a data table with non-negative entries) with I rows and J columns. It is represented by the $I \times J$ matrix \mathbf{X} , whose generic element $x_{i,j}$ gives the number of observations that belong to the i th level of the first nominal variables (*i.e.*, the rows) *and* the j th level of the second nominal variable (*i.e.*, the columns). The grand total of the table is noted N .

Computations

The first step of the analysis is to transform the data matrix into a probability matrix (*i.e.*, a matrix comprising non-negative numbers and whose sum is equal to one) denoted \mathbf{Z} and computed as $\mathbf{Z} = N^{-1}\mathbf{X}$. We denote \mathbf{r} the vector of the row totals of \mathbf{Z} , (*i.e.*, $\mathbf{r} = \mathbf{Z}\mathbf{1}$, with $\mathbf{1}$ being a conformable vector of 1's), \mathbf{c} the vector of the columns totals (*i.e.*, $\mathbf{c} = \mathbf{Z}^T\mathbf{1}$), and $\mathbf{D}_c = \text{diag}\{\mathbf{c}\}$, $\mathbf{D}_r = \text{diag}\{\mathbf{r}\}$. The factor scores are obtained from the following *generalized* singular value decomposition (GSVD, for details on the singular value decomposition see [2; 3; 15; 27; 18; 12; 26]):

$$(\mathbf{Z} - \mathbf{r}\mathbf{c}^T) = \mathbf{P}\mathbf{\Delta}\mathbf{Q}^T \quad \text{with} \quad \mathbf{P}^T\mathbf{D}_r^{-1}\mathbf{P} = \mathbf{Q}^T\mathbf{D}_c^{-1}\mathbf{Q} = \mathbf{I}. \quad (1)$$

Note that the subtraction of the matrix $\mathbf{r}\mathbf{c}^T$ from \mathbf{Z} is equivalent to a double centering of the matrix ([5; 6]). The matrix \mathbf{P} (respectively \mathbf{Q}) contains the left (respectively right) generalized singular vectors of $(\mathbf{Z} - \mathbf{r}\mathbf{c}^T)$, and the diagonal elements of the diagonal matrix $\mathbf{\Delta}$ give its *singular values*. The squared singular values, which are called *eigenvalues*, are denoted λ_ℓ and stored into the diagonal matrix $\mathbf{\Lambda}$. Eigenvalues express the variance extracted by the corresponding factor and their sum is called the total inertia (denoted \mathcal{I}) of the data matrix. With the so called “triplet notation,” ([14]) that is sometimes used as a general framework to formalize multivariate techniques, CA is equivalent to the analysis of the triplet $((\mathbf{Z} - \mathbf{r}\mathbf{c}^T), \mathbf{D}_c^{-1}, \mathbf{D}_r^{-1})$.

From the GSVD, the row and (respectively) column factor scores are obtained as

$$\mathbf{F} = \mathbf{D}_r^{-1}\mathbf{P}\mathbf{\Delta} \quad \text{and} \quad \mathbf{G} = \mathbf{D}_c^{-1}\mathbf{Q}\mathbf{\Delta}. \quad (2)$$

Note that the factor scores of a given set (*i.e.*, the rows or the columns) are pairwise orthogonal when they describe different dimensions and that the variance of the factor scores for a given dimension is equal to the eigenvalue associated with this dimension. So, for example, the variance of the row factor scores is computed as:

$$\mathbf{F}^T \mathbf{D}_r \mathbf{F} = \Delta \mathbf{P}^T \mathbf{D}_r^{-1} \mathbf{D}_r \mathbf{D}_r^{-1} \mathbf{P} \Delta = \Delta \mathbf{P}^T \mathbf{D}_r^{-1} \mathbf{P} \Delta = \Delta^2 = \Lambda . \quad (3)$$

What does correspondence analysis optimize?

In CA the criterion that is maximized is the variance of the factor scores (see [21; 16]). For example, the row first factor \mathbf{f}_1 is obtained as a linear combination of the columns of the matrix $(\mathbf{Z} - \mathbf{r}\mathbf{c}^T)$ taking into account the constraints imposed by the matrices \mathbf{D}_r^{-1} and \mathbf{D}_c^{-1} . Specifically, this means that we are searching for the vector \mathbf{q}_1 containing the weights of the linear combination such as \mathbf{f}_1 is obtained as

$$\mathbf{f}_1 = \mathbf{D}_r^{-1} (\mathbf{Z} - \mathbf{r}\mathbf{c}^T) \mathbf{D}_c^{-1} \mathbf{q}_1 , \quad (4)$$

such that

$$\mathbf{f}_1 = \arg \max_{\mathbf{f}} \mathbf{f}^T \mathbf{D}_r \mathbf{f} , \quad (5)$$

under the constraint that

$$\mathbf{q}_1^T \mathbf{D}_c^{-1} \mathbf{q}_1 = 1 . , \quad (6)$$

The subsequent row factor scores will maximize the residual variance under the orthogonality constraint imposed by the matrix \mathbf{D}_r^{-1} (*i.e.*, $\mathbf{f}_2^T \mathbf{D}_r^{-1} \mathbf{f}_1 = 0$).

How to identify the elements important for a factor

In CA, the rows and the columns of the table have a similar role (and variance) and therefore we can use the same statistics to identify the rows and the columns important for a given dimension. Because the variance extracted by a factor (*i.e.*, its eigenvalue) is obtained as the weighted sum of the factor scores for this factor of either the rows or columns of the table, the importance of a row (respectively a column) is reflected by the ratio of its squared factor score to the eigenvalue of this factor. This ratio is called the *contribution* of the row (respectively column) to the factor. Specifically, the contributions of row i to component ℓ and of column j to component ℓ are obtained respectively as:

$$\text{ctr}_{i,\ell} = \frac{r_i f_{i,\ell}^2}{\lambda_\ell} \quad \text{and} \quad \text{ctr}_{j,\ell} = \frac{c_j g_{j,\ell}^2}{\lambda_\ell} \quad (7)$$

(with r_i being the i th element of \mathbf{r} and c_j being the j th element of \mathbf{c}). Contributions take values between 0 and 1 and their sums for a given factor is equal to one for either the rows or the columns. A convenient rule of thumb is to consider that contributions larger than the average (*i.e.*, $\frac{1}{I}$ for the rows and $\frac{1}{J}$ for the columns) are important for a given factor.

How to identify the important factors for an element

The factors important for a given row or column are identified by computing statistics called *squared cosines*. These statistics are obtained by decomposing the squared distance of an element along the factors of the analysis. Specifically, the vector of the squared (χ^2) distance from the rows and columns to their respective barycenter (*i.e.*, average or center of gravity) are obtained as

$$\mathbf{d}_r = \text{diag} \{ \mathbf{F} \mathbf{F}^\top \} \quad \text{and} \quad \mathbf{d}_c = \text{diag} \{ \mathbf{G} \mathbf{G}^\top \} . \quad (8)$$

Recall that the total inertia (\mathcal{I}) in CA is equal to the sum of the eigenvalues, and, that in CA, this inertia can also be computed as the weighted sum of the squared distances of the rows *or* the columns to their respective barycenter. Formally, the inertia can be computed as:

$$\mathcal{I} = \sum_l^L \lambda_\ell = \mathbf{r}^\top \mathbf{d}_r = \mathbf{c}^\top \mathbf{d}_c . \quad (9)$$

The squared *cosine* between row i and component ℓ and column j and component ℓ are obtained respectively as:

$$\cos_{i,\ell}^2 = \frac{f_{i,\ell}^2}{d_{r,i}^2} \quad \text{and} \quad \cos_{j,\ell}^2 = \frac{g_{j,\ell}^2}{d_{c,j}^2} . \quad (10)$$

(with $d_{r,i}^2$, and $d_{c,j}^2$, being respectively the i -th element of \mathbf{d}_r and the j -th element of \mathbf{d}_c). The sum of the squared cosines over the dimensions for a given element is equal

to one and so the cosine can be seen as the proportion of the variance of an element that can be attributed to a given dimension.

Correspondence Analysis and the Chi-Square test

CA is intimately related to the independence χ^2 test. Recall that the (null) hypothesis stating that the rows and the columns of a contingency table \mathbf{X} are independent can be tested by computing the following χ^2 criterion:

$$\chi^2 = N \sum_i^I \sum_j^J \frac{(z_{i,j} - r_i c_j)^2}{r_i c_j} = N \phi^2, \quad (11)$$

where ϕ^2 is called the mean square contingency coefficient (for a 2×2 table, it is called the squared coefficient of correlation associated to the χ^2 test, in this special case it takes values between 0 and 1). For a contingency table, under the null hypothesis, the χ^2 criterion follows a χ^2 distribution with $(I - 1)(J - 1)$ degrees of freedom and can, therefore, be used (under the usual assumptions) to test the independence hypothesis.

The total inertia of the matrix $(\mathbf{Z} - \mathbf{r}\mathbf{c}^\top)$ (under the constraints imposed on the SVD by the matrices \mathbf{D}_r^{-1} and \mathbf{D}_c^{-1}), can be computed as the sum of the eigenvalues of the analysis or directly from the data matrix as

$$\mathcal{I} = \text{trace} \left\{ \mathbf{D}_c^{-\frac{1}{2}} (\mathbf{Z} - \mathbf{r}\mathbf{c}^\top)^\top \mathbf{D}_c^{-1} (\mathbf{Z} - \mathbf{r}\mathbf{c}^\top) \mathbf{D}_c^{-\frac{1}{2}} \right\} = \sum_i^I \sum_j^J \frac{(z_{i,j} - r_i c_j)^2}{r_i c_j} = \phi^2. \quad (12)$$

This shows that the total inertia is proportional to the independence χ^2 (specifically $\mathcal{I} = N^{-1}\chi^2$) and therefore that the factors of CA perform an orthogonal decomposition of the independence χ^2 where each factor “explains” a portion of the deviation to independence.

The transition formula

In CA, rows and columns play a symmetric role and their factor scores have the same variance. As a consequence of this symmetry, the row factor scores (respectively the

column factor scores) can be derived from the column factor scores (respectively the row factor scores). This can be seen by rewriting Equation 2 taking account Equation 1. For example, the factor scores for the rows can be computed as

$$\begin{aligned}
\mathbf{F} &= \mathbf{D}_r^{-1} \mathbf{P} \mathbf{\Delta} \\
&= \mathbf{D}_r^{-1} (\mathbf{Z} - \mathbf{r} \mathbf{c}^T) \mathbf{D}_c^{-1} \mathbf{Q} \\
&= \mathbf{D}_r^{-1} (\mathbf{Z} - \mathbf{r} \mathbf{c}^T) \mathbf{G} \mathbf{\Delta}^{-1} \\
&= \mathbf{D}_r^{-1} \mathbf{Z} \mathbf{G} \mathbf{\Delta}^{-1} - \mathbf{D}_r^{-1} \mathbf{r} \mathbf{c}^T \mathbf{G} \mathbf{\Delta}^{-1} .
\end{aligned} \tag{13}$$

Because the matrix $(\mathbf{Z} - \mathbf{r} \mathbf{c}^T)$ contains deviations to its row and column barycenters the rows and columns sum are equal to zero and therefore the matrix $\mathbf{D}_r^{-1} \mathbf{r} \mathbf{c}^T \mathbf{G} \mathbf{\Delta}^{-1}$ is equal to $\mathbf{0}$ and therefore Equation 13 can be rewritten as

$$\mathbf{F} = \mathbf{D}_r^{-1} \mathbf{Z} \mathbf{G} \mathbf{\Delta}^{-1} . \tag{14}$$

So, if we denoted by \mathbf{R} the matrix $\mathbf{R} = \mathbf{D}_r^{-1} \mathbf{Z}$ in which each row (whose sum is one) is called a *row profile* and \mathbf{C} the matrix $\mathbf{C} = \mathbf{D}_c^{-1} \mathbf{Z}^T$ in which each column (whose sum is one) is called a *column profile*, the transition formulas are:

$$\mathbf{F} = \mathbf{R} \mathbf{G} \mathbf{\Delta}^{-1} \text{ and } \mathbf{G} = \mathbf{C} \mathbf{F} \mathbf{\Delta}^{-1} . \tag{15}$$

This shows that the factor scores of an element of one set (*e.g.*, a row) are computed as the barycenter of the expression of this element in the other set (*e.g.*, the columns) followed by an expansion (as expressed by the $\mathbf{\Delta}^{-1}$ term) that is inversely proportional to the singular value of each factor.

In ca, the singular values are never larger than one

Note that, together, the two transition formulas from Equation 15 imply that the diagonal terms of $\mathbf{\Delta}^{-1}$ are larger than one (because otherwise the range of each set of factor scores would be smaller than the other one which, in turn, would imply that all

these factor scores are null), and therefore that the singular values in CA are always equal to or smaller than one.

Distributional equivalence

An important property of correspondence analysis is to give the same results when two rows (respectively two columns) that are proportional are merged together. This property, called *distributional equivalence* (or also “distributional equivalency”), also applies approximately: the analysis is only changed a little when two rows (or columns) that are almost proportional are merged together.

How to interpret point proximity

In a CA map when two row (respectively column) points are close to each other this means that these points have similar profiles and when two points have the same profile they will be located exactly at the same place (this is a consequence of the distributional equivalence principle). The proximity between row and column points is more delicate to interpret because of the barycentric principle (see section on the transition formula): the position of a row (respectively column) point is determined from its barycenter on the column (respectively row), and therefore the proximity between a row point and one column point cannot be interpreted directly.

Asymmetric plot: how to interpret row and column proximity

CA treats rows and columns symmetrically and so their rôles are equivalent. In some cases however rows and columns can play different rôles and this symmetry can be misleading. As an illustration, in the example used below (see section “Example”), the participants were asked to choose the color that would match a given piece of music. In this framework, the colors can be considered as a dependent variable and the pieces

of music as independent variable. In this case the rôles are asymmetric and the plots can reflect this asymmetry by normalizing one set such that the variance of its factor scores is equal to 1 for each factor. For example, the normalized to one column factor scores, denoted $\tilde{\mathbf{G}}$ would be computed as (compare with Equation 2)

$$\tilde{\mathbf{G}} = \mathbf{D}_{\mathbf{c}}^{-1}\mathbf{Q} . \quad (16)$$

In the asymmetric plot obtained with \mathbf{F} and $\tilde{\mathbf{G}}$ the distances between rows and columns can now be interpreted meaningfully: the distance from a row point to a column point reflects their association (and a row is positioned exactly at the barycenter of the columns).

Centered vs non-centered analysis

CA is obtained from the generalized singular value decomposition of the centered matrix $(\mathbf{Z} - \mathbf{r}\mathbf{c}^T)$, but it could be also obtained from the same singular value decomposition of matrix \mathbf{Z} . In this case, the first pair of singular vectors is equal to \mathbf{r} and \mathbf{c} and their associated singular value is equal to 1. This property is easily verified from the following relations

$$\mathbf{c}\mathbf{D}_{\mathbf{c}}^{-1}\mathbf{Z} = \mathbf{1}\mathbf{Z} = \mathbf{r} \text{ and } \mathbf{r}\mathbf{D}_{\mathbf{r}}^{-1}\mathbf{Z}^T = \mathbf{1}\mathbf{Z}^T = \mathbf{c} . \quad (17)$$

Because in CA, the singular values are never larger than one, \mathbf{r} and \mathbf{c} having a singular value of 1, are the first pair of singular vectors of \mathbf{Z} . Therefore, the generalized singular value of \mathbf{Z} can be developed as:

$$\mathbf{Z} = \mathbf{r}\mathbf{c}^T + (\mathbf{Z} - \mathbf{r}\mathbf{c}^T) = \mathbf{r}\mathbf{c}^T + \mathbf{P}\mathbf{\Delta}\mathbf{Q}^T . \quad (18)$$

This shows that the ℓ th pair of singular vectors and singular value of $(\mathbf{Z} - \mathbf{r}\mathbf{c}^T)$ are the $(\ell + 1)$ th pair of singular vectors and singular value of \mathbf{Z} .

Supplementary elements

Often in CA we want to know the position in the analysis of rows or columns that were not actually analyzed. These rows or columns are called *illustrative*, *supplementary*, or *out of sample* rows or columns (or supplementary observations or variables). By contrast with the appellation of supplementary (which are not used to compute the factors) the *active* elements are used to compute the factors.

The projection formula, used the *transition* formula (see Equation 15) and is specific to correspondence analysis. Specifically, let $\mathbf{i}_{\text{sup}}^{\text{T}}$ being an illustrative row and \mathbf{j}_{sup} being an illustrative column to be projected (note that in CA, prior to projection, a illustrative row or column is re-scaled such that its sum is equal to one). Their coordinates of the illustrative rows (denoted \mathbf{f}_{sup}) and column (denoted \mathbf{g}_{sup}) are obtained as:

$$\mathbf{f}_{\text{sup}} = (\mathbf{i}_{\text{sup}}^{\text{T}} \mathbf{1})^{-1} \mathbf{i}_{\text{sup}}^{\text{T}} \mathbf{G} \tilde{\mathbf{\Delta}}^{-1} \quad \text{and} \quad \mathbf{g}_{\text{sup}} = (\mathbf{j}_{\text{sup}}^{\text{T}} \mathbf{1})^{-1} \mathbf{j}_{\text{sup}}^{\text{T}} \mathbf{F} \tilde{\mathbf{\Delta}}^{-1}. \quad (19)$$

[note that the scalar terms $(\mathbf{i}_{\text{sup}}^{\text{T}} \mathbf{1})^{-1}$ and $(\mathbf{j}_{\text{sup}}^{\text{T}} \mathbf{1})^{-1}$ are used to ensure that the sum of the elements of \mathbf{i}_{sup} or \mathbf{j}_{sup} is equal to one, if this is already the case, these terms are superfluous].

Programs and packages

CA is implemented in most statistical packages (*e.g.*, SAS, XLSTAT, SYSTAT) with R giving the most comprehensive implementation. Several packages in R are specifically dedicated to correspondence analysis and its variants. The most popular are the packages `ca`, `FactoMineR` ([19]), `ade4`, and `ExPosition` (this last package was used to analyze the example presented below). MATLAB programs are also available for download from www.utdallas.edu/~herve.

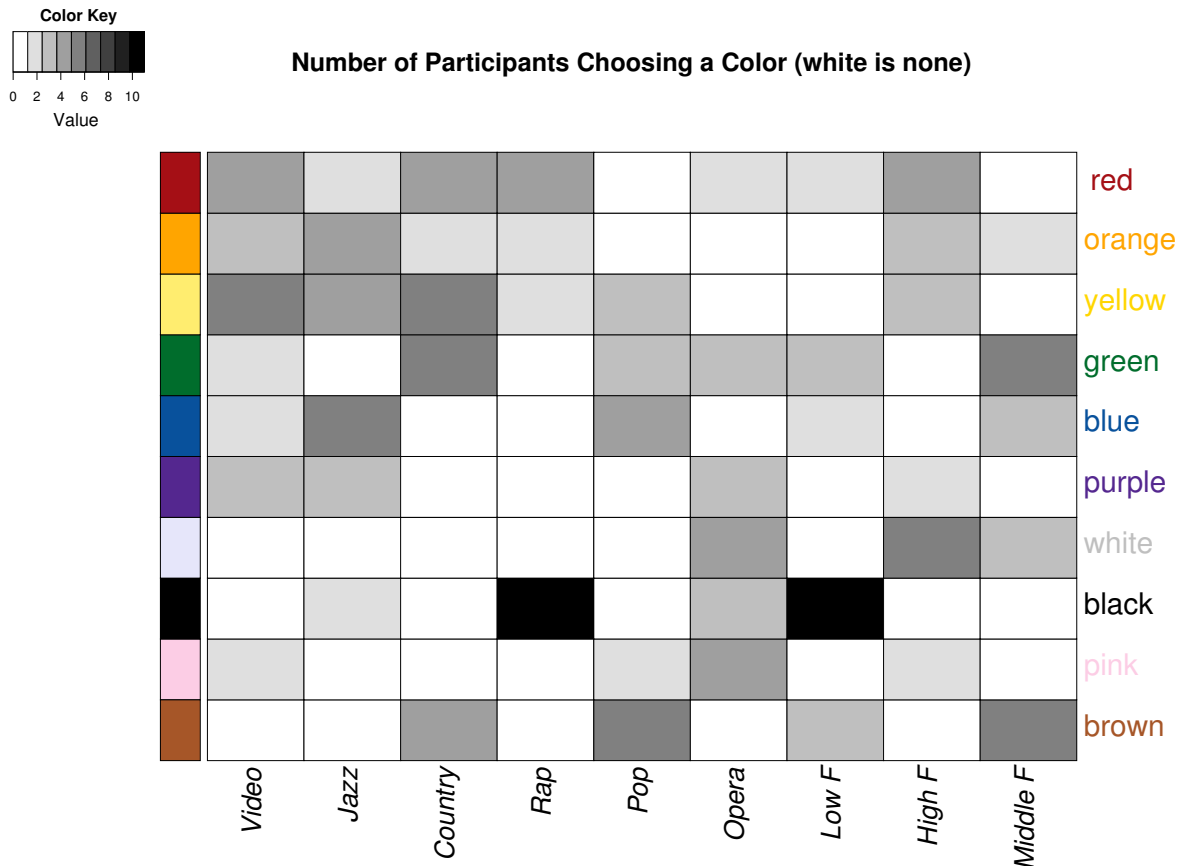


Fig. 1. CA The Colors of Music. A heat map of the data from Table 1

Example: The Colors of Sounds

To illustrate CA we have collected data (as part of a science project!) from twenty-two participants who were presented with nine “pieces of music” and asked to associate one of ten colors to each piece of music. The pieces of music were 1) the music of a video game (video), 2) a Jazz song (jazz) 3) a country and western song (country), 4) a rap song (rap), 5) a pop song (pop) 6) an extract of the opera Carmen (opera), 7) the low F note played on a piano (low F), 8) the middle F note played on the same piano, and, finally, 9) the high F note still played on the same piano. The data are shown in Table 1 where the columns are the pieces of music, the rows are the colors, and the numbers at the intersection of the rows and the columns give the number of times the color in the row was associated to the piece of music in the column. A graphics representation of the data from Table 1 is given by the “heat map” displayed in Figure 1.

Table 2. CA The Color of Music. Factor scores, contributions, mass, mass \times squared factor scores, inertia to barycenter, and squared cosines for the rows. For convenience, squared cosines and contributions have been multiplied by 1000 and rounded.

	F_1	F_2	ctr ₁	ctr ₂	r_i	F_1^2	F_2^2	$d_{r,i}^2$	\cos_1^2	\cos_2^2
red	-0.026	0.299	0	56	.121	.000	.011	.026	3	410
orange	-0.314	0.232	31	25	.091	.009	.005	.030	295	161
yellow	-0.348	0.202	53	27	.126	.015	.005	.057	267	89
green	-0.044	-0.490	1	144	.116	.000	.028	.048	5	583
blue	-0.082	-0.206	2	21	.096	.001	.004	.050	13	81
purple	-0.619	0.475	87	77	.066	.025	.015	.050	505	298
white	-0.328	0.057	26	1	.071	.008	.000	.099	77	2
black	1.195	0.315	726	75	.146	.208	.014	.224	929	65
pink	-0.570	0.300	68	28	.061	.020	.005	.053	371	103
brown	0.113	-0.997	5	545	.106	.001	.105	.108	12	973
Σ	—	—	1000	1000	—	.287	.192	.746		
						λ_1	λ_2	\mathcal{I}		
						39%	26%			
						τ_1	τ_2			

A CA of Table 1 extracted eight components. We will, here, only consider the first two components which together account for 65% of the total inertia (with eigenvalues of .287 and .192, respectively). The factor scores of the observations (rows) and variables (columns) are shown in Tables 2 and 3 respectively. The corresponding map is displayed in Figure 2.

We can see from Figures 2 (symmetric plot) and 3 (asymmetric plot) that the first component isolates the color black from all the other colors and that black is mostly associated with two pieces of music: Rap and the low F. The association with

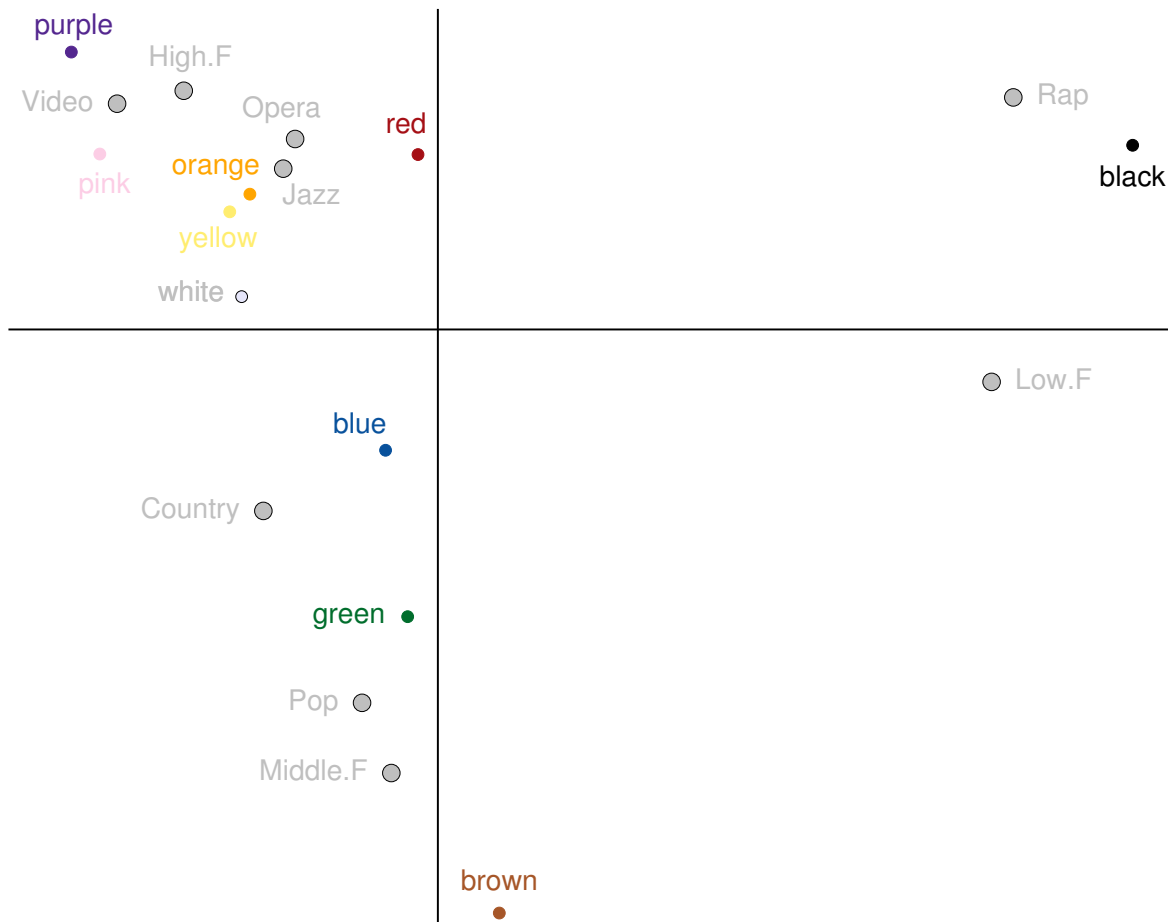


Fig. 2. CA The Colors of Music. Symmetric Plot: The projections of the rows and the columns are displayed in the same map. $\lambda_1 = .287$, $\tau_1 = 39$; $\lambda_2 = .192$, $\tau_2 = 26$. In this plot the proximity between rows and columns cannot be directly interpreted.

the low note reflects a standard association between pitch and color (high notes are perceived as bright and low notes as dark); by contrast, the association of rap music with the black color is likely to reflect a semantic association. The squared cosines show that the first component accounts for most of the variance of the black color (93%, see Table 2). The second component separates the colors brown and (to a lesser extent) green from the other colors (in particular purple) and that brown and green as associated with Pop and Middle F. On the other side of the second dimension we find the color purple and a quartet of pieces of music (Video, High F, Opera, and Jazz).

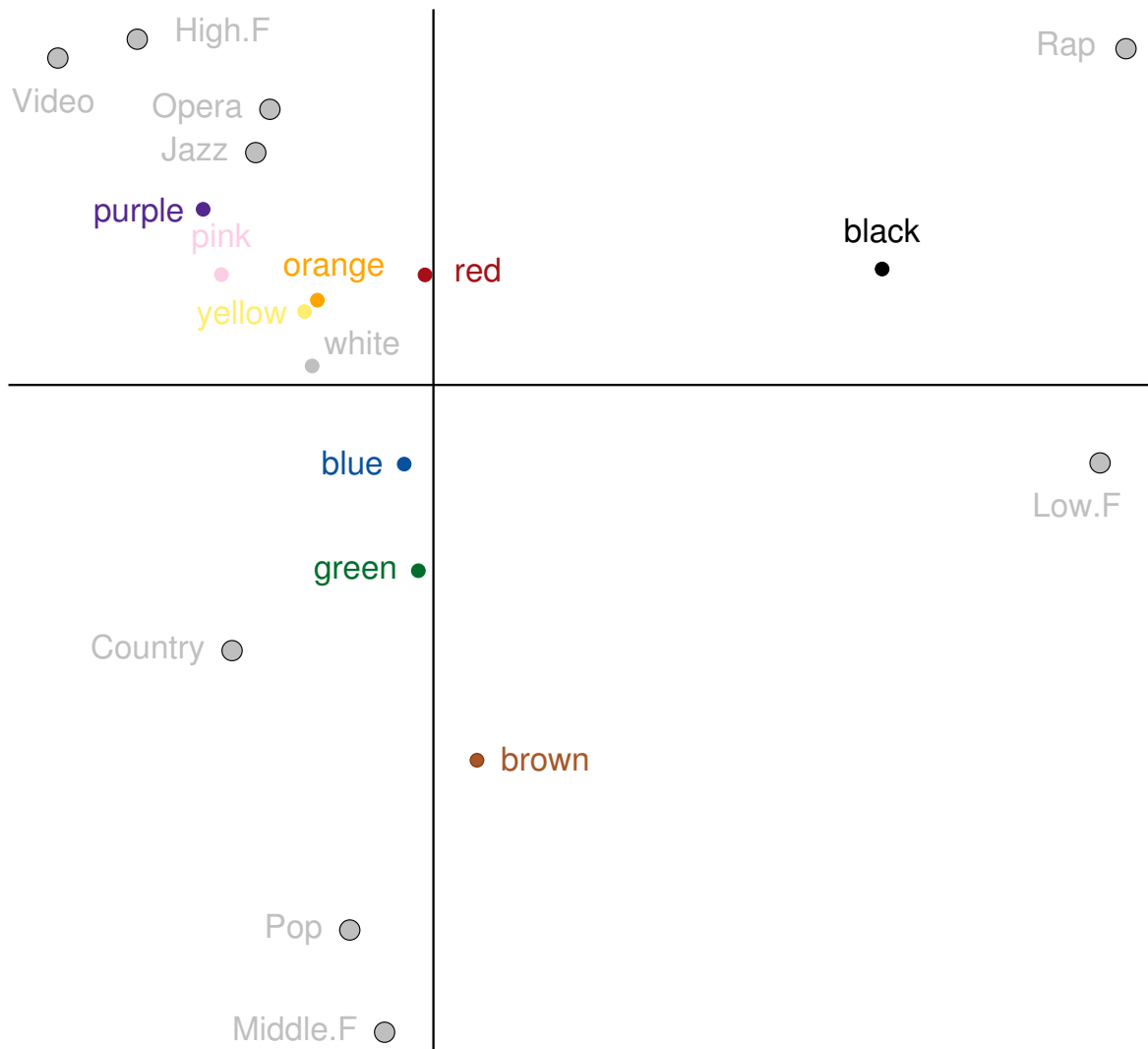


Fig. 3. CA The Colors of Music. Asymmetric Plot: The projections of the rows and the columns are displayed in the same map. The inertia of the projections of the column factor scores is equal to one for each dimension and the inertia of the projections of the row factor scores are $\lambda_1 = .287$, $\tau_1 = 39$; $\lambda_2 = .192$, $\tau_2 = 26$. In this plot the proximity between rows and columns can be directly interpreted.

Conclusion

CA is a very versatile and popular technique in multivariate analysis. In addition to the basics presented here CA includes numerous variants such as multiple correspondence analysis (to analyze several nominal variables, see [7; 21; 17]) discriminant correspondence analysis [to assign observation to a-priori defined groups, see *e.g.*, 23; 25; 4], multi-block correspondence analysis (when the variables are structured into blocks, see [7; 21; 17; 10; 28]).

Table 3. CA The Colors of Music. Factor scores, contributions, mass, mass \times squared factor scores, inertia to barycenter, and squared cosines for the columns. For convenience, squared cosines and contributions have been multiplied by 1000 and rounded.

	G_1	G_2	\tilde{G}_1	\tilde{G}_2	ctr ₁	ctr ₂	c_j	G_1^2	G_2^2	$d_{c,j}^2$	\cos_1^2	\cos_2^2
Video	-0.541	0.386	-1.007	0.879	113	86	.111	.032	.017	.071	454	232
Jazz	-0.257	0.275	-0.478	0.626	25	44	.111	.007	.008	.069	105	121
Country	-0.291	-0.309	-0.541	-0.704	33	55	.111	.009	.011	.066	142	161
Rap	0.991	0.397	1.846	0.903	379	91	.111	.109	.017	.133	822	132
Pop	-0.122	-0.637	-0.227	-1.450	6	234	.111	.002	.045	.064	26	709
Opera	-0.236	0.326	-0.440	0.742	22	61	.111	.006	.012	.079	78	149
Low.F	0.954	-0.089	1.777	-0.203	351	5	.111	.101	.001	.105	962	8
High.F	-0.427	0.408	-0.795	0.929	70	96	.111	.020	.018	.074	271	249
Middle.F	-0.072	-0.757	-0.134	-1.723	2	330	.111	.001	.064	.084	7	759
Σ	—	—	—	—	1000	1000	—	.287	.192	.746		
								λ_1	λ_2	\mathcal{I}		
								39%	26%			
								τ_1	τ_2			

Cross-References

Bibliometrics, Network Analysis, and Knowledge Generation, Barycentric Discriminant Analysis, Clustering Algorithms, Co-Inertia Analysis, Data Mining, Distance and Similarity Measures, Eigenvalues, Singular Value Decomposition, Machine Learning for Social Networks, Matrix Algebra, Basics of Matrix Decomposition, Network Analysis in French Sociology and Anthropology, Network Models, Principal Component Analysis, Probability Matrices, Similarity Metrics on Social Networks,

References

1. Abdi, H. (2003). Multivariate analysis. In M. Lewis-Beck, A. Bryman, and T. Futing (Eds), *Encyclopedia for research methods for the social sciences* pp. 699–702 . Thousand Oaks, CA: Sage.
2. Abdi, H. (2007a). Singular Value Decomposition (SVD) and Generalized Singular Value Decomposition(GSVD). In N.J. Salkind (Ed), *Encyclopedia of measurement and statistics* (pp. 907–912). Thousand Oaks: Sage.
3. Abdi, H. (2007b). Eigen-decomposition: eigenvalues and eigenvectors. In N.J. Salkind (Ed), *Encyclopedia of measurement and statistics* (pp. 304–308). Thousand Oaks: Sage.
4. Abdi, H. (2007d). Discriminant correspondence analysis. In N.J. Salkind (Ed), *Encyclopedia of measurement and statistics* (pp. 270–275). Thousand Oaks: Sage.
5. Abdi, H. (2007e). Metric multidimensional scaling. In N.J. Salkind (Ed), *Encyclopedia of measurement and statistics* (pp. 598–605). Thousand Oaks: Sage.
6. Abdi, H. (2007f). Z-scores. In N.J. Salkind (Ed), *Encyclopedia of measurement and statistics* (pp. 1057–1058). Thousand Oaks: Sage.
7. Abdi, H., and Valentin, D. (2007a). Multiple correspondence analysis. In N.J. Salkind (Ed), *Encyclopedia of measurement and statistics* (pp. 651–657). Thousand Oaks, CA: Sage.
8. Abdi, H. and Williams L.J. (2010a). Principal component analysis. *Wiley Interdisciplinary Reviews: Computational Statistics*, 2, 433–459.
9. Abdi, H., and Williams, L.J. (2010b). Correspondence analysis. In N.J. Salkind (Ed), *Encyclopedia of research design*. Thousand Oaks: Sage.
10. Abdi, H., Williams, L.J., and Valentin, D. (2013). Multiple factor analysis: Principal component analysis for multi-table and multi-block data sets. *Wiley Interdisciplinary Reviews: Computational Statistics*, 5, 149–179.
11. Benzécri, J.-P. (1973). *L'analyse des données, Vols. 1 and 2*. Paris: Dunod.
12. Eckart, C., and Young, G. (1936). The approximation of a matrix by another of a lower rank. *Psychometrika*, 1, 211–218.
13. Escofier, B., and Pagès, J. (1990). *Analyses factorielles simples et multiples: objectifs, méthodes, interprétation*. Dunod: Paris.

14. Escoufier, Y. (2007). Operators related to a data matrix: a survey. In COMPSTAT: Proceedings in Computational Statistics; 17th Symposium Held in Rome, Italy, 2006. New York: Physica Verlag, (pp. 285–297).
15. Good, I., (1969). Some applications of the singular value decomposition of a matrix. *Technometrics*, 11, 823–831.
16. Greenacre, M.J. (1984). *Theory and applications of correspondence analysis*. London: Academic Press.
17. Greenacre, M.J. (2007). *Correspondence analysis in practice (2nd Edition)*. Boca Raton (FL): Chapman & Hall/CRC.
18. Hotelling, H. (1933). Analysis of a complex of statistical variables into principal components. *Journal of Educational Psychology*, 25, 417–441.
19. Husson, F., Lê, and Pagès (2011) *Exploratory Multivariate Analysis by Example Using R*. Boca Raton (FL): Chapman & Hall/CRC.
20. Hwang, H., Tomiuk, M. A., and Takane, Y. (2010). Correspondence analysis, multiple correspondence analysis and recent developments. In R. Millsap and A. Maydeu-Olivares (Eds.). *Handbook of quantitative methods in psychology*. London: Sage Publications.
21. L. Lebart, A. Morineau, and K. M. Warwick (1984). *Multivariate Descriptive Statistical Analysis: Correspondence Analysis and Related Techniques for Large Matrices*. London: Wiley.
22. Lebart, L., and Fénélon, J.P. (1975). *Statistique et informatique appliquées*. Paris: Dunod.
23. Nakache, J.P., Lorente, P., Benzécri, J.P., and Chastang, J.F. (1977). Aspect pronostics et thérapeutiques de l'infarctus myocardique aigu. *Les Cahiers de l'Analyse des Données*, 2, 415–534.
24. Pearson, K. (1901). On lines and planes of closest fit to systems of points in space. *Philosophical Magazine*, 6, 559–572.
25. Saporta, G, and Niang, N. (2006). Correspondence analysis and classification. In M. Greenacre and J. Blasius (Eds), *Multiple correspondence analysis and related methods*. (pp. 371–392). Boca Raton: Chapman & Hall.
26. Stewart, G.W. (1993). On the early history of the singular value decomposition. *SIAM Review*, 35, 551–566.
27. Takane, Y. (2002). Relationships among various kinds of eigenvalue and singular value decompositions. In Yanai, H., Okada, A., Shigemasu, K., Kano, Y., and Meulman, J. (Eds.), *New developments in psychometrics* (pp. 45–56). Tokyo: Springer Verlag.

28. Williams, L.J., Abdi, H., French, R., and Orange, J.B. (2010). A tutorial on Multi-Block Discriminant Correspondence Analysis (MUDICA): A new method for analyzing discourse data from clinical populations. *Journal of Speech Language and Hearing Research*, 53, 1372–1393.