

REPORT DOCUMENTATION PAGE			Form Approved OMB NO. 0704-0188	
<small>Public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comment regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302, and to the Office of Management and Budget, Paperwork Reduction Project (0704-0188), Washington, DC 20503.</small>				
1. AGENCY USE ONLY (Leave blank)		2. REPORT DATE March 1997		3. REPORT TYPE AND DATES COVERED Technical - 97-06
4. TITLE AND SUBTITLE An Alternative to Correspondence Analysis Using Hellinger Distance			5. FUNDING NUMBERS DAAH04-96-1-0082	
6. AUTHOR(S) C.R. Rao				
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) Center for Multivariate Analysis Department of Statistics 417 Thomas Building Penn State University University Park, PA 16802			8. PERFORMING ORGANIZATION REPORT NUMBER	
9. SPONSORING / MONITORING AGENCY NAME(S) AND ADDRESS(ES) U.S. Army Research Office P.O. Box 12211 Research Triangle Park, NC 27709-2211			10. SPONSORING / MONITORING AGENCY REPORT NUMBER	
11. SUPPLEMENTARY NOTES The views, opinions and/or findings contained in this report are those of the author(s) and should not be construed as an official Department of the Army position, policy or decision, unless so designated by other documentation.				
12a. DISTRIBUTION / AVAILABILITY STATEMENT  Approved for public release; distribution unlimited.			12 b. DISTRIBUTION CODE	
13. ABSTRACT (Maximum 200 words)  In this paper, a general theory of canonical coordinates is developed for reduction of dimensionality in multivariate data, assessing the loss of information and plotting higher dimensional data in two or three dimensions for visual displays. The theory is applied to data in two way tables with variables in one category and samples (individual or populations) in the other. The method is applicable to data with continuous measurements on the variables as well as to frequencies of attributes. An alternative distance is suggested. The new method has some attractive features and does not suffer from some inherent drawbacks resulting from the use of the chi-square distance and variable sample sizes for the populations in the correspondence analysis. The technique of biplots where the populations and the variables are represented on the same chart is discussed.				
14. SUBJECT TERMS Canonical coordinates, Chisquare distance, Contingency tables, Correspondence analysis, Hellinger distance, Matrix approximation, Principal component analysis.			15. NUMBER OF PAGES 19	
			16. PRICE CODE	
17. SECURITY CLASSIFICATION OR REPORT UNCLASSIFIED	18. SECURITY CLASSIFICATION OF THIS PAGE UNCLASSIFIED	19. SECURITY CLASSIFICATION OF ABSTRACT UNCLASSIFIED	20. LIMITATION OF ABSTRACT UL	

# DISCLAIMER NOTICE



**THIS DOCUMENT IS BEST  
QUALITY AVAILABLE. THE  
COPY FURNISHED TO DTIC  
CONTAINED A SIGNIFICANT  
NUMBER OF PAGES WHICH DO  
NOT REPRODUCE LEGIBLY.**

**AN ALTERNATIVE TO CORRESPONDENCE ANALYSIS  
USING HELLINGER DISTANCE**

**C. Radhakrishna Rao**

Technical Report 97-06

March 1997

Center for Multivariate Analysis  
417 Thomas Building  
Penn State University  
University Park, PA 16802

The research work of the author is supported by the Army Research Office under Grant DAAHO4-96-1-0082. The United States Government is authorized to reproduce and distribute reprints for governmental purposes notwithstanding any copyright notation hereon.

19970521 176

Typeset by  $\mathcal{A}\mathcal{M}\mathcal{S}$ - $\mathcal{T}\mathcal{E}\mathcal{X}$

## AN ALTERNATIVE TO CORRESPONDENCE ANALYSIS USING HELLINGER DISTANCE

C. Radhakrishna Rao

**ABSTRACT.** In this paper, a general theory of canonical coordinates is developed for reduction of dimensionality in multivariate data, assessing the loss of information and plotting higher dimensional data in two or three dimensions for visual displays. The theory is applied to data in two way tables with variables in one category and samples (individual or populations) in the other. The method is applicable to data with continuous measurements on the variables as well as to frequencies of attributes. An alternative to the usual correspondence analysis of contingency tables based on Hellinger rather than the chisquare distance is suggested. The new method has some attractive features and does not suffer from some inherent drawbacks resulting from the use of the chi-square distance and variable sample sizes for the populations in the correspondence analysis. The technique of biplots where the populations and the variables are represented on the same chart is discussed.

### 1. Canonical Coordinates

The concept of canonical variates (coordinates) was introduced in an early paper by the author (Rao (1948)) for graphical representation of taxonomical units characterized by multiple measurements. This was, perhaps, the first attempt to reduce high dimensional data to two or three dimensions using an objective criterion for purposes of graphical displays. Since then, graphical representation of multivariate data for visual examination of clusters, outliers and other structures in the data has been an active field of research. Some of the developments are biplots (Gabriel (1971), Gifi (1990), Nishisato (1980), Gower (1993), Greenacre (1993)), multidimensional scaling (Kruskal and Wish (1978)), correspondence analysis (Benzécri (1992), Greenacre (1984)), Chernoff's faces (Chernoff (1973)) and parallel coordinates (Mahalanobis, Mazumdar and Rao (1949), Wegman (1990)). Cavalli-Sforza (1991) uses canonical coordinates (variables) in interpreting the evolution of human populations.

The object of the present paper is to briefly review the concept of canonical coordinates as originally introduced in 1948 and later elaborated in Rao (1964, 1979, 1980, 1985) in the light of modern developments and present an alternative to the current practice of correspondence analysis, which seems to have some attractive properties.

In Section 2 we consider the general problem of transforming the points of a  $p$ -dimensional vector space endowed with a specified inner product to a lower dimensional Euclidean space with the usual definition of inner product and distance. The solution to the problem is considered in a more general set up than what is possible through the use of Eckart and Young (1936) theorem. In Section 3, some measures are introduced to assess the loss of information in reduction of dimensionality. The role of biplots and their interpretation are also discussed. An alternative to correspondence analysis applied to contingency tables based on Hellinger rather than the chisquare distance is given in Section 4.

It is argued that the chisquare distance used in correspondence analysis is not an intrinsic measure of the difference between two given population distributions as it depends to some extent on the whole set of populations considered in the study, and also on the sample sizes available for the estimation of population distributions. In such a case, the configuration of a subset of the populations as revealed by correspondence analysis may depend on what other populations are included in the analysis. An example is given to show how anomalies can arise in correspondence analysis based on the chisquare distance. On the other hand no such anomalies arise with the use of Hellinger distance.

---

1991 *Mathematics Subject Classification.* 62H30, 62H17.

*Key words and phrases.* Canonical coordinates, Chisquare distance, Contingency tables, Correspondence analysis, Hellinger distance, Matrix approximation, Principal component analysis.

The research work of this paper is sponsored by the US Army Research Office under Grant DAAH04-96-1-0082.

## 2. Reduction of Dimensionality

The problem we consider may be stated as follows. Let  $X = (X_1 : \dots : X_m)$  be a  $p \times m$  data matrix, with the  $i$ -th column vector  $X_i$  representing measurements of  $p$  variables made on the  $i$ -th population (individual or unit). The column vector  $X_i$  will be referred to as the  $i$ -th population profile (PP). The PP's can be represented as  $m$  points in a  $p$ -dimensional vector space  $R^p$  with a specified inner product and the associated norm

$$(x, y) = x' M y, \quad x, y \in R^p \quad (2.1)$$

$$\|x\| = (x, x)^{1/2}, \quad x \in R^p \quad (2.2)$$

where  $M$  is a positive definite matrix. We may call this the Mahalanobis or  $M$ -space. In practical situations, it may be necessary to attach a weight  $w_i \geq 0$  to the  $i$ -th PP, the exact use of which will be detailed in the following discussion. We represent the vector  $(w_1, \dots, w_m)'$  by  $w$  and the diagonal matrix with  $w_i$  as the  $i$ -th diagonal element by  $W$ . The  $M$ -space with weight as an additional dimension will be referred to as  $WM$ -space. [In our treatment we consider  $W$  as a general positive definite matrix to cover more general applications].

The problem is to find a  $k \times m$  matrix

$$Y = (Y_1 : \dots : Y_m) \quad (2.3)$$

with  $k < p$  for representing the PP's in a  $k$ -dimensional Euclidean space ( $E^k$ ) with the usual inner product,  $x'y$  for  $x, y \in E^k$ , and the  $k$ -vector  $Y_i$  as the profile of the  $i$ -th population, in such a way that the relative positions of the PP's in the  $M$ -space (in terms of distances between profiles) are preserved to the extent possible in  $E^k$ . For this purpose, we need to have a criterion for measuring the loss of information in reducing the dimension of the profile space, by minimizing which we obtain an optimum solution for (2.3).

The relative positions of the PP's in the  $M$ -space can be described by what may be called a *configuration matrix*

$$C = (X - \xi 1')' M (X - \xi 1') = ((X_i - \xi)' M (X_j - \xi)) = (c_{ij}) \quad (2.4)$$

where  $\xi$  is some chosen reference (profile) vector and the  $c_{ij}$ 's represent the distances and angles between profiles.

The corresponding configuration about the origin in the reduced space  $E^k$  is  $Y'Y$ . The problem then reduces to minimizing

$$\|C - Y'Y\| \quad (2.5)$$

with respect to  $Y$ , a  $k \times m$  matrix as defined in (2.3), for a suitably chosen matrix norm. The following theorem proved in Rao (1979, 1980, 1985) provides the solution.

**THEOREM 1.** Consider the s.v.d. (singular value decomposition)

$$M^{1/2}(X - \xi 1')W^{1/2} = \lambda_1 U_1 V_1' + \dots + \lambda_p U_p V_p' \quad (2.6)$$

with singular values  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$ , where  $M^{1/2}$  and  $W^{1/2}$  are symmetric square roots of  $M$  and  $W$ . Then the choice

$$Y = Y_{(k)} = \begin{pmatrix} \lambda_1 V_1' W^{-1/2} \\ \vdots \\ \lambda_k V_k' W^{-1/2} \end{pmatrix} \quad (2.7)$$

or conventionally written in the transposed form

$$\lambda_1 W^{-1/2} V_1, \lambda_2 W^{-1/2} V_2, \dots, \lambda_k W^{-1/2} V_k \quad (2.8)$$

where the components of the  $i$ -th  $m$ -vector are the  $i$ -th canonical coordinates (i.e., the coordinates in the  $i$ -th dimension of the reduced space) for the different populations, minimizes (2.5) for any  $(W, W)$ -invariant norm as defined in Note 2.1. We call these coordinates the canonical coordinates for populations (CCP).

*Note 2.1.*  $(A, B)$ -invariant norm of an  $m \times n$  matrix is the usual norm (satisfying the postulates of a norm) with the additional property

$$\|C \cdot D\| = \|\cdot\| \quad \text{for any } C, D \text{ such that } C'AC = A, D'BD = B \quad (2.9)$$

where  $C$  is an  $m \times m$  matrix,  $D$  is an  $n \times n$  matrix, and  $A$  and  $B$  are positive definite matrices of orders  $m$  and  $n$  respectively. This is a generalization given in Rao (1980) of a unitarily invariant norm defined by von Neumann (1937) with  $A$  and  $B$  as unit matrices.

*Note 2.2.* In our applications, we indicate some choices of the reference vector  $\xi$ . However, we note that a further minimization of (2.5) with respect to  $\xi$  leads to the choice

$$\hat{\xi} = (1'W1)^{-1}XW1 \quad (2.10)$$

where  $1$  is the column vector of unities.

*Note 2.3.* Using the notation

$$\begin{aligned} \Lambda_{(i)} &= \text{Diag}(\lambda_1, \dots, \lambda_i) \\ U_{(i)} &= (U_1 : \dots : U_i), \quad V_{(i)} = (V_1 : \dots : V_i) \end{aligned}$$

we may write the solution  $Y$  given in (2.7) in the concise form

$$Y_{(k)} = \Lambda_{(k)} V'_{(k)} W^{-1/2}. \quad (2.11)$$

*Note 2.4.* In the expression (2.6), a symmetric square root of a positive definite matrix is used. It can be computed in a simple way as follows. If  $A$  is a positive definite matrix of order  $p$  with the spectral decomposition

$$A = \Sigma \lambda_i^2 Q_i Q_i' = Q \Lambda^2 Q'$$

where  $Q = (Q_1 : \dots : Q_p)$ , then

$$\begin{aligned} A^{1/2} &= \Sigma \lambda_i Q_i Q_i' = Q \Lambda Q' \\ A^{-1/2} &= \Sigma (\lambda_i)^{-1} Q_i Q_i' = Q \Lambda^{-1} Q' \end{aligned} \quad (2.12)$$

We may look at the problem in a slightly different way by defining what is called the dispersion matrix between profiles

$$B = (X - \xi 1')W(X - \xi 1')' = (b_{ij}) \quad (2.13)$$

where  $b_{ii}$  is the weighted variance of the  $i$ -th variable and  $b_{ij}$  is the weighted covariance between the  $i$ -th and  $j$ -th variables across the profiles. Consider an approximation,  $Z_i \in R^p$  to  $(X_i - \xi)$ , with the restriction that  $Z_1, \dots, Z_m$  lie in a  $k$  dimensional subspace of  $R^p$ , in which case we have the representation

$$Z = (Z_1 : \dots : Z_m) = AC \quad (2.14)$$

where  $A$  is a  $p \times k$  matrix whose columns span the subspace and  $C$  is a  $k \times m$  matrix. Without loss of generality we may choose  $A$  to satisfy the condition  $A'MA = I$  (i.e., the columns of  $A$  are orthonormal in the  $M$ -space). The dispersion matrix between profiles in the reduced space is  $ACWC'A'$ , and we choose  $A$  and  $C$  such that

$$\|B - ACWC'A'\| \quad (2.15)$$

is a minimum for an appropriate norm of the matrix. The solution is given in Theorem 2, which is proved on the same lines as in Theorem 1.

**THEOREM 2.** Consider the same s.v.d. as in Theorem 1

$$M^{1/2}(X - \xi 1')W^{1/2} = \lambda_1 U_1 V_1' + \dots + \lambda_p U_p V_p'$$

Then the optimum choice of  $AC$  which minimizes (2.15) for any  $(M, M)$ -invariant norm is

$$AC_{(k)} = M^{-1/2}(\lambda_1 U_1 V_1' + \dots + \lambda_k U_k V_k')W^{-1/2} \quad (2.16)$$

where the suffix  $(k)$  is introduced to indicate the dimension of the reduced space. We may choose

$$\begin{aligned} A &= M^{-1/2}(U_1 : \dots : U_k) = M^{-1/2}U_{(k)} \\ C_{(k)} &= \begin{pmatrix} \lambda_1 V_1' W^{-1/2} \\ \vdots \\ \lambda_k V_k' W^{-1/2} \end{pmatrix} = \Lambda_{(k)} V_{(k)}' W^{-1/2}. \end{aligned} \quad (2.17)$$

*Note 2.5.* We may represent the profiles by plotting the columns of  $C_{(k)}$  in a  $k$ -dimensional Euclidean space, which is the same solution as that obtained in Theorem 1.

A geometric approach to the problem of reduction of dimensionality is to fit a  $k$ -dimensional plane to the data. A set of  $m$  points on a  $k$ -plane can be written as

$$\xi 1' + AC \quad (2.18)$$

where  $A$  is a  $p \times k$  matrix and  $C$  is a  $k \times m$  matrix. We determine  $A, C, \xi$  such that

$$\|X - \xi 1' - AC\| \quad (2.19)$$

is a minimum for a suitably chosen norm. The solution is given in Theorem 3, which is proved on the same lines as in Theorems 1 and 2.

**THEOREM 3.** *Consider the same s.v.d. as in Theorem 1. Then the choices of  $A$  and  $C$  as in Theorem 2 and  $\xi = (1'W1)^{-1}XW1$  as in (2.10) minimize any  $(M, W)$ -invariant norm of (2.19).*

*Note 2.6.* We may also look at the problem in some other ways. Let  $T$  be a  $k \times p$  matrix providing a transformation of the column vectors of  $X$  to  $Y = TX$  in a  $k$ -dimensional space with the induced inner product matrix  $TM^{-1}T'$ . The squared distance between the  $i$ -th and  $j$ -th profiles is

$$D_{ij}^2 = (X_i - X_j)'M(X_i - X_j) \quad (2.20)$$

in the full space, and

$$D_{ij(k)}^2 = (X_i - X_j)'T'(TM^{-1}T')^{-1}T(X_i - X_j) \quad (2.21)$$

in the reduced space. By definition  $D_{ij(k)}^2 \leq D_{ij}^2$ . We may then choose  $T$  by minimizing some function of the differences or ratios of  $D_{ij}^2$  and  $D_{ij(k)}^2$ .

One of the functions suggested in Rao (1948) was the difference in the weighted sum of all possible differences

$$\Sigma \Sigma w_i w_j (D_{ij}^2 - D_{ij(k)}^2) \quad (2.22)$$

which leads to the same solution for  $Y = TX$  as in Theorems 1, 2 and 3.

Another method is to choose  $T$  by maximizing the minimum of  $D_{ij(k)}^2$  over all  $i$  and  $j$  as suggested by Eslava-Gomez and Marriott (1993), or by maximizing the minimum of the ratios  $D_{ij(k)}^2/D_{ij}^2$ . Both these methods are computationally very complex, but can be managed when  $p$  is small.

*Note 2.7.* The choices of  $M$  and  $W$  as inputs in the analysis for canonical coordinates need some discussion. The choice of  $M$  is related to the distance measure between profiles appropriate to a given investigation. In taxonomical classification,  $M$  is generally chosen as the inverse of the variance-covariance (dispersion) matrix of the measurements on units within taxa leading to Mahalanobis (1936) distance (see Rao (1945, 1947)). The matrix  $W$  is taken to be diagonal with the  $i$ -th diagonal element  $w_i$  proportional to the number of individuals sampled from the  $i$ -th taxa to estimate its profile. For a chosen  $M$ , the configuration of the profiles in the reduced space will depend on  $W$ , but is likely to be robust provided the  $w_i$ 's are not widely different. In the study reported in Rao (1948), all the  $w_i$ 's were chosen as equal although the sample sizes for different populations were different. However, the choice of  $w_i$ 's as proportional to sample sizes enables us to test hypotheses on goodness of fit of lower dimensional planes to the observed profiles. For details, the reader is referred to Rao (1973, pp. 556-560, 1985).

If we desire that the configuration of a subset of profiles to be better preserved in the reduced space than the others, then we have to give bigger weights to those profiles.

*Note 2.8.* In many situations we have a data matrix  $X$  giving the measurements of  $p$  variables made on  $m$  individuals without any further information to guide us in the choices of the  $M$  and  $W$  matrices. In such cases, the usual choices of  $M$  and  $W$  are the unit matrices and the resulting canonical coordinate analysis is the Principal Component Analysis (PCA) introduced by Hotelling. Some characterizations of the principal components and their applications are given in papers by Rao (1958, 1964, 1987). It is also the practice to apply PCA on  $CX$ , i.e., after a suitable scaling of the measurements. One choice of  $C$  is a diagonal matrix with the  $i$ -th diagonal element  $c_i = s_{ii}^{-1/2}$ , where  $s_{ii}$  is the  $i$ -th diagonal element of the matrix

$$(X - \bar{X}1')(X - \bar{X}1')'. \quad (2.23)$$

This procedure is equivalent to using the canonical coordinate analysis choosing  $M = C$  and  $W = I$ . Another possibility which has not been considered before is the choice,  $c_i = 1/m_i$  where  $m_i$  is a measure of location such as the mean or median of the measurements on the  $i$ -th variable.

*Note 2.9.* A more general problem not considered in this paper is as follows. The basic space is somewhat general with a specified nonnegative proximity index between any two points. Given a set of points with the matrix of proximity indices between points, the problem is to transform the points to a low dimensional Euclidean space such that the inequality relationships between proximity indices are maintained to the extent possible in the corresponding Euclidean distances. Such a transformation is achieved through the algorithm for multidimensional scaling as developed by Kruskal and Wish (1978).

### 3. Loss of information

The representation of the PP's in a lower dimensional space will entail some loss of information depending on the object of statistical analysis. However, we provide some general criteria for assessing the amount of distortion in the configuration of the profiles due to reduction of dimensionality.

In Theorems 1 and 2 of Section 2, it is shown that the best approximation to  $X$  in the reduced space is

$$\hat{X} = \xi 1' + M^{-1/2} U_{(k)} \Lambda_{(k)} V'_{(k)} W^{-1/2} \quad (3.1)$$

so that the matrix

$$D_1 = X - \hat{X} = M^{-1/2} (\lambda_{k+1} U_{k+1} V'_{k+1} + \dots + \lambda_p U_p V'_p) W^{-1/2} \quad (3.2)$$

gives a complete account of the errors in individual profiles due to reduction.

The configuration of the profiles in the reduced space is

$$C_{(k)} = W^{-1/2} V_{(k)} \Lambda_{(k)}^2 V'_{(k)} W^{-1/2} \quad (3.3)$$

so that the matrix

$$D_2 = C_{(p)} - C_{(k)} = W^{-1/2} (\lambda_{k+1}^2 V_{k+1} V'_{k+1} + \dots + \lambda_p^2 V_p V'_p) W^{-1/2} \quad (3.4)$$

measures the distortion in the configuration, where  $C_{(p)} = C$  as defined in (2.4). An overall (weighted) measure of loss of information is the ratio of

$$\text{trace } W^{1/2} D_2 W^{1/2} = \lambda_{k+1}^2 + \dots + \lambda_p^2, \quad (3.5)$$

to the total variation  $(\lambda_1^2 + \dots + \lambda_p^2)$ , which can be written as

$$1 - \left( \sum_{i=1}^k \lambda_i^2 \right) / \left( \sum_{i=1}^p \lambda_i^2 \right). \quad (3.6)$$

It is more important to assess the distortions in the inter profile squared distances. The matrix of these squared distances denoted by  $S$  can be computed from the configuration matrix  $C$  using the formula

$$S = c1' + 1c' - 2C \quad (3.7)$$

where  $c$  is the vector of the diagonal elements of  $C$ . The corresponding matrix in the reduced space is

$$S_{(k)} = c_{(k)} 1' + 1c'_{(k)} - 2C_{(k)} \quad (3.8)$$



so that the matrix

$$D_3 = S - S_{(k)} = (d_{ij}^*) \quad (3.9)$$

measures the deficiencies in the distances due to reduction of dimensionality. An over all measure of deficiency is

$$\Sigma \Sigma w_i w_j d_{ij}^* = \lambda_{k+1}^2 + \dots + \lambda_p^2 \quad (3.10)$$

which is the same as in (3.5).

The dispersion matrix between profiles in the whole space, as introduced in (2.13) is

$$B = (X - \xi 1') W (X - \xi 1')' = (b_{ij}) \quad (3.11)$$

while the corresponding matrix in the reduced  $k$ -dimensional space is

$$B_{(k)} = M^{-1/2} U_{(k)} \Lambda_{(k)}^2 U_{(k)}' M^{-1/2} = (b_{ij(k)}) \quad (3.12)$$

The proportion of the between profile variance in the  $i$ -th variable explained by the first  $k$  canonical variates (coordinates) is

$$b_{ii(k)}/b_{ii}, i = 1, \dots, p. \quad (3.13)$$

For an interpretation of the canonical coordinates in different dimensions it would be useful to compute the proportion of variance in each variable explained by each of the canonical variates, i.e., to obtain a decomposition of (3.13) in terms of canonical variates. For this purpose, we introduce the matrices

$$E_1 = M^{-1/2} (\lambda_1 U_1 : \dots : \lambda_p U_p) = (e_{ij}) \quad (3.14)$$

$$E_2 = (e_{ij}/\sqrt{b_{ii}}) = (f_{ij}) \quad (3.15)$$

where  $b_{ii}$  is as defined in (3.11). Let  $E_{i(k)}$  be the matrix obtained by retaining only the first  $k$  columns in  $E_i$  for  $i=1,2$ . Then it is seen that

$$E_1 E_1' = B, E_{1(k)} E_{1(k)}' = B_{(k)}. \quad (3.16)$$

Let us consider the matrix  $E_{1(k)}$  and define what may be called canonical coordinates for variables (CCV) in  $k$  dimensions as follows.

TABLE 1. Canonical coordinates for variables

variable	dim 1	dim 2		dim $k$
1	$e_{11}$	$e_{12}$	$\dots$	$e_{1k}$
2	$e_{21}$	$e_{22}$	$\dots$	$e_{2k}$
.	.	.	$\dots$	.
$p$	$e_{p1}$	$e_{p2}$	$\dots$	$e_{pk}$

If we plot the variables as points in  $E^k$  using the row coordinates in different dimensions, then the scalar products of the vectors representing the variables are the elements of  $B_{(k)}$ , the best  $k$ -dimensional approximation to  $B$ .

There is some advantage in plotting the variables using the standardized coordinates  $(f_{ij})$  defined in (3.15) as shown in Table 2.

The magnitudes in the right hand block of Table 2 indicate the influence of different variables in each dimension (canonical variate) in the reduced space. This may enable us to associate each dimension with certain variables. We may plot the variables using the standardized CCV's in the same chart as the canonical coordinates for the profiles. It is seen that all variable points lie inside the unit sphere in  $E^k$ , and the variables close to the surface of the sphere have greater influence on the canonical variates.

It may also be mentioned that it is the usual practice in a biplot to represent the  $i$ -th variable as a directed line using the direction cosines proportional to the  $i$ -th row elements in the matrix

$$E_{1(k)} = M^{-1/2} (U_1 : \dots : U_k) \quad (3.17)$$

TABLE 2. Standardized CCV's and the variance explained by each canonical variate

Variable	Standardized	Proportion of variance	
	coordinates	explained	
	dim 1 ... dim $k$	dim 1 ... dim $k$	total
1	$f_{11} \dots f_{1k}$	$f_{11}^2 \dots f_{1k}^2$	$\Sigma f_{1i}^2$
$\vdots$	$\vdots \quad \vdots$	$\vdots \quad \vdots$	$\vdots$
$p$	$f_{p1} \dots f_{pk}$	$f_{p1}^2 \dots f_{pk}^2$	$\Sigma f_{pi}^2$

in which case the projections of a profile point in these directions are proportional to the approximate coordinates of the profile in the original space (see Gabriel (1971) and Greenacre (1993)).

*Note 3.1.* We may consider the  $k$  columns in Table 1 of the CCV's as  $k$  points in the  $p$ -dimensional variable space. These points were termed as *typical profiles* in Rao (1964), in the sense that the variance-covariance matrix of the variables computed from them provides the best approximation to that computed from all the original profiles.

*Note 3.2.* The standardized CCV's are not the coordinates for row profiles. They are used for interpreting the CC's of column profiles. If a representation of row profiles is needed, we consider the matrix  $X'$  with appropriate choices of the  $M$  and  $W$  matrices (which may be different from those used for column profiles) and repeat the analysis indicated in (2.6)-(2.8).

#### 4. Application to two way contingency tables

We consider dichotomous categorical data with  $s$  rows and  $m$  columns and  $n_{ij}$  observations in the  $(i, j)$ -th cell. Define

$$\begin{aligned}
 N &= (n_{ij}), n_{i.} = \sum_{j=1}^m n_{ij}, n_{.j} = \sum_{i=1}^s n_{ij}, n = \sum_{i=1}^s \sum_{j=1}^m n_{ij} \\
 R &= \text{Diag } (n_{1.}/n, \dots, n_{s.}/n), C = \text{Diag } (n_{.1}/n, \dots, n_{.m}/n) \\
 P &= n^{-1}NC^{-1} = \begin{pmatrix} p_{1|1} & \dots & p_{1|m} \\ \vdots & \dots & \vdots \\ p_{s|1} & \dots & p_{s|m} \end{pmatrix}, \quad \text{column profiles} \\
 Q &= n^{-1}R^{-1}N = \begin{pmatrix} q_{1|1} & \dots & q_{m|1} \\ \vdots & \dots & \vdots \\ q_{1|s} & \dots & q_{m|s} \end{pmatrix}, \quad \text{row profiles}
 \end{aligned} \tag{4.1}$$

$$(p_1, \dots, p_s)' = p = R1, \quad q = C1 = (q_1, \dots, q_m)'. \tag{4.2}$$

The problem is to represent the column (row) profiles as points in  $E^k, k < s$ , such that the Euclidean distances between points reflect specified affinities between the corresponding column (row) profiles.

The technique developed for this purpose by Benzécri (1992) is known as correspondence analysis (CA) which can be identified as canonical coordinate analysis. For instance, for representing the column profiles by this method, one chooses

$$X = P, \quad M = R^{-1}, \quad W = C \tag{4.3}$$

and applies the analysis described in Theorem 1 (equation (2.6)). Thus one finds the s.v.d. of

$$R^{-1/2}(P - p1')C^{1/2} = \lambda_1 U_1 V_1' + \dots + \lambda_s U_s V_s' \tag{4.4}$$

giving the coordinates for the column profiles in  $E^k$

$$\lambda_1 C^{-1/2} V_1, \lambda_2 C^{-1/2} V_2, \dots, \lambda_k C^{-1/2} V_k \tag{4.5}$$

where the components of  $i$ -th vector are the coordinates of the profiles in the  $i$ -th dimension. The standardized canonical coordinates in  $E^k$  for the rows, as described in (3.15), obtained from the same s.v.d. as in (4.4) are

$$\lambda_1 \Delta^{-1} R^{1/2} U_1, \lambda_2 \Delta^{-1} R^{1/2} U_2, \dots, \lambda_k \Delta^{-1} R^{1/2} U_k \quad (4.6)$$

where the components of the  $i$ -th vector are the coordinates of the rows in the  $i$ -th dimension and  $\Delta$  is a diagonal matrix with the  $i$ -th diagonal element as the square root of the  $i$ -th diagonal element of the matrix

$$R^{1/2}(\lambda_1^2 U_1 U_1' + \dots + \lambda_s^2 U_s U_s') R^{1/2} = (P - p1')C(P - p1')'. \quad (4.7)$$

The coordinates (4.6) do not represent the row profiles but are useful in interpreting the different dimensions of the column profiles. The coordinates for representing the row profiles in correspondence analysis are given in (4.11).

Implicit in this analysis is the choice of measure of affinity between the  $i$ -th and  $j$ -th profiles as the squared distance (with  $p_1, \dots, p_s$  as defined in (4.2))

$$d_{ij}^2 = \frac{(p_{1|i} - p_{1|j})^2}{p_1} + \dots + \frac{(p_{s|i} - p_{s|j})^2}{p_s} \quad (4.8)$$

which is the chisquare distance. The squared Euclidean distance in  $E^k$ , the reduced space, between the points representing the  $i$ -th and  $j$ -th profiles is an approximation to (4.8). Thus the clusters we see in the Euclidean representation is based on the affinities as measured by the chisquare distance (4.8).

Why should one choose the chisquare distance to measure the affinities between profiles? Some of the advantages mentioned by Benzécri and Greenacre are as follows.

1. Note that the expression in (4.4)

$$R^{-1/2}(P - p1')C^{1/2} = R^{1/2}(Q - 1q')C^{-1/2} = T \quad (4.9)$$

so that if we need a representation of the row (as population) profiles in  $E^k$ , we use the same s.v.d. as in (4.4)

$$R^{1/2}(Q - 1q')C^{-1/2} = \lambda_1 U_1 V_1' + \dots + \lambda_s U_s V_s' \quad (4.10)$$

leading to the row (population) coordinates

$$(\lambda_1 R^{-1/2} U_1 : \dots : \lambda_k R^{-1/2} U_k) \quad (4.11)$$

so that no extra computations are needed if we want a representation of the row profiles also. In correspondence analysis it is customary to plot the points (4.5) and (4.11) in the same chart. Then the standardized coordinates for the columns (as variables) are

$$\lambda_1 \Delta_1^{-1} C^{1/2} V_1, \dots, \lambda_k \Delta_1^{-1} C^{1/2} V_k \quad (4.12)$$

where  $\Delta_1$  is the diagonal matrix with the  $i$ -th diagonal element as the square root of the  $i$ -th diagonal element of  $(Q - 1q')'R(Q - 1q')$ .

2. It is easy to see that

$$\begin{aligned} n(\lambda_1^2 + \dots + \lambda_k^2) &= n \text{ trace } TT', \text{ with } T \text{ as in (4.9)} \\ &= \sum \sum \frac{(n_{ij} - np_i q_j)^2}{np_i q_j} \end{aligned}$$

which is the Pearson chisquare statistic for testing independence between the attributes in a contingency table. Thus the computations involved in CA automatically allow us to test for independence, and also tests for the dimensionality of the space of profiles using statistics of the type

$$n(\lambda_i^2 + \dots + \lambda_s^2), i = 1, 2, \dots \quad (4.13)$$

as discussed in Rao (1973, pp. 556-560).

3. CA is only an exploratory data analysis to examine the configuration of row and column profiles in a general way, so that a particular convenient choice of the distance measure can serve the purpose.

On the other hand, there seem to be some drawbacks in using the chisquare distance.

1. The chisquare distance (4.8) is not a function of the  $i$ -th and  $j$ -th column profiles only. It involves the marginal profile which is a weighted average of the individual column profiles. The weights depend on the observed numbers of individuals in the column categories. These numbers may not have any relevance to the problem under study, especially when the columns represent different populations from each of which some individuals are chosen and classified according to row categories. In such a case the marginal profile depends on the actual sample sizes chosen or realized for different populations. The examples discussed in the sequel show that the derived configurations in the reduced space may be sensitive to the sample numbers.

2. The marginal profile depends on the set of populations included in CA. The CA's based on a given set of populations ( $S_1$ ) and an extended set of populations ( $S_1, S_2$ ) may provide different configurations to the subset  $S_1$ .

3. There is no particular advantage in plotting the row and column profiles in the same chart. Indeed one could use different distance measures for column and row profiles and study configurations of the column and row profiles separately.

4. Since the chisquare distance uses the marginal proportions in the denominator, undue emphasis is given to the categories with low frequencies in measuring affinities between profiles.

An alternative to the chisquare distance which has some advantages is the Hellinger Distance (HD) between the  $i$ -th and  $j$ -th column profiles defined by

$$d_{ij}^2 = (\sqrt{p_{1|i}} - \sqrt{p_{1|j}})^2 + \dots + (\sqrt{p_{s|i}} - \sqrt{p_{s|j}})^2 \quad (4.14)$$

which depends only on the  $i$ -th and  $j$ -th column profiles. In such a case, the Euclidean distance in the reduced space between the  $i$ -th and  $j$ -th column profiles is an approximation to (4.14). For the derivation of canonical coordinates of the column profiles (considered as population) we choose

$$X = \begin{pmatrix} \sqrt{p_{1|1}} & \dots & \sqrt{p_{1|m}} \\ \vdots & \dots & \vdots \\ \sqrt{p_{s|1}} & \dots & \sqrt{p_{s|m}} \end{pmatrix}$$

$$M = I, \quad W = C = \text{Diag} (n_{.1}/n, \dots, n_{.m}/n)$$

and consider the s.v.d.

$$(X - \xi 1')C^{1/2} = \lambda_1 U_1 V_1' + \dots + \lambda_s U_s V_s' \quad (4.15)$$

We may choose  $\xi' = (\xi_1, \dots, \xi_s)$  as

$$\xi_i = \sqrt{p_{i.}} = \sqrt{n_{i.}/n}, \quad \text{or} \quad (4.16)$$

$$= n^{-1}(n_{.1}\sqrt{p_{i|1}} + \dots + n_{.m}\sqrt{p_{i|m}}). \quad (4.17)$$

The canonical coordinates in  $E^k$  for the column profiles choosing  $\xi$  as in (4.16) or (4.17) are

$$\lambda_1 C^{-1/2} V_1, \lambda_2 C^{-1/2} V_2, \dots, \lambda_k C^{-1/2} V_k \quad (4.18)$$

where the components of the  $i$ -th vector are the coordinates of the  $m$  column (population) profiles in the  $i$ -th dimension. The standardized coordinates in  $E^k$  for the variables, i.e., the row categories, obtained as described in (3.15) from the same s.v.d. as in (4.15) are

$$\lambda_1 \Delta^{-1} U_1, \lambda_2 \Delta^{-1} U_2, \dots, \lambda_k \Delta^{-1} U_k \quad (4.19)$$

where  $\Delta$  is a diagonal matrix with the  $i$ -th diagonal element as the square root of the  $i$ -th diagonal element of

$$\lambda_1^2 U_1 U_1' + \dots + \lambda_s^2 U_s U_s' = (X - \xi 1')C(X - \xi 1')'. \quad (4.20)$$

The  $s$  components of  $\lambda_i \Delta^{-1} U_i$  in (4.19) are the coordinates of the variables in the  $i$ -th dimension.

It can be shown that the statistic

$$4n(\lambda_1^2 + \dots + \lambda_s^2) \quad (4.21)$$

is distributed asymptotically as chisquare on  $(s-1)(m-1)$  degrees of freedom to test independence in the two way contingency table. Further, hypotheses specifying the dimensions of the subspace in which the profiles can be represented can also be tested in the same way as in (4.13) using the residual singular values.

The advantages in using HD between profiles are the following.

1. The measure depends only on the profiles of the concerned pair. It is not altered when an extended set of profiles is considered.
2. The measure does not depend on the sample sizes on which the profiles are estimated.
3. If a representation of the row profiles is also needed we take  $X = \text{sqrt}(Q')$ , i.e., the elements of  $X$  are the square roots of the elements of  $Q'$  where  $Q$  is the matrix defined in (4.2) and compute the s.v.d.

$$(X - \eta 1')R^{1/2} = \mu_1 A_1 B'_1 + \dots + \mu_s A_s B'_s \quad (4.22)$$

leading to the canonical coordinates for row profiles

$$\mu_1 R^{-1/2} B_1, \mu_2 R^{-1/2} B_2, \dots, \mu_k R^{-1/2} B_k.$$

The corresponding standardized coordinates for the columns considered as variable are

$$\mu_1 \Delta_c^{-1} A_1, \mu_2 \Delta_c^{-1} A_2, \dots,$$

where  $\Delta_c$  is the diagonal matrix with  $i$ -th diagonal element as the square root of the  $i$ -th diagonal element of

$$\mu_1^2 A_1 A'_1 + \dots + \mu_k^2 A_s A'_s.$$

4. If we choose  $\xi$  as in (4.16), then the matrix in (4.15) is

$$(X - \xi 1')C^{1/2} = \left( \sqrt{\frac{n_{ij}}{n}} - \sqrt{\frac{n_{i.}}{n} \frac{n_{.j}}{n}} \right)$$

which is symmetric in  $i$  and  $j$ . Then, the same s.v.d. as in (4.15) could be used for computing the canonical coordinates

$$\lambda_1 R^{-1/2} U_1, \lambda_2 R^{-1/2} U_2, \dots, \lambda_k R^{-1/2} U_k$$

for the row profiles, as in the case of CA.

#### Example 4.1.

We consider the data (from Greenacre (1993)) on 796 scientific researchers classified according to their scientific discipline (as populations) and funding category (as variables) as shown in Table 3.

TABLE 3. Scientific disciplines by research funding categories

Scientific discipline		Funding category					Total
		a	b	c	d	e	
Geology	$G$	3	19	39	14	10	85
Biochemistry	$B_1$	1	2	13	1	12	29
Chemistry	$C$	6	25	49	21	29	130
Zoology	$Z$	3	15	41	35	26	120
Physics	$P$	10	22	47	9	26	114
Engineering	$E$	3	11	25	15	34	88
Microbiology	$M_1$	1	6	14	5	11	37
Botany	$B_2$	0	12	34	17	23	86
Statistics	$S$	2	5	11	4	7	29
Mathematics	$M_2$	2	11	37	8	20	78
Total		31	128	310	129	198	796

The canonical coordinates for the scientific disciplines (considered as populations) in the first three dimensions and percentage of variance explained by each are given in Table 4 for the analyses based on the chisquare distance (correspondence analysis) and the Hellinger distance (alternative). The formula (4.10) is used for the analysis based on chisquare and the formula (4.15) for that based on Hellinger distance. For Hellinger distance analysis, the central point is chosen according to the formula (4.17).

TABLE 4. Canonical coordinates for the scientific disciplines in the first three dimensions

Subjects	Chisquare Distance			Hellinger distance		
	dim1	dim2	dim3	dim1	dim2	dim3
<i>G</i>	.076401	.302569	-.087749	-.031140	.167408	-.048245
<i>B<sub>1</sub></i>	.179892	-.454996	-.151716	-.129374	-.242174	-.077614
<i>C</i>	.037644	.073353	.042371	-.021144	.040433	.028254
<i>Z</i>	-.327365	.102283	.064515	.138850	.045255	.056894
<i>P</i>	.315552	.026997	.108688	-.165340	.010679	.023844
<i>E</i>	-.117495	-.291712	.107330	.049451	-.129906	.082901
<i>M<sub>1</sub></i>	.012766	-.109656	-.041435	-.004913	-.052588	-.008439
<i>B<sub>2</sub></i>	-.178695	-.038501	-.129055	.151404	-.036559	-.108025
<i>S</i>	.124638	.014162	.107190	-.066639	.011763	.052571
<i>M<sub>2</sub></i>	.106751	-.061316	-.175688	-.050307	-.037572	-.078006
% var.	47.20	36.66	13.11	45.87	34.10	16.57

The plots of the scientific disciplines (subjects) using the canonical coordinates based on the chisquare and Hellinger distances are given in Figures 1 and 2 respectively. The coordinates in the third dimension are plotted on a line on the right hand side of the two dimensional plot. This will be of help in visualizing the plot in three dimensions and in interpreting the distances in the two dimensional plot. Thus, although *B<sub>2</sub>* and *E* appear to be close to each other in the two dimensional chart, they are clearly separated in the third dimension. No additional distances in the third dimension are involved in the case of *P*, *C*, *S*, *Z* and *E*.

It is of interest to note in this example that the configuration of the scientific disciplines in three dimensions obtained by both the methods are very similar. The percentage variance explained in each dimension is nearly the same for both the methods.

The standardized canonical coordinates for the funding categories (considered as variables) are computed using the formula (4.12) for the chisquare analysis and the formula (4.19) for the Hellinger distance analysis. These are obtained from the same s.v.d. used to compute the canonical coordinates for the scientific disciplines. Table 5 gives the standardized canonical coordinates for the funding categories, a, b, c, d, e, using the two methods.

TABLE 5. Standardized canonical coordinates for funding categories (variables) in the first three dimensions

Funding category	Chisquare Distance				Hellinger Distance			
	dim1	dim2	dim3	%var	dim1	dim2	dim3	%var
a	.758	.114	-.619	97.1	-.796	-.164	-.573	98.9
b	.535	.728	-.137	83.5	-.438	-.766	-.008	77.9
c	.583	.352	.694	94.6	-.501	-.327	.759	93.4
d	-.426	.331	-.172	99.8	.888	-.358	-.285	99.7
e	-.108	-.909	-.081	99.6	.088	.978	-.159	98.9

The standardized canonical coordinates for the funding categories are plotted in Figure 3 (for chisquare distance) and in Figure 4 (for Hellinger distance). It may be noted that all the points lie within the unit circle. It is customary to represent the canonical coordinates for the subjects and variables in one chart. We are using separate charts in order to explain the salient features of the configuration of the variables. The following interpretations emerge from the study of Table 5 and Figures 3 and 4.

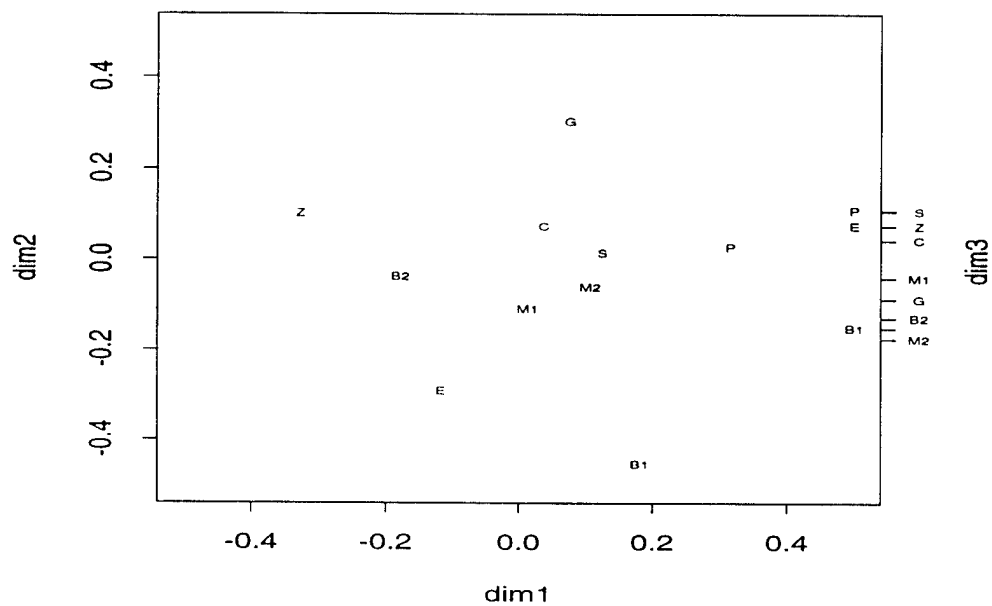


Figure 1. Configuration of scientific disciplines  
using Chisquare distance  
(Correspondence Analysis)

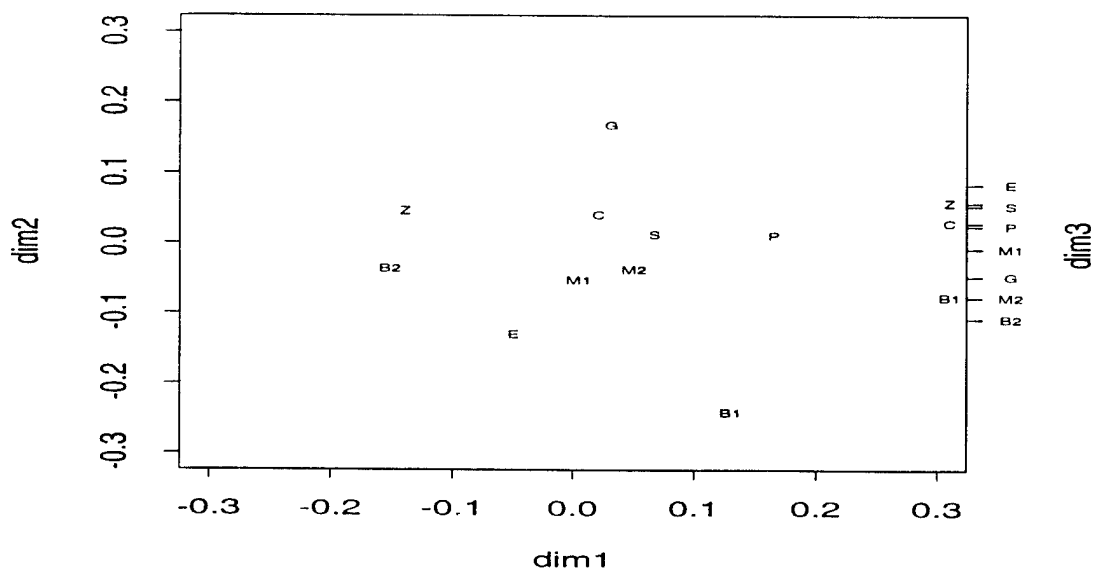


Figure 2. Configuration of scientific disciplines  
using Hellinger distance  
(Alternative to Correspondence Analysis)

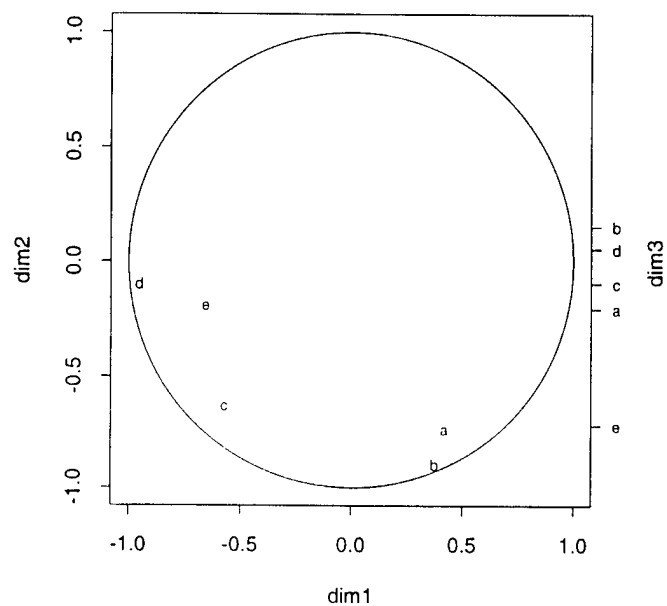


Figure 3. Configuration of funding categories using standardized canonical coordinates based on Chisquare distance

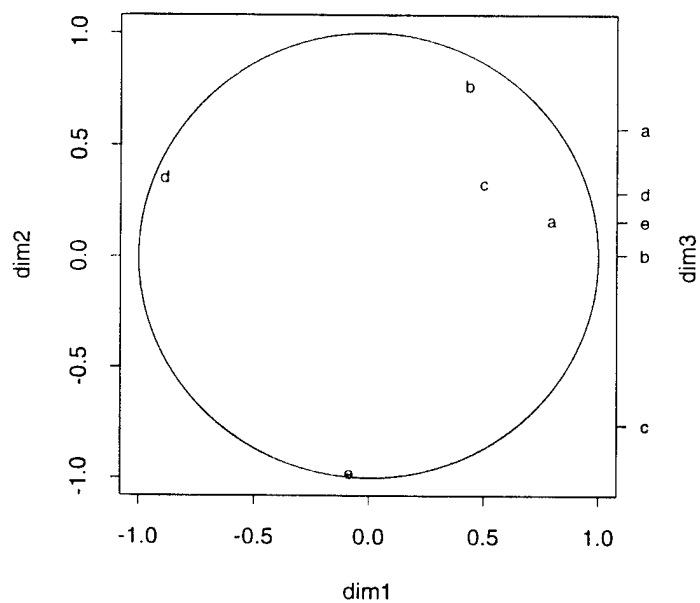


Figure 4. Configuration of funding categories using standardized canonical coordinates based on Hellinger distance



1. The configurations of the funding categories as shown in Figures 3 and 4 obtained by using chisquare and Hellinger distances are very similar.

This is not generally the case, although in some examples studied by the author a good deal of robustness was observed in the choice of the distance measure and relative sample sizes for the populations under study. However, addition or deletion of some populations may affect the configuration of the populations when correspondence analysis is used. Example 4.2 discussed below throws more light on the problem and shows that the analysis based on Hellinger's distance is more robust to relative sample sizes.

2. All most all the variation in the funding categories  $a, d$  and  $e$  is captured in the first three canonical coordinates of the scientific disciplines. A large percentage of variation in  $b$  and  $c$  is explained by the first three coordinates.
3. The first dimension is strongly influenced by  $a, d$ , the second dimension by  $b, c$  and the third dimension by  $a, c$ .

Thus the use of standardized coordinates for variables enables us to interpret the different dimensions in terms of observed variables. There are other ways of plotting the coordinates of the variables as mentioned in the paragraphs below Table 2. Such biplots having a different interpretation are discussed in Gabriel (1971), Gifi (1990), Gower (1993) and Greenacre (1993).

*Note 4.2.* In computing the canonical coordinates based on Hellinger distance (HD) using the formula (4.15), we chose the relative sample sizes as the weights to be attached to the populations. We could have chosen an alternative set of weights if we wanted distances between a specified subset of populations to be better preserved in the reduced space than the others. In particular, we could have chosen uniform weights for all populations. In fact such an option could be exercised if the sample sizes of different populations were widely different. Unfortunately no such options are available in correspondence analysis.

### Example 4.2.

In the example 4.1, there was a perfect match between the plots based on CA and HD. This probably demonstrates that the method of derivation of canonical coordinates is somewhat robust to the choice of the distance measure as well as to the weights. However the choice of HD provides an insurance against possible distortion due to variations in sample sizes for the populations as the following example shows.

Table 6, reproduced from Gifi (1990), gives the distributions of the pages devoted to different topics denoted by  $A, B, C, D, E, F$  and  $G$  in 20 books on Multivariate analysis designated as  $a, b, \dots, t$ . Gifi (1990) did correspondence analysis on the data and drew some conclusions based on the first three canonical coordinates which explain a high percentage of variation. The first three canonical coordinates for the profiles of the books based on CA and HD approaches are given in Table 7.

It may be noted that the total number of pages of a book depends on the font size of the print, while its profile in terms of proportions of pages used on different topics remain the same for all sizes. Table 8 gives the data on books having the same profiles as in Table 6 with the total number of pages altered for the books  $d, f, g, h, j$  and  $n$ .

The three dimensional canonical coordinates based on CA and HD approaches are given in Table 9. Using the coordinates one can obtain the mutual distances between the books in the three dimensional reduced Euclidean space. Figure 5 gives a plot comparing the squared distances between books based on CA using the data of Tables 6 and 8. Figure 6 gives the corresponding plot for the squared distances based on the HD approach. It is seen that the three dimensional representation of the data of Tables 6 and 8 are more similar under HD analysis than that under CA. The relative positions of the books are influenced by the font size in printing when CA is used, although the profiles of the books are not altered. There appears to be greater stability with the HD analysis which provides insurance against different choice of sample sizes. Further, one can exercise the option of using a common weight for all the books in the HD analysis when the differences in book sizes are large.

TABLE 6. Number of pages by topics

Books	A	B	C	D	E	F	G
a	31	0	0	0	0	164	11
b	0	16	54	18	27	13	14
c	0	40	32	10	42	60	0
d	19	0	35	19	28	163	52
e	14	7	35	22	17	0	56
f	20	69	72	33	55	0	32
g	74	0	86	14	0	84	48
h	78	0	80	5	17	105	60
i	74	19	33	12	26	0	0
j	80	68	67	15	29	0	0
k	108	48	4	10	46	108	0
l	109	13	5	17	39	32	46
m	16	35	69	24	0	26	41
n	26	86	60	6	48	48	28
o	290	10	6	0	8	0	2
p	184	48	82	42	134	0	0
q	29	0	0	0	41	211	32
r	0	19	56	0	39	75	0
s	0	22	45	42	60	230	59
t	30	128	90	28	48	0	0

TABLE 7. Canonical coordinates

	Chisquare Distance			Hellinger Distance		
	dim 1	dim 2	dim 3	dim 1	dim 2	dim 3
a	-1.10857	-0.61445	-0.33902	0.64632	0.36299	0.12879
b	0.07397	0.70254	0.25265	-0.01661	-0.48923	-0.12388
c	-0.21153	0.46054	-0.49228	0.10998	-0.42185	0.32822
d	-0.77795	-0.11074	0.15556	0.46658	-0.01597	-0.10284
e	0.02781	0.40651	1.06135	-0.19193	-0.15180	-0.45570
f	0.35780	0.69602	0.09284	-0.37016	-0.29359	-0.14451
g	-0.16412	-0.15719	0.46353	0.23979	0.16911	-0.35829
h	-0.25023	-0.19626	0.39002	0.26103	0.14730	-0.23804
i	0.72788	-0.19452	-0.04749	-0.50899	0.14292	0.02936
j	0.68403	0.24337	-0.17956	-0.53320	-0.01242	0.04724
k	0.02729	-0.36648	-0.44297	0.03996	0.21098	0.36189
l	0.26802	-0.44749	0.28287	-0.00524	0.27070	-0.06481
m	0.02188	0.50893	0.51719	0.01506	-0.19266	-0.34080
n	0.12052	0.48459	-0.19476	-0.04555	-0.20966	0.04945
o	1.08308	-1.32602	0.03206	-0.39476	0.66357	-0.00090
p	0.64959	-0.07081	-0.13268	-0.49299	0.09097	0.08510
q	-0.98347	-0.39273	-0.25019	0.58910	0.21379	0.19442
r	-0.40006	0.32919	-0.33826	0.21605	-0.35929	0.28764
s	-0.74726	0.08101	-0.00508	0.43349	-0.30134	0.03139
t	0.56547	0.81454	-0.35256	-0.51167	-0.27874	0.10162

TABLE 8. Number of pages by topics

Books	A	B	C	D	E	F	G
a	31	0	0	0	0	164	11
b	0	16	54	18	27	13	14
c	0	40	32	10	42	60	0
d	190	0	350	190	280	1630	520
e	14	7	35	22	17	0	56
f	10	34	36	17	28	0	16
g	740	0	860	140	0	840	480
h	780	0	800	50	170	1050	600
i	74	19	33	12	26	0	0
j	40	34	33	8	15	0	0
k	108	48	4	10	46	108	0
l	109	13	5	17	39	32	46
m	16	35	69	24	0	26	41
n	13	43	30	3	24	24	14
o	290	10	6	0	8	0	2
p	184	48	82	42	134	0	0
q	29	0	0	0	41	211	32
r	0	19	56	0	39	75	0
s	0	22	45	42	60	230	59
t	30	128	90	28	48	0	0

TABLE 9. Canonical coordinates

	Chisquare Distance			Hellinger Distance		
	dim 1	dim 2	dim 3	dim 1	dim 2	dim 3
a	-0.62310	0.30413	-0.53463	-0.35082	-0.04632	0.42925
b	0.63345	0.41500	0.44316	0.25096	-0.37625	-0.40565
c	0.90486	0.78379	-0.17802	0.22985	-0.60374	-0.08540
d	-0.36427	0.36470	-0.12611	-0.23454	-0.16742	0.01739
e	0.20621	0.05647	0.54414	0.34853	0.01585	-0.32928
f	1.23299	0.45214	0.27035	0.59591	-0.20402	-0.28683
g	-0.16974	-0.27626	0.25537	-0.08445	0.24917	-0.09573
h	-0.18352	-0.18729	0.09783	-0.07908	0.11818	0.00148
i	0.85607	-0.49586	-0.21672	0.73058	0.07206	0.07026
j	1.35943	-0.06808	0.07459	0.77422	-0.01788	-0.04783
k	0.61122	0.01365	-0.60327	0.28646	-0.20830	0.40764
l	0.32350	-0.41447	-0.39537	0.23929	0.02213	0.24724
m	0.58680	0.20860	0.65792	0.17707	0.01371	-0.36016
n	1.20448	0.51816	0.07444	0.32243	-0.27238	-0.09222
o	0.47616	-1.60929	-0.73075	0.61206	0.41339	0.46017
p	0.87199	-0.26044	-0.37985	0.72806	-0.03491	0.09187
q	-0.44497	0.44631	-0.57047	-0.27936	-0.25639	0.41829
r	0.34209	0.56913	-0.04981	0.10278	-0.49844	-0.07991
s	-0.10036	0.60181	-0.15749	-0.14779	-0.45659	-0.10069
t	1.87687	0.57376	0.28038	0.78353	-0.24748	-0.19823

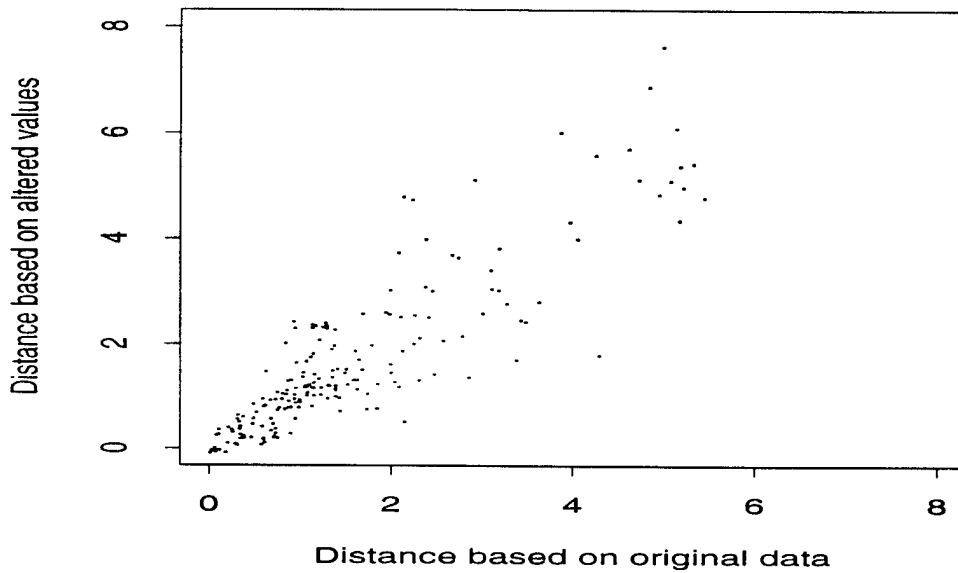


Figure 5. A comparative plot of squared distances between all pairs of books in the reduced spaces based on correspondence analysis

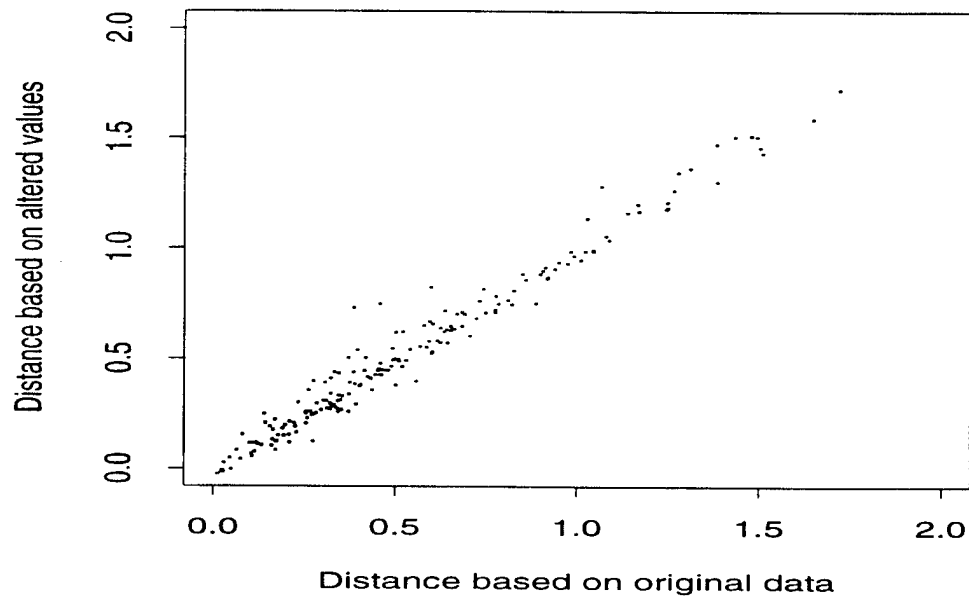


Figure 6. A comparative plot of squared distances between all pairs of books in the reduced spaces based on Hellinger distance analysis

### 5. Concluding remarks

A general theory is developed for plotting high dimensional "population by variable" data, i.e. measurements made on a set of characteristics of given populations, in a low dimensional Euclidean space. A first step in such a problem is the specification of the basic metric space in which the populations can be represented as points using the entire data, and a characterization of the configuration of the points in terms of distances between points. The second is the development of methodology for transforming the points from the basic space to a low dimensional Euclidean space with the usual definition of distance preserving the configuration of points to the extent possible. The choices of the basic space and the distance function between points have to be made on practical considerations depending on the problem under investigation. A closed form solution is obtained when the basic space is a vector space endowed with an inner product and the associated norm. Some examples are given involving measurements on discrete variables.

When we have data in the form of frequencies of individuals of a population under different categories of an attribute, a well known method for dimensionality reduction for representing, say the populations, is correspondence analysis. The basic space in this case is a vector space where each population is represented by the vector of relative frequencies of the different categories of an attribute and distance between vectors is defined by a chisquare type formula. Such a distance function is not an intrinsic measure of difference between two populations as it depends not only on the differences between their relative frequencies, but also on the average relative frequencies computed from the set of populations under study. Thus the configuration of any subset of populations depends on what other populations are included in the analysis, and also on the relative numbers of individuals observed from each population. An alternative approach of representing a population by the vector of the square roots of relative frequencies and defining distance between two populations by the Hellinger formula does not have the drawbacks associated with the chisquare type formula. In addition, the new analysis has the same advantage of providing tests of significance for homogeneity of the populations as in correspondence analysis based on the chisquare formula.

It may be contended that CA is meant to be used for the analysis of contingency tables with dichotomized data using two attributes like hair color and eye color (as originally demonstrated by R.A. Fisher), and not for the analysis of population by variable data where anomalies of the type described in the paper may occur. However, one finds in published literature more examples of the latter type of data analyzed through CA. Further, even with attribute data, if the configurations of the column (or row) profiles for two different populations (with possibly different marginal distributions) are to be compared, HD analysis is more appropriate than the CA. It is the author's opinion that the choice of a distance measure between populations (row or column profiles) must depend on the nature of the data and the purpose of analysis. Prescription to use a particular distance as in the CA in all problems may be misleading. Distance measures other than the chisquare and Hellinger types may be more appropriate in some situations. For a purely exploratory data analysis, it is possible that a wide variety of distance measures reveal similar configurations of the populations in terms of clustering and inter cluster relationships.

Between the choices of chisquare and Hellinger distances, the latter seems to offer some advantages, as the latter has similar theoretical properties as the former and in addition it is defined as an intrinsic function of two population profiles independent of what other populations are included in a study.

A recent technical report by Rios, Villarroya and Oller (1994) discusses the same problem as in the present paper, viz., simultaneous representation of populations and random variables, under the assumption of an underlying parametric model.

The method, referred there as *Intrinsic Data Analysis*, is based on the Riemannian structure given by the Fisher information metric and its corresponding distance, the *Rao distance*. The statistical populations are viewed as points on a Riemannian manifold and the random variables with finite expectation, as vector fields, namely, the gradient of the random variable mean value, or, by integration, a bundle of curves on the manifold.

Then, assuming certain additional regularity conditions, a reference point on the manifold is selected as the statistical populations Riemannian center of mass, and the points representing the populations and the curves representing the variables are mapped, through the inverse of the Riemannian exponential map, into the tangent space at the center of mass, which has a Euclidean vector space structure. Then, classical dimension reduction techniques such as principal component analysis can be used to obtain a

low dimensional Euclidean space which allows an optimal population representation. Finally, the curves in the tangent space are projected into the low dimensional space obtained.

This method is applied to multivariate normal and multinomial distributions. In the multinomial case, the Rao distance,  $\rho$ , between two populations  $p_1, \dots, p_n$  and  $q_1, \dots, q_n$ , is proportional to the Bhattacharyya distance

$$\rho = 2 \arccos \sum_{j=1}^n \sqrt{p_j q_j}$$

which is a monotone transformation of the Hellinger distance, and thus this method will share some properties with the latter.

### References

1. J.P. Benzécri, *Correspondence Analysis Handbook*, Marcel Dekkar, Inc., New York, (1992).
2. L.L. Cavallis-Sforza, *Genes, peoples and languages*, Scientific American, **265** (1991), 104-110.
3. H. Chernoff, *The use of faces to represent points in k-dimensional space graphically*, J. Amer. Statist. Assoc. **68** (1973), 361-368.
4. C. Eckart and G. Young, *The approximation of one matrix by another of lower rank*, Psychometrika **1** (1936), 211-308.
5. G. Eslava-Gomez and F.H.C. Marriott, *Criteria to represent groups in the plane when the grouping is unknown*, Biometrics **49** (1993), 1088-1098.
6. K.R. Gabriel, *The biplot graphical display of matrices with applications to principal component analysis*, Biometrika, **58** (1971), 453-467.
7. A. Gifi, *Nonlinear Multivariate Analysis*, New York: John Wiley (1981, 1990).
8. J.C. Gower, *Recent advances in biplot methodology*, In Multivariate Analysis: Future Directions 2 (Eds. C.M. Cuadras and C.R. Rao), North Holland (1993), 295-325.
9. M.J. Greenacre, *Theory and Applications of Correspondence Analysis*, London: Academic Press (1984).
10. ———, *Biplot in correspondence analysis*, J. Applied Statistics **20** (1993), 251-269.
11. J.B. Kruskal and M. Wish, *Multidimensional Scaling*. Sage Publications (1978).
12. P.C. Mahalanobis, *On the generalized distance in statistics*, Proc. Nat. Inst. Sci., India **12** (1936), 49-55.
13. P.C. Mahalanobis, D.N. Mazumdar. and C.R. Rao, *Anthropometric survey of United Provinces, 1941. A statistical study*, Sankhya **9** (1949), 90-324.
14. S. Nishisato, *Analysis of Categorical Data: Dual Scaling and its Applications*, University of Toronto Press, Toronto, Canada (1980).
15. C.R. Rao, *Information and accuracy attainable in the estimation of statistical parameters*, Bull. Cal. Math. Soc. **37** (1945), 81-91.
16. ———, *The problem of classification and distance between two populations*, Nature **159** (1947), 30.
17. ———, *The utilization of multiple measurements in problems of biological classification (with discussion)*, J. Roy. Statist. Soc. Series B10 (1948), 159-193.
18. ———, *Some statistical methods for comparison of growth curves*, Biometrics **14** (1958), 1-17.
19. ———, *The use and interpretation of principal component analysis in applied research*, Sankhya **26** (1964), 329-357.
20. ———, *Linear Statistical Inference and its Applications*, 2nd Edition, New York : Wiley (1973).
21. ———, *Separation theorems for singular values of matrices and their applications in multivariate analysis*, J. Multivariate Analysis **9** (1979), 362-377.
22. ———, *Matrix approximations and reduction of dimensionality in multivariate statistical analysis*, In Multivariate Analysis V (Ed. P.R. Krishnaiah), Amsterdam: North Holland, (1980), 3-22.
23. ———, *Tests for dimensionality and interaction of mean vectors under general and reducible covariance structures*, J. Multivariate Analysis **16** (1985), 173-184.
24. ———, *Prediction of future observations in growth curve type models*, J. Statistical Science **2** (1987), 434-471.
25. M. Rios, A. Villarroya and J.M. Oller, *Intrinsic Data Analysis: A method for the simultaneous representation of populations and variables*, Mathematics preprint series **160** (1994), Universitat de Barcelona.
26. J. Von Neumann, *Some matrix inequalities and metrization of metric spaces*, Tomsk. Univ. Rev. **1** (1937), 286-299.
27. E.J. Wegman, *Hyperdimensional data analysis using parallel coordinates*, J. Amer. Statist. Assoc. **85** (1990), 664-675.

E-mail address: crr1@psuvm.psu.edu

CENTER FOR MULTIVARIATE ANALYSIS, DEPARTMENT OF STATISTICS, PENNSYLVANIA STATE UNIVERSITY, UNIVERSITY PARK, PA 16802, USA