



**HAL**  
open science

## The biproportional factorial analysis

L. de Mesnard

► **To cite this version:**

L. de Mesnard. The biproportional factorial analysis. [Research Report] Laboratoire d'analyse et de techniques économiques(LATEC). 1993, 10 p., figures, bibliographie. hal-01545703

**HAL Id: hal-01545703**

**<https://hal.science/hal-01545703>**

Submitted on 23 Jun 2017

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

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

# LATEC

## LABORATOIRE D'ANALYSE ET DE TECHNIQUES ÉCONOMIQUES

UMR 5601 CNRS

DOCUMENT DE TRAVAIL



Pôle d'Économie et de Gestion

2, boulevard Gabriel - 21000 DIJON - Tél. 03 80 39 54 30 - Fax 03 80 39 54 43

ISSN : 1260-8556

n° 9306

**THE BIPROPORTIONAL FACTORIAL  
ANALYSIS**

**L. de MESNARD\***

**September 1993**

*Professeur à l'Université de Bourgogne, Faculté de Science Economique et de Gestion  
LATEC (URA 342 CNRS)*

# THE BI-PROPORTIONAL FACTORIAL ANALYSIS

L. de Mesnard\*

May 1993

---

**Abstract.** Beyond the Factorial Analysis of Correspondences, the paper presents a new method of data analysis: the Biproportional Factorial Analysis. In the Factorial Analysis of Correspondences, the matrix to be diagonalised is the product of the two matrices of profiles, row and columns: this matrix is not symmetrical. In the Biproportional Factorial Analysis, the matrix to be diagonalised is the symmetrical product of an intermediate matrix over itself; this intermediate matrix is calculated as the biproportion of the data matrix over normalised margins. This provides a full symmetry between rows and columns. After recalling the Factorial Analysis of Correspondences, the paper recall what it is biproportion and then presents the Biproportional Factorial Analysis and discuss it.

**Keywords.** Biproportion, RAS, Factorial Analysis.

## 1. RECALL: THE FACTORIAL ANALYSIS OF CORRESPONDENCES

In the factorial analysis methods, like the Factorial Analysis of Correspondences, the data table  $N$  is a contingency table, with dimensions  $(n,m)$ . In the space of  $\mathcal{R}^m$  (when points are rows), to eliminate the effect of the size of the margins, the terms  $N_{ij}$  of the rows of this matrix are divided by the margins of the rows  $N_{i.}$  giving  $X_{Rij} = \frac{N_{ij}}{N_{i.}}$ , that is to say<sup>1</sup>, if  $N_R$  is the diagonal matrix of the terms  $N_{i.}$ ,  $X_R = N_R^{-1} N$  (this matrix is called the matrix of row profiles). The row margins of  $X_R$  are unitary. The distance between two points is calculated using the  $\chi^2$  metric, multiplying the squares of the differences  $\frac{N_{ij}}{N_{i.}} - \frac{N_{i'j}}{N_{i'}}$  by  $\frac{N}{N_{j.}}$ : if  $N_C$  is the diagonal matrix of the terms  $N_{j.}$ , the metric is  $M_R = N_{..} N_C^{-1}$ . Data are weighted by  $\frac{N_{i.}}{N_{..}}$ , that

---

\* Professor at University of Dijon, Faculty of Economics, 4 Boulevard Gabriel, 21000 DIJON, France.

<sup>1</sup> We use a notation system identical to those of G. Saporta [ SAPORTA 1990, pp. 203-205 ].

the terms  $N_j$ , the metric is  $M_R = N_{..} N_C^{-1}$ . Data are weighted by  $\frac{N_i}{N_{..}}$ , that is to say  $P_R = \frac{N_R}{N_{..}}$ .

Thus, the matrix to be diagonalised to obtain the principal factors  $v$  is  $M_R X_R' P_R X_R = N_C^{-1} N' N_R^{-1} N$ . To obtain the principal components  $c$  (the projection of the  $n$  individuals upon the  $m$  principal axis), we calculate  $c = X_R v$ .

At this point, the method is similar to the simple Principal Components Analysis Method. However, it is not necessary to have centred data<sup>2</sup>. Suppose that data are not centred. The vector of coordinates of the gravity center of the set of points  $X_R = N_R^{-1} N$ , weighted by

$$P_R = \frac{N_R}{N_{..}}, \text{ is } g_R = X_R' P_R \mathbf{1} = \frac{1}{N_{..}} N_R^{-1} N N_R \mathbf{1} = \begin{pmatrix} N_1/N_{..} \\ N_2/N_{..} \\ \vdots \\ N_m/N_{..} \end{pmatrix}, \text{ where } \mathbf{1} \text{ is the sum vector: } \mathbf{1} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}.$$

The variance matrix of the points is  $V = X_R' P_R X_R - g g'$ . Then  $M_R X_R' P_R X_R = M_R V + M_R g g'$

$$\Rightarrow g' M_R X_R' P_R X_R = g' M_R V + g' M_R g g'$$

As the vector  $Og$  is orthogonal to the set of points, it is a principal axis:  $g$  is eigenvector of the matrix  $M_R V$ , with the eigenvalue  $\lambda = 0$ :  $M_R V g = 0$ . Then  $g' M_R X_R' P_R X_R = 0 + \|g\|_{X^2}^2 g' = g'$ . The diagonalisation of  $M_R X_R' P_R X_R = N_C^{-1} N' N_R^{-1} N$  provides a first trivial eigenvalue equal to 1.

The same dual calculus can be made in the space of  $\mathcal{R}^n$  (when points are columns), the terms  $N_{ij}$  of the columns of this matrix are divided by the margins of the columns  $N_j$  giving

$$X_{Cij} = \frac{N_{ij}}{N_j}, \text{ or } X_C = N N_C^{-1} \text{ (this matrix is called the matrix of column profiles). The column}$$

margins of  $X_C$  are unitary. The metric is  $M_C = N_{..} N_R^{-1}$ . The distances are weighted by

$$P_C = \frac{N_C}{N_{..}}. \text{ To found the principal factors, the matrix to be diagonalised is}$$

$$M_C X_C P_C X_C' = N_R^{-1} N N_C^{-1} N'. \text{ Eigenvalues are the same than in the preceding case.}$$

Note that Factorial Analysis of Correspondences looks like a particular case of the general Principal Components Analysis (with metric not unitary and with weights). Also, the matrix to be diagonalised is the product of the two matrices of profiles.

## 2. THE BIPROPORTIONAL FACTORIAL ANALYSIS

The starting problem of Correspondence Analysis appears to give unitary margins to the data table  $N$ , alternatively to rows and to columns. The result depends on that fact: the correspondence between these double result is studied. However, when the row margins are equal (unitary), they are not equal (unitary) in columns, and reciprocally. As there is duality

<sup>2</sup> We follow here [ SAPORTA pp. 200-204 ].

between these double results, in fact they brings the same information seen on a different point. The information obtained when row and column margins are simultaneously equal could be interesting and it eliminates simultaneously the effect of the size of the margins of rows and columns. Then it looks to be interesting to give the same margins to rows and the same margins to columns simultaneously, that is to say to give such margins to  $N$  using a biproportion and then to make a "Correspondence-like Analysis".

## 2.1. Recall: biproportion

A biproportion of a matrix  $N$  on the margins of a matrix  $L$  is the matrix  $X = A N B$ .

Note that margins respects a global equilibrium:  $\sum_{i=1}^n l_i = \sum_{j=1}^m l_j$ .

Different algorithms may be used to calculate  $A$  and  $B$ . For example:

$$a_i = \frac{l_i}{\sum_{j=1}^m b_j n_{ij}}, \forall i \quad \text{and} \quad b_j = \frac{l_j}{\sum_{i=1}^n a_i n_{ij}}, \forall j$$

$$\text{or, } \begin{cases} A = \text{Diag}(A^*) \\ A^* = [\text{Diag}(N B^*)]^{-1} L_R \mathbf{1} \end{cases} \quad \text{and} \quad \begin{cases} B = \text{Diag}(B^*) \\ B^* = [\text{Diag}(N' A^*)]^{-1} L_C \mathbf{1} \end{cases}$$

where  $\text{Diag}$  is the matrix operation who transforms a vector into a diagonal matrix<sup>3</sup>,  $L_R$  is the diagonal matrix of the terms  $l_i$  and  $L_C$  is the diagonal matrix of the terms  $l_j$ .

This algorithm is convergent, the solution of is unique. Moreover, we proved that, whatever be the algorithm, the result is the same: biproportion is so unique than proportion is [de MESNARD].

Remark: to apply a biproportion, data must not be negative, what excludes centred data: it is not a problem because diagonalisation may be done over not centred data, as shown.

## 2.2. The new method

1) First step. There are two cases.

a) The data matrix is square  $(n,n)$ . We calculate a biproportion of the matrix  $N$  on any matrix  $L$  with unitary margins, that is to say, we calculate a matrix  $X = A N B$ , where,

$$\begin{cases} A = \text{Diag}(A^*) \\ A^* = [\text{Diag}(N B^*)]^{-1} \mathbf{1} \end{cases} \quad \text{and} \quad \begin{cases} B = \text{Diag}(B^*) \\ B^* = [\text{Diag}(N' A^*)]^{-1} \mathbf{1} \end{cases}$$

---

<sup>3</sup> This compact and easy-to-read writing is not efficient in computing: one must use a special algorithm to inverse diagonal matrices without using the standard matrix inversion routine.

b) The data matrix is rectangular  $(n,m)$ . We calculate a biproportion of the matrix  $N$  on any matrix  $L$  with unitary margins in one side and with equal margins in the other side. Suppose

that  $m < n$  (without loss of generality); then  $L$  is a constant matrix:  $L_{n,m} = \begin{bmatrix} 1/m & \dots & 1/m \\ \vdots & & \vdots \\ 1/m & \dots & 1/m \end{bmatrix}$  and

$L_R = \mathbf{1}$  and  $L_C = \begin{pmatrix} n/m \\ \vdots \\ n/m \end{pmatrix} = \frac{n}{m} \mathbf{1}$ . We are obliged to do so to respect the global equilibrium

$\sum_{i=1}^n l_i = \sum_{j=1}^m l_j$ . However, the important is that each rows have the same margins and each columns have the same margins, and not a unitary value of margins (multiplying every margins by  $k$  multiplies by  $k^2$  the eigenvalues of the diagonalised matrix <sup>4</sup>).

Then, we calculate a matrix  $X = A N B$ , where,

$$\begin{cases} A = \text{Diag}(A^*) \\ A^* = [\text{Diag}(N B^*)]^{-1} \mathbf{1} \end{cases} \text{ and } \begin{cases} B = \text{Diag}(B^*) \\ B^* = \frac{n}{m} [\text{Diag}(N' A^*)]^{-1} \mathbf{1} \end{cases}$$

Besides, it is equivalent to get  $L = \begin{bmatrix} 1 & \dots & 1 \\ \vdots & & \vdots \\ 1 & \dots & 1 \end{bmatrix}$ ; then  $L_R = m \mathbf{1}$  and  $L_C = n \mathbf{1}$ . Thus,

$$\begin{cases} A = \text{Diag}(A^*) \\ A^* = m [\text{Diag}(N B^*)]^{-1} \mathbf{1} \end{cases} \text{ and } \begin{cases} B = \text{Diag}(B^*) \\ B^* = n [\text{Diag}(N' A^*)]^{-1} \mathbf{1} \end{cases}$$

There is a quasi normalisation simultaneously in rows and columns.

2) Second step. We diagonalise the matrix  $X' X$  in the space of  $\mathfrak{R}^m$  (or  $X X'$  in the space of  $\mathfrak{R}^n$ : classically, first eigenvalues are the same in one case and in the another case). A trivial eigenvalue is found:  $\lambda = \frac{n}{m}$ , associated to the eigenvector  $\mathbf{g}$ , the gravity center. Consider the

gravity center:  $\mathbf{g} = X' \mathbf{1}$ . With the chosen margins, we get  $\mathbf{g} = \frac{n}{m} \mathbf{1}_{(m,1)}$ . Then

$$X' X \mathbf{g} = X' \left( \frac{n}{m} X \mathbf{1}_{(m,1)} \right) = \frac{n}{m} X' X \mathbf{1}_{(m,1)} = \frac{n}{m} \mathbf{g}; \text{ thus } \lambda = \frac{n}{m} \text{ with } \mathbf{g} \text{ as eigenvector.}$$

It is possible to go back to the original data. Classically, the reconstruction of the matrix  $X$  from the eigenvalues and the eigenvectors remains possible. And it is possible to retrieve the matrix  $N$  from the matrix  $X$ , making a biproportion of  $X$  on the margins of  $N$ . So, the only condition is to know the margins of  $N$ .

---

<sup>4</sup> See later.

### 2.3. About metric and weights

Here, we use a simple Euclidean distance ( $M = I$ ). There are two cases if we apply another metric:

- The metric is used before applying biproportion. A theorem of invariance shows that this transformation (premultiply or postmultiply  $N$  by a diagonal matrix) ) does not affects the result of biproportion [ de MESNARD ] :

$$\text{if } Y = NM^{1/2}, \text{ then } X^* = A^* Y B^* = A^* N M^{1/2} B^* = A N B,$$

$$\text{with } A = A^* \text{ and } B = M^{1/2} B^* .$$

For the same reason, there is no need to weight data before applying biproportion.

- The metric is used after applying biproportion, or data are weighted after applying biproportion. The matrix to be diagonalised is  $M X' P X$  .

It is not logical to weight after applying biproportion because biproportion on  $N$  replaces the weighting on either rows or columns of  $N$  : every rows takes the same importance, every columns takes the same importance.

To use a non unitary metric is more logical. However, the  $\chi^2$  metric used by Factorial Analysis of Correspondences is a little artificial. An argument is that unitary distance between rows (reciprocally columns) gives more importance to the bigger columns (reciprocally rows). Another and good argument is the "distributional invariance" of the  $\chi^2$  metric. Biproportion allows to take into account the first argument but not the second: if this "distributional invariance" is essential, then one must get the  $\chi^2$  metric.

### 2.4. Advantages and disadvantages of the new method

1. In Biproportional factorial Analysis, there is only one matrix  $X$  to calculate the symmetric matrix  $X' X$  to be diagonalised. In the Factorial Analysis of Correspondences, there are two matrices  $X_R = N_R^{-1} N$  and  $X_C = N N_C^{-1}$  , and no symmetry.

In the biproportional factorial analysis, there is not such a strict duality (the classical duality between  $\mathfrak{R}^m$  and  $\mathfrak{R}^n$  remains), but there is an absolute symmetry between rows and columns. The new method does not requires to study correspondences: its interpretation is simple as in Principal Components Analysis.

There is no need to weight the distances in the Biproportional Factorial Analysis: this operation looks a little artificial even if it is absolutely necessary in the Factorial Analysis of Correspondences to obtain a strict duality of the analysis between rows and columns, called "correspondence". In the Factorial Analysis of Correspondences, the formulas  $X_R M_R X_R' P_R = N_R^{-1} N N_C^{-1} N'$  and  $X_R M_R X_R' P_R = N_R^{-1} N N_C^{-1} N'$  show that the  $N_{ij}$  are normalised one time in rows and one time in columns, that is to say, divided by  $N_{i.}$  in one case and by  $N_{.j}$  in one another case. In the Biproportional Factorial Analysis, data are simultaneously (quasi) normalised in rows and columns.



2. With this new method, we shall lose the "distributional invariance" characterising the  $\chi^2$  metric: if we replace two rows or columns by only one, we must recalculate the biproportion. It will be a disadvantage. However, we shall get a better analysis of data, because of the absolute symmetry between rows and columns.

### 3. ANNEXE

#### 3.1. The diagonalisation in factorial analysis

When data are centred, minimising the sum of squares of Euclidean distance is equivalent to maximise the sum of projections of points along the axis of regression. Generally, in the space of  $\mathcal{R}^m$ , if  $X$  is the matrix of data, and if  $\mathbf{u}$  is the unitary vector supported by this axis (with  $m$  elements),  $X \mathbf{u}$  is the projection of the points along the axis. Then, we must solve,

$$\text{Max } (X \mathbf{u})' (X \mathbf{u}) = \mathbf{u}' X' X \mathbf{u} , \text{ under constraint: } \mathbf{u}' \mathbf{u} = 1 .$$

$$\text{The Lagrangian is } L = \mathbf{u}' X' X \mathbf{u} - \lambda (1 - \mathbf{u}' \mathbf{u}) = 0$$

Deriving the Lagrangian with respect of  $\mathbf{u}$  and  $\lambda$ , we obtain:

$$\frac{\partial L}{\partial \mathbf{u}} = X' X \mathbf{u} - \lambda \mathbf{u} = 0$$

$$\frac{\partial L}{\partial \lambda} = 1 - \mathbf{u}' \mathbf{u} = 0$$

The second equation verifies the constraint, and the first equation provides:  $X' X \mathbf{u} = \lambda \mathbf{u}$ .  $\lambda$  is the eigenvalue and  $\mathbf{u}$  is the eigenvector associated to it. As we maximise, we select eigenvalues  $\lambda$  from the greater to the smaller. The associated eigenvector gives the direction of the axis. The cosine of the angle between a variable  $X_i$  and the axis is the correlation coefficient for this variable.

### 3.2. Example

We shall construct a small example with a 7x5 matrix.

$$N = \begin{bmatrix} 1 & 2 & 4 & 5 & 3 \\ 4 & 2 & 2 & 4 & 6 \\ 12 & 3 & 12 & 18 & 6 \\ 24 & 12 & 12 & 32 & 16 \\ 25 & 10 & 5 & 45 & 35 \\ 5 & 8 & 4 & 9 & 12 \\ 3 & 16 & 21 & 7 & 18 \end{bmatrix}$$

#### 3.2.1. Factorial Analysis of Correspondences

In  $\mathcal{R}^5$  :

$$M_R X'_R P_R X_R = N_C^{-1} N' N_R^{-1} N = \begin{bmatrix} 0.21329 & 0.11024 & 0.12259 & 0.32892 & 0.22495 \\ 0.15393 & 0.16267 & 0.17716 & 0.25543 & 0.25082 \\ 0.15120 & 0.15649 & 0.21711 & 0.25162 & 0.22359 \\ 0.20283 & 0.11281 & 0.12581 & 0.32780 & 0.23075 \\ 0.17340 & 0.13847 & 0.13974 & 0.28844 & 0.25995 \end{bmatrix}$$

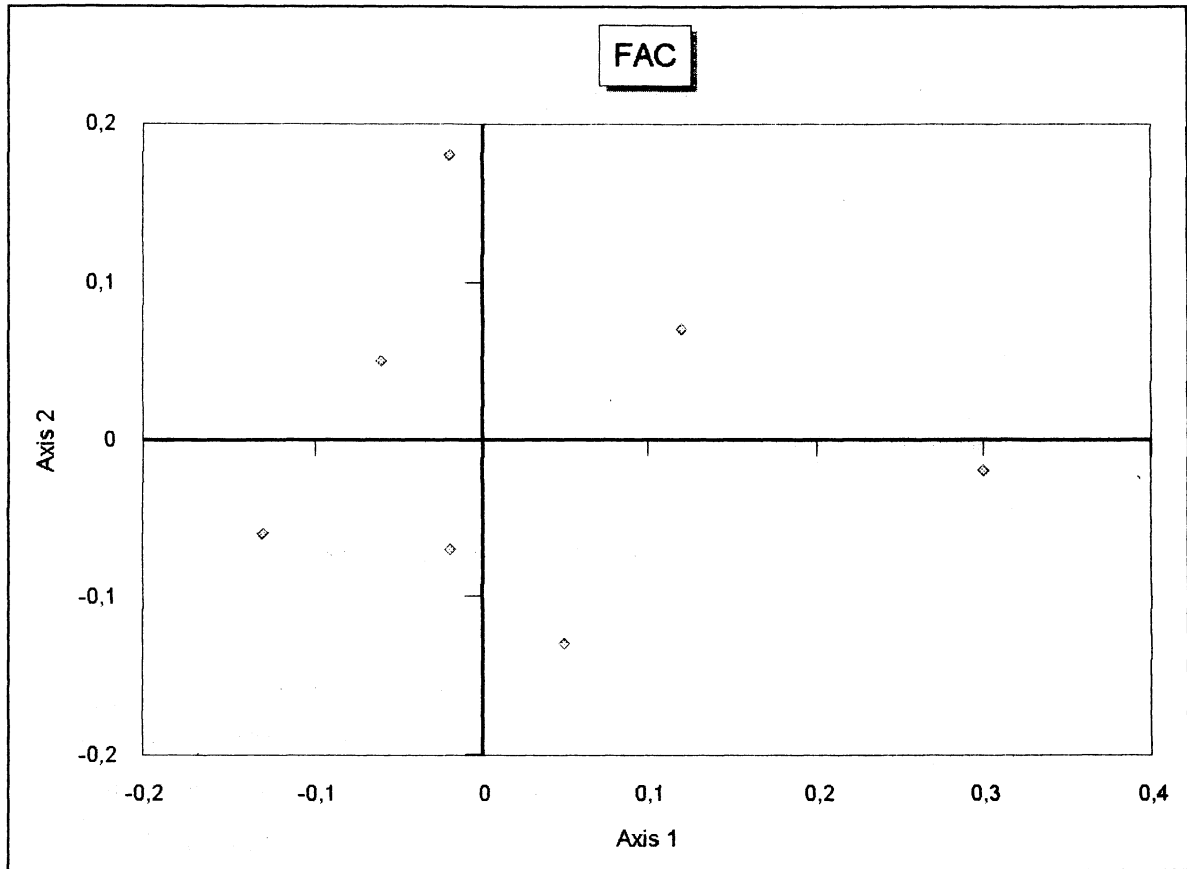
Eigenvalues are ( $\lambda_0 = 1$  is trivial):  $\lambda_1 = 0.127$ ,  $\lambda_2 = 0.04271$ ,  $\lambda_3 = 0.00766$ ,  $\lambda_4 = 0.00344$ .  
The axis explain respectively 70.2%, 23.6%, 4.2%, 1.9% .

Principal factors are:

$$v_1 = \begin{pmatrix} -0.40835 \\ 0.48708 \\ 0.69063 \\ -0.34242 \\ 0.04223 \end{pmatrix}, v_2 = \begin{pmatrix} -0.27837 \\ 0.36322 \\ -0.63762 \\ -0.15080 \\ 0.60106 \end{pmatrix}, v_3 = \begin{pmatrix} 0.54197 \\ 0.68898 \\ -0.26857 \\ -0.26594 \\ -0.29785 \end{pmatrix}, v_4 = \begin{pmatrix} -0.50683 \\ 0.57534 \\ -0.22489 \\ 0.47044 \\ -0.37445 \end{pmatrix}$$

Principal components are:

$$v_1 = \begin{pmatrix} 0.11620 \\ -0.02190 \\ -0.02081 \\ -0.06197 \\ -0.13179 \\ 0.05375 \\ 0.29900 \end{pmatrix}, v_2 = \begin{pmatrix} 0.07022 \\ -0.07449 \\ 0.17667 \\ 0.05398 \\ -0.06447 \\ -0.12682 \\ -0.02077 \end{pmatrix}, v_3 = \begin{pmatrix} -0.09184 \\ 0.00877 \\ -0.02405 \\ 0.04975 \\ -0.02747 \\ 0.03104 \\ -0.00328 \end{pmatrix}, v_4 = \begin{pmatrix} 0.06488 \\ -0.09396 \\ -0.01634 \\ 0.01150 \\ 0.00019 \\ 0.02394 \\ -0.00746 \end{pmatrix},$$



### 3.2.2. Biproportional Factorial analysis

$$X = \begin{bmatrix} 0.08019 & 0.19818 & 0.31807 & 0.23728 & 0.16628 \\ 0.26723 & 0.16510 & 0.13249 & 0.15814 & 0.27704 \\ 0.28299 & 0.08742 & 0.28061 & 0.25120 & 0.09779 \\ 0.29732 & 0.18369 & 0.14741 & 0.23459 & 0.13699 \\ 0.26843 & 0.13267 & 0.05323 & 0.28593 & 0.25973 \\ 0.15398 & 0.30442 & 0.12215 & 0.16402 & 0.25542 \\ 0.04985 & 0.32853 & 0.34604 & 0.06884 & 0.20674 \end{bmatrix}$$

In  $\mathcal{R}^5$  (rows are points):

$$X'X = \begin{bmatrix} 0.34458 & 0.23823 & 0.23450 & 0.30756 & 0.27513 \\ 0.23823 & 0.32612 & 0.29445 & 0.24866 & 0.29254 \\ 0.23450 & 0.29445 & 0.35669 & 0.26057 & 0.25379 \\ 0.30756 & 0.24866 & 0.26057 & 0.31284 & 0.27036 \\ 0.27513 & 0.29254 & 0.25379 & 0.27036 & 0.30817 \end{bmatrix}$$

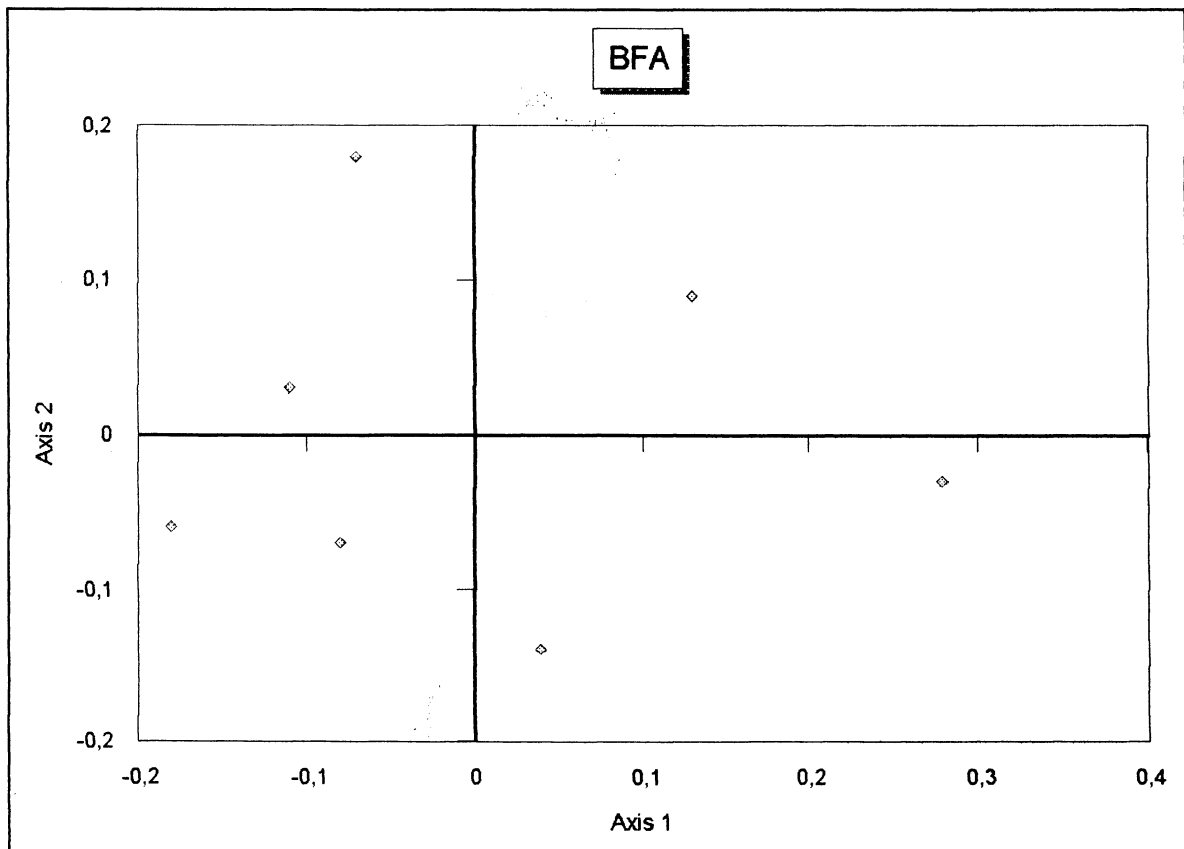
Eigenvalues are ( $\lambda_0 = 1.4$  is trivial):  $\lambda_1 = 0.15052$ ,  $\lambda_2 = 0.07252$ ,  $\lambda_3 = 0.01587$ ,  $\lambda_4 = 0.00948$ . The axis explain respectively 65.6%, 29.2%, 6.4%, 3.8%. The last third axis are stronger than in Factorial Analysis of Correspondences.

Principal factors are:

$$v_1 = \begin{pmatrix} -0.61566 \\ 0.42591 \\ 0.56340 \\ -0.34859 \\ -0.02507 \end{pmatrix}, v_2 = \begin{pmatrix} -0.12260 \\ 0.45084 \\ -0.61993 \\ -0.27537 \\ 0.56706 \end{pmatrix}, v_3 = \begin{pmatrix} -0.63592 \\ -0.14593 \\ -0.12910 \\ 0.72262 \\ 0.18832 \end{pmatrix}, v_4 = \begin{pmatrix} 0.03933 \\ 0.62773 \\ -0.28564 \\ 0.28367 \\ -0.66509 \end{pmatrix}$$

Principal components are:

$$v_1 = \begin{pmatrix} 0.12736 \\ -0.08163 \\ -0.06891 \\ -0.10697 \\ -0.18495 \\ 0.04010 \\ 0.27501 \end{pmatrix}, v_2 = \begin{pmatrix} 0.08872 \\ -0.07309 \\ 0.18296 \\ 0.03194 \\ -0.06245 \\ -0.14232 \\ -0.02576 \end{pmatrix}, v_3 = \begin{pmatrix} 0.08180 \\ -0.04469 \\ -0.02900 \\ -0.03958 \\ 0.05860 \\ 0.00851 \\ -0.03564 \end{pmatrix}, v_4 = \begin{pmatrix} -0.00658 \\ -0.06310 \\ -0.00793 \\ 0.06033 \\ -0.01300 \\ 0.03891 \\ -0.00862 \end{pmatrix}$$



In this example, results are enough close to those of Correspondence Analysis.

## REFERENCES

- J.P. Benzécri, *L'analyse des données*, DUNOD, PARIS, 1973.
- L. Lebart, A. Morineau, J.-P. Fénelon, *Traitement des données statistiques*, DUNOD, PARIS, 2e ed., 1982.
- L. de Mesnard, "Unicity of biproportion", to be published in *SIAM Journal on Matrix Analysis and Applications*.
- G. Saporta, *Probabilités, analyse des données et statistique*, EDITIONS TECHNIP, PARIS, 1990.
- M. Volle, *Analyse des données*, ECONOMICA, PARIS, 3e ed, 1985.