

**MULTI-SET NONLINEAR CANONICAL CORRELATION ANALYSIS
VIA THE BURT-MATRIX**

Robert Tijssen

Department of Data Theory

University of Leiden

Jan de Leeuw

Departments of Psychology and Mathematics

University of California, Los Angeles

MULTI-SET NONLINEAR CANONICAL CORRELATION ANALYSIS VIA THE BURT-MATRIX

Robert Tijssen*, Jan de Leeuw**

*Department of Data Theory, Middelstegegracht 4, 2312 TW Leiden,
University of Leiden, The Netherlands,

**Departments of Psychology and Mathematics, University of California - Los Angeles,
405 Hilgard Avenue, Los Angeles CA 90024, USA.

Canonical Correlation Analysis (CCA) is a data analysis method in which the correlation between two sets (of linear combinations) of categorical variables is maximized. CCA is discussed in geometrical and matrix-algebraic terms as an introduction to a multi-set CCA with optimal scaling properties of category values. Instead of using a rectangular (objects x variables)-data matrix, this technique operates on the square Burt-matrix, which contains the whole of bivariate relations. Basics of the underlying statistical theory are discussed. A real-life application is presented.

1. INTRODUCTION

Canonical Analysis can be regarded as a generalization of regression theory and is found in many disguises amongst multivariate analysis methods (cf. Gittens, 1984). *Canonical Correlation Analysis* (CCA) is a CA method offering the product-moment correlation as a measure of resemblance between sets of categorical variables. This measure is derived from both the between and within-set relations. Linear CCA was introduced in classical multivariate statistical analysis by Hotelling (1936). The linearity of classical CCA refers to the fact that its results are invariant under linear transformations of the category values. In the years to follow, Classical CCA has become an established multivariate statistical method for relating two sets of variables (e.g., Thomson, 1947). Multi-set extensions of CCA have been developed and successfully applied during the last decades (cf. Kettenring, 1971). More recently, nonlinear generalizations of multi-set CCA have also been introduced (Van der Burg et al., 1988).

2. LINEAR TWO-SETS CCA

The basic strategy of CCA can be described as finding an optimal linear combination -a weighted sum- of the variables in each set, in such a way that the correlation between each set of variables is maximized. A $(n \times m)$ data matrix \mathbf{H} , with measurements of a number of objects or persons i ($i=1, \dots, n$) on discretely-valued categorical variables j ($j=1, \dots, l, \dots, m$) with r_j ($r=1, \dots, k_j$) categories, is partitioned in k ($k=1, 2$) sets of variables. The first $(n \times m_1)$ -sized set of variables will be denoted by a m_1 -dimensional vector $\mathbf{h}_1=(\mathbf{h}_{11}, \dots, \mathbf{h}_{1m_1})$, the second $(n \times m_2)$ -sized set by a m_2 -dimensional vector $\mathbf{h}_2=(\mathbf{h}_{21}, \dots, \mathbf{h}_{2m_2})$. The space \mathbb{L} spanned by the $m=m_1+m_2$ variables, with dimensionality $s = (1, \dots, p)$ where $p \leq \min.(m_1, m_2)$, can be partitioned in such a way that the variables in each set span a corresponding linear subspace \mathbb{L}_1 and \mathbb{L}_2 . The aim of CCA can be reformulated as finding directions in \mathbb{L}_1 and \mathbb{L}_2 , the *canonical axes* or *canonical variates*, which are as similar as possible. We therefore

need measure to indicate the 'goodness-of-fit'. The cosine of the angle between the *canonical variates*, the *canonical correlation*, is used for this purpose. Linear combinations of the variables in each set are thus formed by the *canonical weights* in the ($m_1 \times s$) matrix A_1 and the ($m_2 \times s$) matrix A_2 . Conceptually, the general CCA problem of relating two sets of variables can thus be seen as finding weight matrices A and weighted columns

$$h_1 a^s_1 = h_{11} a^{s1}_1 + h_{12} a^{s1}_2 + \dots + h_{m1} a^{s1}_{m1} \quad (1a)$$

$$h_2 a^s_2 = h_{21} a^{s2}_1 + h_{22} a^{s2}_2 + \dots + h_{m2} a^{s2}_{m2} \quad (1b)$$

which represent an orthogonal basis for the subspaces L_1 and L_2 of dimensionality $s < p$ which yield an optimal canonical correlation for each dimension.

Shifting the focus from a geometrical point-of-view to an interpretation CCA in matrix-algebraic terms, boils down to finding maximally related axes in subspaces L_1 and L_2 of the common space L , by means of *eigenvalue-eigenvector decomposition* (Wilkinson, 1965) of the partitioned datamatrix $H = (H_1 | H_2)$. Since the canonical correlation is invariant under linear scaling of h_j we will require that h_j is centered with unit-variance (i.e., $u'h_j = 0$, $h_j'h_j = 1$; u denotes a unity-vector). Solutions for the canonical weights can be found in such a way that the correlations between the weighted combinations $H_1 A_1$ and $H_2 A_2$ are optimal. In terms of the sample-based correlation matrix $R_{1,2} = H_1' H_2$, the problem of finding the optimal set of canonical weights is equivalent to solving the following canonical equations (Anderson, 1958):

$$[(R_{1,1})^{-1} R_{1,2} (R_{2,2})^{-1} R_{2,1} - \lambda^2 I] A_1 = 0, \quad (2a)$$

$$[(R_{2,2})^{-1} R_{2,1} (R_{1,1})^{-1} R_{1,2} - \lambda^2 I] A_2 = 0, \quad (2b)$$

where I denotes the identity matrix and λ^2 represents the first and largest eigenvalue of the characteristic equations

$$|(R_{1,1})^{-1} R_{1,2} (R_{2,2})^{-1} R_{2,1} - \lambda^2 I| = 0, \quad (3a)$$

$$|(R_{2,2})^{-1} R_{2,1} (R_{1,1})^{-1} R_{1,2} - \lambda^2 I| = 0. \quad (3b)$$

The largest eigenvalue, either of the product matrix $(R_{1,1})^{-1} R_{1,2} (R_{2,2})^{-1} R_{2,1}$ or of the matrix $(R_{2,2})^{-1} R_{2,1} (R_{1,1})^{-1} R_{1,2}$, is now equal to the highest squared canonical correlation in the first dimension. The accompanying first pair of eigenvectors a_1 and a_2 of the n -orthonormal basis, will

yield the highest canonical correlation between all possible linear combinations of weighted variables within sets. Since the dimension-wise weight vectors \mathbf{a}_1 and \mathbf{a}_2 are interchangeable, via $\mathbf{a}_1=[(\mathbf{R}_{1,1})^{-1}\mathbf{R}_{1,2}\mathbf{x}_2]/\lambda$ and $\mathbf{a}_2=[(\mathbf{R}_{2,2})^{-1}\mathbf{R}_{2,1}\mathbf{a}_1]/\lambda$, only one characteristic equation needs to be solved. The second pair of canonical variates (i.e., \mathbf{a}_2^1 and \mathbf{a}_2^2) are uncorrelated with the first pair, etc. The non-paired canonical variates are also mutually uncorrelated between sets.

3. NONLINEAR CCA

A generalization of CCA to enable the incorporation of variables with less restrictive (nonlinear) measurement levels - generally classified either as *nominal* or *ordinal*, dates back to Young et al. (1976). A nonlinear version of two-sets CCA is also proposed by Van der Burg and De Leeuw (1983), with an accompanying alternating least squares program CANALS, which yields CCA-results which are invariant under certain nonlinear transformations of the variables. One must, however, keep in mind that nonlinearly transformed variables do again define a linear space \mathbb{L} . This type of *optimal scaling* (Young, 1981) of category values, i.e. rescaled according to the constraints of the respective measurement levels, leads to category values $y_j = t_j(r_j)$ referred to as *category quantifications* (Gifi, 1981). In case of a nonrestricted nominal variable one can visualize the set of possible values y_j as a k_j -dimensional space \mathbb{S}_j . The set of possible transformations of an ordinal variable j defines a regression problem of the datavector \mathbf{h}_j on a polyhedral convex cone \mathbb{K}_j in the k_j -space. Numerical (linear) variables define a regression on a one-dimensional subset of \mathbb{S}_j .

CCA now consists of two computational subproblems which are dependent upon each other and can be solved by means of an algorithm which maximizes the canonical correlation, while simultaneously imposing the proper measurement restrictions on variables. Assume a $(n \times k_j)$ binary *indicator matrix* \mathbf{G}_j as a basis for each \mathbb{S}_j , with $\mathbf{G}_j\mathbf{u}=\mathbf{u}$, and $g_{jk}=1$ if $\mathbf{h}_{ij}=k$; $g_{jk}=0$ if $\mathbf{h}_{ij}\neq k$ (cf. De Leeuw, 1984). Thus each category defines a binary variable and each individual is represented in only one category per variable. The indicator matrices of the variables j thus contain k_j independent columns and span an orthogonal basis for each variable. The expression $\mathbf{q}_j=\mathbf{G}_j\mathbf{y}_j$ defines a transformed variable (Gifi 1981). Unit-normalized vectors \mathbf{q}_j define a correlation matrix $\mathbf{R}(\mathbf{Q})=\mathbf{R}(\mathbf{q}_1,\dots,\mathbf{q}_m)$, with elements $r_{jl}=\mathbf{q}_j'\mathbf{q}_l$. An induced correlation based on the bivariate matrix \mathbf{C}_{jl} can now be defined as $\mathbf{R}(\mathbf{q}_1,\dots,\mathbf{q}_m)=\mathbf{y}_j'\mathbf{C}_{jl}\mathbf{y}_l$ with $\mathbf{C}_{jl}=\mathbf{G}_j'\mathbf{G}_l$; a diagonal matrix with univariate marginals \mathbf{D}_{jl} and $\mathbf{y}_j'\mathbf{D}_{jl}\mathbf{y}_l=1$; $\mathbf{u}'\mathbf{D}_j\mathbf{y}_j=0$.

Nonlinear two-sets CCA can now be rephrased as a technique which computes the optimally scaled variables \mathbf{q}_j , and canonical variates $\mathbf{Q}_k\mathbf{A}_k$ maximizing an optimality-criterion function $\phi_p(\mathbf{R}(\mathbf{Q}_k\mathbf{A}_k))$, or at least computes a stationary value of this function. In general, we are searching for all stationary values of $\phi_p(\mathbf{R}(\mathbf{q}_1,\dots,\mathbf{q}_m))$ on \mathbb{S}_j .

4. MULTI-SET NONLINEAR CCA

The generalization of nonlinear two-sets CCA to more than two sets $k(k=1, \dots, K)$ is - essentially - quite straightforward. Applying this extension to the generalized canonical solution in terms of matrices $\mathbf{Q}=(\mathbf{Q}_1 | \dots | \mathbf{Q}_K)$, results in a geometrical solution with eigenvectors in \mathbf{A}_k "bundling" around a *mean canonical variate*, equal to the sum vector \mathbf{QA} . The mean canonical variate thus contains the overall information of the canonical variates of each canonical solution and forms a orthogonal basis for the common canonical space \mathbb{L} , providing similar interpretations with respect to canonical correlations between the various sets. In K -sets CCA the criteria ϕ_p are also based on the maximization or minimization of properties of the eigenvalues of the correlation matrix $\mathbf{R}(\mathbf{Q}_k \mathbf{A}_k)$. Using more than two sets provides us with opportunities to extend the possible ways of computing CCA-solutions. The canonical correlations can now be generalized to measures of relatedness, which are derived under various conditions and optimality criteria, each combination producing somewhat different results (cf. Kettenring, 1971; Gifi, 1981; Van de Geer, 1984; Meulman, 1986). Van der Burg et al. (1988) introduce an alternating least squares multi-set program OVERALS, as a particular generalization of nonlinear two-sets CCA.

5. MULTI-SET NONLINEAR CCA BASED ON BIVARIATE CROSSTABLES

In some cases one might not have or want to analyze the usual $(n \times m)$ matrix \mathbf{H} , but instead an aggregated $(m \times k \times m \times k)$ Burt matrix \mathbf{C} , with crosstables \mathbf{C}_{jl} between all variables (Burt, 1950). In the sequel it will be shown that it is still possible to devise a technique which can perform a K -sets nonlinear CCA based on \mathbf{C} via $\mathbf{R}(\mathbf{Q}_k \mathbf{A}_k)$ (De Leeuw, 1983; Tijssen, 1985).

It has been shown that a linear CCA-solution is computed on \mathbf{R} . In case of a nonlinear CCA based on \mathbf{C} certain requirements must be fulfilled to obtain an optimal CCA-solution, because the induced correlations are only optimal association measures in case of linear and homoscedastical regression between variables. CCA-solutions based on \mathbf{R} , given the bivariate relations in \mathbf{C} , which also handle nonlinear variables, thus need an optimal approximation of the linearity property of correlations. Maximizing the largest eigenvalues of an \mathbf{R} , induced from \mathbf{C} via category quantifications \mathbf{y} , is such a method. The resulting \mathbf{R} will now be as one-dimensional (linear) as possible. Suppose we want to optimize a function of the correlation matrix \mathbf{R} , written as $\phi(\mathbf{R})$. An example would be the largest eigenvalue of \mathbf{R} , which leads to Multiple Correspondence Analysis, or the sum of the p largest eigenvalues, which leads to Principal Components Analysis. Compare De Leeuw (1986) for a more complete discussion of the statistical concepts behind this data-analytic approach.

In the formalization of this approach, we use the fact that $\mathbf{r}_{jl} = \mathbf{y}'_j \mathbf{C}_{jl} \mathbf{y}_l$, if $\mathbf{y}'_j \mathbf{D}_{jl} \mathbf{y}_l = 1$. The stationary equations for this optimization problem are

$$\sum_l \delta\phi / \delta \mathbf{r}_{jl} \mathbf{C}_{jl} \mathbf{y}_l = \lambda_j \mathbf{D}_j \mathbf{y}_j, \quad (4)$$

where δ denotes the derivative.

If category values y_j can be found such that the equality $C_{jl} y_l = r_{jl} D_j y_j$ holds for all variables j , then the y_j 's linearize the bivariate regressions. Hence the stationary equations are satisfied no matter how the function $\phi(\mathbf{R})$ is defined.

Studying a K -sets CCA-problem requires a partitioning of \mathbf{R} according to the allocation of the variables over the sets. Such a division can be accomplished by creating a conveniently arranged $(m \times m)$ binary matrix \mathbf{E} , with elements $e_{jl} = 1$ if variables belong to the same set and $e_{jl} = 0$, otherwise. The Hadamard product (Styan, 1973) of \mathbf{E} and \mathbf{R} results in a $(m \times m)$ matrix - denoted as $\mathbf{E} \wedge \mathbf{R}$, with $r_{jl} = e_{jl} \wedge r_{jl}$ if the corresponding variables belong to the same set. The nonlinear K -sets CCA-problem is to maximize the sum of the first p eigenvalues of the pair $(\mathbf{R}, \mathbf{E} \wedge \mathbf{R})$. This is of the form $\phi(\mathbf{R})$. In order to find the eigenvalues one has to solve the *generalized eigenvalue problem* :

$$\mathbf{R}\mathbf{A} = \lambda^2 (\mathbf{E} \wedge \mathbf{R})\mathbf{A}, \quad (5)$$

with the $(p \times p)$ diagonal matrix with generalized λ^2 -values. The result of this particular eigenvalue problem will provide us with a common basis \mathbf{L} for the sets based on eigenvectors \mathbf{A} which are now orthonormal with respect to the within-sets correlations, i.e. $\mathbf{A}'(\mathbf{E} \wedge \mathbf{R})\mathbf{A} = \mathbf{I}$.

In this paper we shall study a somewhat more general class of criteria ϕ which depend on \mathbf{R} through the generalized eigenvalues λ^2 . One could also use the product of the generalized values, for instance, or the sum of squares of their deviations from unity.

By applying the chain rule, it is well known that

$$\delta \lambda^2_s / \delta r_{jl} = [1 - (\lambda^2_s e_{jl})] a_{js} a_{ls}. \quad (6)$$

Combining Eq. (4) and Eq. (6) provides us with two $(\sum_j k_j \times \sum_j k_j)$ matrices for each pair of variables:

$$\mathbf{T}_{jl} = \sum_s \delta \phi / \delta \lambda^2_s (\lambda^2_s C_{jl} a_{js} a_{ls} e_{jl}), \quad (7a)$$

$$\mathbf{U}_{jl} = \sum_s \delta \phi / \delta \lambda^2_s (C_{jl} a_{js} a_{ls}). \quad (7b)$$

Hence \mathbf{T} is based on the within-sets bivariate relations, whereas \mathbf{U} is based on both the within and between sets bivariate relations.

We can now define stationary equations as

$$\mathbf{T}\mathbf{y} - \mathbf{U}\mathbf{y} = \mathbf{0}, \quad (8)$$

which suggests an iteration scheme

$$\mathbf{y}^* = \mathbf{y}^{(n+1)} = \mathbf{U} + \mathbf{T} \mathbf{y}^{(n)}, \quad (9)$$

with the superscript + denoting the Moore-Penrose inverse, in which alternations between successive steps (n) and ($n+1$) will lead to category quantifications \mathbf{y}^* with (approximate) linearizing properties. The quantifications are normalized in each iteration step, thus inducing an updated optimal correlation matrix \mathbf{R}^* in each step via

$$\mathbf{R}^* = \{\mathbf{r}^*_{jl}\} = \mathbf{y}^*_{j'} \mathbf{C}_{jl} \mathbf{y}^*_{l}. \quad (10)$$

The resulting quantifications \mathbf{y}^* are thus an optimal function of the bivariate within-sets structure in combination with the total bivariate structure. Categories with a relatively high bivariate frequency will tend to have, on the whole, more similar quantifications.

The iteration process therefore in fact consists of two subprocesses. The inner-iterations [Eq. (9)] produce optimal quantifications \mathbf{y}^* with linearizing properties, which are used in the outer-iteration steps [Eqs. (5), (10)] to compute the corresponding \mathbf{R}^* . The generalized eigenvalue-eigenvector decomposition of $(\mathbf{R}^*, \mathbf{E} \wedge \mathbf{R}^*)$ subsequently computes the corresponding eigenvectors \mathbf{A} for the K -sets CCA solution. The elements of the Burt matrix are weighted with the corresponding elements in \mathbf{A} , creating updates of \mathbf{U} and \mathbf{T} , which provide a new update of \mathbf{y}^* . These quantifications are then used to induce a new \mathbf{R}^* , etcetera ... until convergence is reached. Assuming this iteration process leads to a stable value of the function $\phi_p(\mathbf{R})$ with a corresponding optimally linearized correlation matrix \mathbf{R}^* , one obtains optimal category quantifications in the sense that they induce an optimally linear matrix \mathbf{R}^* from \mathbf{C} , while incorporating the K -sets structure in \mathbf{R}^* .

The (canonical) correlations between the various obtained quantifications are derived by introducing two m -sized columns of the identity matrix for each variable; the vectors \mathbf{v} and \mathbf{w} with a one in the position of the respective variable(s) in question. The covariances, variances and correlations of the quantifications are now expressed as, respectively:

$$\mathbf{Cov}_{\mathbf{v}\mathbf{w}} = \mathbf{A}'_S (\mathbf{R}^* \wedge (\mathbf{v}\mathbf{w}')) \mathbf{A}_S, \quad (13a)$$

$$\text{Var}_{\mathbf{v}\mathbf{v}} = \mathbf{A}'_s (\mathbf{R}^{\wedge}(\mathbf{v}\mathbf{v}')) \mathbf{X}_{s'} = \text{Var}_{\mathbf{w}\mathbf{w}} = \mathbf{X}'_s (\mathbf{R}^{\wedge}(\mathbf{w}\mathbf{w}')) \mathbf{A}_s, \quad (13b)$$

$$\text{Cor}_{\mathbf{v}\mathbf{w}} = (\text{diag. } \text{Var}_{\mathbf{v}})^{1/2} \text{Cov}_{\mathbf{v}\mathbf{w}} (\text{diag. } \text{Var}_{\mathbf{w}})^{1/2}. \quad (13c)$$

With the appropriate filling of vectors \mathbf{v} and \mathbf{w} one can now compute for example: the canonical correlations - with ones in \mathbf{v} and \mathbf{w} on the positions indicating the variables of the sets, implying $\mathbf{Z}_{\mathbf{v}\mathbf{w}} = \mathbf{A}'\mathbf{R}^*\mathbf{A} = \mathbf{I}$. In a similar fashion one can compute the correlations between a quantified variable and the respective canonical variate (also referred to as *canonical loadings*), reflecting the extent to which a variable contributes in the variate.

6. ALGORITHMIC FEATURES

De Leeuw (1983b) describes the outline of an algorithm which can perform multiset nonlinear CCA based on the joint bivariate analysis as described in the previous section. A workable algorithm is presented in Tijssen (1985). Some specific features of this algorithm are treated in some detail in this section.

6.1 MULTIPLE CATEGORY QUANTIFICATIONS VIA THE COPIES-METHOD

In the Gifi-system the types of regressions/transformations leading to optimally scaled categories \mathbf{q}_j can be applied to variables in two different ways. First, if the variables j are transformed *single*, we demand that the transformations per dimension are linearly related. Hence one set of quantifications suffices. CCA will provide us with a $(k_j * 1)$ vector \mathbf{y}_j of category quantifications, the so-called *single* quantifications. Secondly, in case of a nonlinear transformations - applicable to nominal and ordinal variables - we have the additional possibility to solve the regression problem for each dimension separately, which provides us with a $(k_j * p)$ matrix \mathbf{Y}_j with so-called *multiple* category quantifications. However, CCA implies a fixed p -dimensional eigenvalue-eigenvector solution. A vector of independent (orthogonal) category quantifications are therefore not computable within this framework. Analogous to other CCA-programs within the Gifi-system (e.g., CANALS, OVERALS) one might want to obtain such multiple quantifications for categories. A remedy for this problem has been developed by De Leeuw (1983a): if we introduce $t = \min.(p, k_j - 1)$ *copies* of a variable into the set in question, we obtain t separate quantifications of the categories. In this case maximizing the sum of the first p dimensions yields t rank-one (or single) quantifications \mathbf{y}_j . Different numbers of copies, ranging between 1 (the single quantification) and t (the multiple quantification), and applying different measurement levels to different copies enables one to arrive at tailored category quantifications.

6.2 A RANK REDUCING POLYNOMIAL BASIS

The problem concerning possible singularities in \mathbf{C} can be avoided by means of the following rank reducing method (De Leeuw, 1985). Suppose \mathbf{Z}_j is a $(k_j * k_j)$ matrix, satisfying $\mathbf{Z}_j' \mathbf{D}_j \mathbf{Z}_j = \mathbf{I}$. Moreover, suppose that the elements in the first column of \mathbf{Z}_j are all equal to 1. The last $k_j - 1$ columns

of \mathbf{Z}_j span the subspace of all quantifications in deviations of the mean. The first column of \mathbf{Z}_j is thus trivial within a proper solution and can therefore be deleted. A vector \mathbf{y}_j derived from \mathbf{Z}_j can be written as a linear combination of the remaining k_j-1 columns of \mathbf{Z}_j . In case of the bivariate matrix we can write $\underline{\mathbf{C}}_{jl} = \mathbf{Z}_j' \mathbf{C}_{jl} \mathbf{Z}_l$, reducing \mathbf{C} to a rank $\sum_j k_j - m$ and, in general, implying non-singularity.

Orthogonal polynomials based on the category scores will be used for the k_j-1 vectors in \mathbf{Z}_j - matrices, replacing indicator matrices \mathbf{G}_j (see section 3). We start with $(k_j * k_j)$ matrices \mathbf{V}_j , with a first column with unity-elements and subsequent column vectors containing linear, quadratic, cubic, fourth and fifth degree functions of the original category scores. Additional vectors will consist of k_j pseudo-random values within the range $[0, \max(\mathbf{D}_j)]$.

A modified Gram-Schmidt (GRAM)-orthogonalization of \mathbf{V}_j provides us with \mathbf{Z}_j with k_j column vectors, each representing an orthogonal polynomial weighted with respect to the marginal frequencies \mathbf{D}_j : $\mathbf{Z}_j = \text{GRAM}(\mathbf{D}_j^{1/2} \mathbf{V}_j)$, with $\mathbf{Z}_j' \mathbf{Z}_j = \mathbf{I}$. This is followed by the step $\mathbf{Z}_j = \mathbf{D}_j^{-1/2} \mathbf{Z}_j$, with $\mathbf{Z}_j' \mathbf{D}_j \mathbf{Z}_j = \mathbf{I}$.

The matrices \mathbf{Z}_j now form a basis for \mathbb{L}_j consisting of k_j orthogonal polynomials in the \mathbf{D} -metric with

$$\underline{\mathbf{C}} = \mathbf{Z}' \mathbf{C} \mathbf{Z}, \tag{14}$$

while $\mathbf{Z}_j' \underline{\mathbf{C}}_j \mathbf{Z}_j = \mathbf{I}$.

Submatrices $\underline{\mathbf{C}}_{jl}$ now possess a similar first row and column. Removing these first rows and columns means reducing the order of the subspaces \mathbb{L}_j to a dimensionality of k_j-1 , creating a non-singular reduced $(\sum_j k_j - m * \sum_j k_j - m)$ matrix $\underline{\mathbf{C}}_{jl}$.

The algorithm allows one to choose an individual polynomial basis per variable (including copies), implying a dimensionality of \mathbf{U}_j and thus \mathbf{Z}_j - matrices dependent on the particular choice. It is, for instance, possible to use a quadratic basis for a particular (nonlinear) variable and a cubic polynomial basis for another (nonlinear) variable.

6.2 THE ITERATION PROCESS.

Using copies of single variables to obtain multiple variables has an additional advantage with respect to our choice of using orthogonal polynomials. Suppose we perform a p -dimensional CCA with variables j , consisting of k_j categories with $t = 1$ copies. We can define a $(\sum_j k_j - m * mp)$ matrix $\mathbf{W} = \mathbf{w}_1 + \dots + \mathbf{w}_m$, consisting of the $(k_j-1 * p)$ submatrices \mathbf{w}_j with a one on the respective position of the category, representing the initial estimations of the weights of categories for j (in case of $t > 1$ copies one has t binary orthogonal columns for j). Removing the first - superfluous - column of \mathbf{Z}_j , yields k_j-1 single category quantifications \mathbf{y}_j by means of

$$y_j = Z_j w_j, \quad (15)$$

with the $(k_j * k_j - 1)$ matrix Z_j and the $(k_j - 1 * 1)$ vector w_j , containing the category weights. Category quantifications of additional copies are thus computed by the corresponding columns of the weight vectors.

In the previous chapter we described the main and inner-iteration proces in terms of the category quantifications y_j . In the actual iteration proces we will, however, only use the weight vectors w_j . The optimal category quantifications will be computed through the bases Z_j as in Eq. (15), with use of optimal weight values. To obtain these weights we follow the procedure analogous to computing the optimal quantifications as described in Eq. (9). In the iteration process these weights will thus become a function of the between and within-set bivariate structure. This means that our initial (and orthogonal) estimates weights w_j will loose their orthogonality, but will now be become optimal in the sense that they are a optimal function of both the between and the within-sets relations of the variables. Due to the normalization $Z_j' C_j Z_j = I$, we now obtain the corresponding $(m * m)$ covariance matrix F , induced by w_j by means of

$$F = \{f_{ij}\} = w_j' C_{jl} w_l, \quad (16)$$

yielding the updated correlation matrix

$$R^* = \{r^*_{jl}\} = f_{jj}^{-1/2} f_{jl} f_{ll}^{-1/2}. \quad (17)$$

Assuming we only have multiple variables (i.e., with multiple category quantifications) we now obtain a $(mp * mp)$ matrix R^* with correlations between t copies of m variables, based on the weights w_j . If we only use single variables we obtain a $(m * m)$ -sized R^* in which each element represents the linear relationship between the variables.

The subsequent generalized eigenvalue problem of $(R^*, E^* R^*)$ is solved by the reduction to a standard symmetric eigenvalue problem in the following way. Let $E^* R^* = S S'$ be the Choleski decomposition of the matrix $E^* R^*$, where S is a non-singular lower triangular matrix (Wilkinson, 1965). This decomposition is extremely stable when the matrix $E^* R^*$ is symmetric and positive semi-definite. Equation (5) is now equivalent to

$$S^{-1} R^* (S^{-1})' S X = S X \lambda, \quad (18)$$

Computing A and B simply consist of weighing the cells of the non-singular matrix C . This enables us to compute separate Choleski decompositions of the k non-zero submatrices $A_k = V_k' V_k$, without

singularity problems. Equation (9) is solved via $\mathbf{U}^{(n+1)} = \mathbf{T}^{-1} (\mathbf{T} \mathbf{U}^{(n)})$ with $\mathbf{B} \mathbf{U}^{(n)} = \mathbf{T}'$. The orthogonality of copies within a variable is restored by means of a Gram-Schmidt orthogonalization after each inner-iteration. Due to the special properties of \mathbf{C}_{jl} we can now compute an update of the correlation matrix by simply repeating the steps main and inner-iteration steps with optimal $\mathbf{U}^{(n+1)}$ updates.

An optimal number of inner-iteration steps can not be given at this stage of development. However, investigations of the algorithm performance have shown that a few inner-iterations can produce convergence to the similar stable values with less main-iterations. Reducing the number of inner-iterations to one seems, for the time being, the safest policy and will, in general, provide a sufficiently stable canonical solution. The sum of the p -largest eigenvalues corresponding to each subsequent generalized eigenvalue solution will be referred to as the "fit" of the CCA-solution. The difference between the fit of each pair of subsequent main-iteration steps is used as a criterion to terminate the main-iteration process.

Although a theoretical proof of convergence of this algorithm is not available so far, empirical studies of this iteration process have shown that the convergence to a stable value is attained in all applied cases. This algorithm can also be troubled by sub-optimal solutions, the so-called *local minima /maxima*.

6.4 ADDITIONAL OUTPUT

Eventually the iteration process provides us with an optimal set of k_j-1 category weights \mathbf{w}_j corresponding to the k_j-1 dimensions in each \mathbb{L}_j , spanned by the orthogonal polynomials. The resulting category quantifications \mathbf{y}_j will then be composed of a optimal linear part, a quadratic part, ..., etcetera, due to the specific properties of the bases \mathbf{Z}_j . This particular approach enables a computation of the total variance of a variable j via

$$\mathbf{y}'_j \mathbf{D}_j \mathbf{y}_j = \mathbf{w}'_j \mathbf{z}'_j \mathbf{D}_j \mathbf{z}_j \mathbf{y}_j = \mathbf{w}'_j \mathbf{w}_j, \quad (19)$$

in terms of the weights corresponding to the specific polynomials. This enables us to partition the total variance of a variable in (a percentage) linear, quadratic and other existing polynomial contributions. Hence insight is provided with respect to the extent in which the transformations from category score to category quantification have been linearized.

Obtaining the appropriate category quantifications according to the restrictions imposed by measurement levels (see section 3) is achieved as follows: nominal variables, of course, do not require any additional restrictions in the algorithm - the category quantifications are simply the resulting optimal transformations of the initial category scores. The numerical measurement level -with its linearity property- can be incorporated into the algorithm by simply constraining the first element of \mathbf{w}_j

to unity and the additional elements to zero, before each iteration. Thus one only uses the bivariate relations in \underline{C} which correspond to the linear (first-degree) polynomial. The quantifications resulting from the iteration process are now always a linear function of the first orthogonal polynomial and thus a linear function of the original category scores. Ordinal variables can be computed by appropriate monotonic transformations on the orthogonal polynomials. Implementation of this type of regression is still a subject of investigation.

7. AN APPLICATION

7.1. INTRODUCTION

An example is given of a data-analytic situation in which the sheer magnitude of the matrix ($n=2666$; $m=77$) matrix may warrant the use of a K -sets CCA method based on a smaller Burt-matrix. The data originate from a survey carried out by the Educational Department and the Educational Research Center of the University of Leiden to compare pupil characteristics of the following three groups of elementary vocational educations (in order of increasing educational level): VBO/MLK, IBO and LBO (Van Putten, 1987). The study focussed on the assessment of characteristics of the IBO pupils. The Dutch acronym "IBO" refers to a number of individualized elementary vocational educations, with the aim to provide a suitable education for slow learners from the primary school or special educations, in the pupil age of 12-16 years.

These pupils were subjected to a number of tests, in order to assess (1) - cognitive capacities, (2) - reading abilities, and (3) - arithmetic abilities. They also received (4) - a questionnaire with items concerning their attitude on some relevant socio-emotional aspects of the school/classroom-situation, (5) - two teachers were asked -independently- to rate the classroom behavior of each pupil, and (6) - background information was obtained on the pupils by means of a parent/caretaker-questionnaire. Each test/questionnaire contained well over 400 categories.

7.2. DATA PRE-PROCESSING

Data-reduction was achieved by applying the nonmetric Principal Components Analysis program PRINCALS (Gifi, 1985) to the above mentioned groups of variables. The results provided us with "condensed" data in the form of the following independent and relatively clear-cut defined components, on which each pupil obtained a score. These metricized pupil-quantifications were divided into a number of categories (on the average, about 5 categories). In the case of above groups (1)-(5), the low-valued categories indicate the positive/desirable side of the attribute, for instance a high degree of accuracy or no fear of failure, whereas high category scores indicate the opposite qualification.

The original survey variables were reduced to the following (composite) variables :

- | | | |
|---------------------------|---------------------------------|---------------------------------|
| 1-Logical reasoning (LR) | 2-Workspace (WP) | 3-Accuracy (AC) |
| 4-Verbal abilities (VA) | 5-Technical reading (TR) | 6-Mental arithmetic (MA) |
| 7-Applied arithmetic (AA) | 8-Involvement and interest (II) | 9-Problematic behaviour (PB) |
| 10-Independence (IN) | 11-Motivation (MO) | 12-Relations fellow-pupils (RP) |
| 13-Fear of failure (FF) | | |

The following background variables were selected for the CCA analysis:

- | | |
|--|-------------------------------------|
| 14-Previous education (PRE) | 15-Promotion previous year (PPY) |
| 16-Gender (GEN) | 17-Age (AGE) |
| 18-Occupational family type (OFT) | 19-Number of brothers/sisters (NBS) |
| 20-Birth rank among brothers/sisters (BBS) | |

The above variables were placed in the following sets: (I) a cognitive set (variables 1-7), (II) a socio-emotional set (variables 8-13), and (III) the background information (variables 14-20), respectively.

The pupil scores in the cognitive and socio-emotional set were analyzed with a numerical measurement level. The optimal (linear) properties of these scores provides a plausible basis for such a restriction in the category quantification. The remaining (background) variables are of a more qualitative nature were analyzed with the nominal measurement level, with the exception of a pupil's age.

7.3. SOME CCA ANALYSIS RESULTS

The CCA analysis was done for 2 dimensions. The fit of the analysis solution was equal to 1.034, with the eigenvalues .530 and .504 for the first and second dimension, respectively. Both dimensions can thus be regarded as almost equally important and - on the whole - "explaining" about half of the existing variance (the maximum fit is equal to 2). The highest canonical correlations were found between the cognitive set and the set of background variables (with values equal to .357 and .405 for the first and second dimension, respectively).

A general description of the relations between the separate variables, given the canonical structure, can be given on the basis of the configuration of the correlations between variables and the mean canonical variate (see Fig. 1). Projecting the points on the axes reveals a structure, in which the first dimension is determined by applied arithmetic (AA), previous education (PRE), gender (GEN) and - to a lesser extent - by variables such as problematic behaviour (PB). The second dimension is also largely determined by GEN and PRE, but now in combination with technical reading (TR), mental arithmetic (MA), verbal abilities (VA) and motivation (MO).

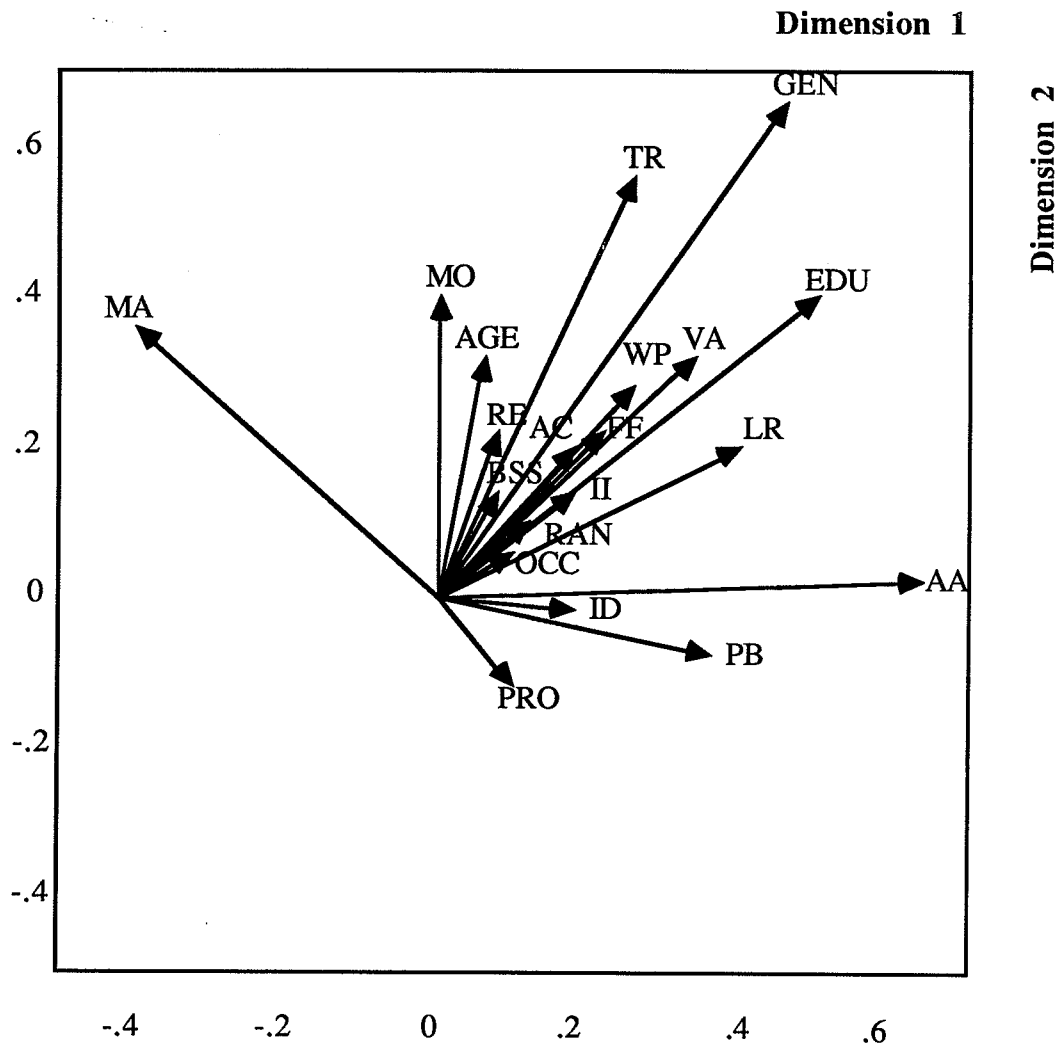


Fig. 1- Correlations of variables with mean canonical variate

Further insight in the analysis results can be obtained by projecting the linear category quantifications in the mean canonical space (see Fig. 2). These quantifications present an image of the structure of relations on a category-level; for each variable the categories are located on a straight line through the origin. Both GEN-categories, the 'extreme' categories of PRE and AGE and the lowest- and highest valued category quantifications are shown for the cognitive and socio-emotional variables. With respect to the latter, the variable problematic behaviour (PB) has, for example, a label PB+ for non-problematic behaviour and PB- for problematic behaviour. The linear quantification of two (extreme) categories of GEN, PRE and AGE are shown connected with a line.

As for the interpretation of these results, consider the perpendicular lines of AGE and PRE which reveal a boy-girl distinction on the cognitive- and socio-emotional variables, which is independent of the previous educational level. Differences between boys and girls are found in both dimensions, for example: boys are -on the average- older than the girls; tend to have relatively better marks on applied arithmetic-tests (AA+); show less fear of failure (FF+); more independent behavior (IN+); less accuracy (AC-) and are found to have less classroom involvement and interest (II-).

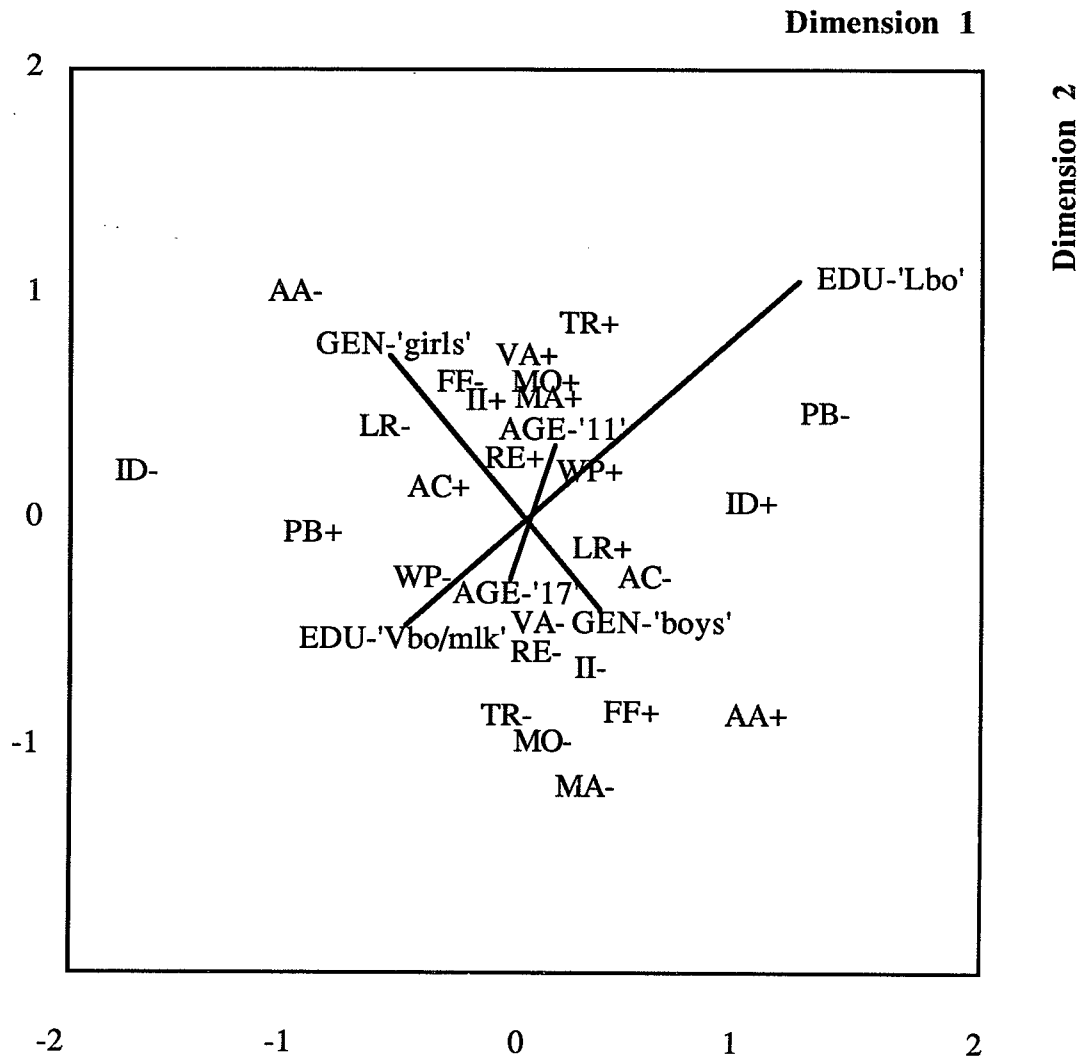


Fig. 2 - Selected rank-one category quantifications.

Information on the nature of the quantifications can be found in the partitioning of the variance accounted by these quantifications in their polynomial parts (see section 6.4). Table 1 displays the partitioning for the background variables.

Table 1. - Partitioned variance (in percentages)

Variable	Linear	Quadratic	Cubic
NBS	63.3	17.0	19.7
BBS	0.0	3.3	96.6
AGE	100.0		
OFT	98.0	1.0	1.0
GEN	100.0		
PPY	100.0		
EDU	30.7	57.1	12.2

For example, it is shown that a relatively small linearization of the transformations from category score to quantification is found for the education-variable EDU; in this $p=2$ case the lowest and highest educational types received more extremely valued category quantifications in the second dimension. The other (numerical) variables are of course all linearized and thus have 100% linear variance. This is obviously also the case for the binary variables GEN and PPY.

8. CONCLUSIONS

The presented technique based on a joint bivariate analysis seems a potentially fruitful approach to nonlinear K -sets CCA in case of large matrices; when $\sum_j k_j \ll n$ such an algorithm, should be computationally more efficient, in terms of core-memory and CPU processing time. Some pilot applications have indicated that the algorithm is indeed a relatively fast way of computing an - at least locally - stable CCA solution. Although a theoretical proof of convergence is not available (so far), empirical studies of the iteration process have shown that convergence was obtained in all cases. The technique is an interesting alternative to a nonlinear K -sets CCA approach as in the program OVERALS, which is based on the (objects \times variables) data matrix.

REFERENCES

- Anderson, T. W. (1958). *Introduction to Multivariate Statistical Analysis*. New York: Wiley.
- Burt, C. (1950). The factorial analysis of qualitative data. *British Journal of Statistical Psychology*, 3, 166-185.
- Gifi, A. (1981). *Nonlinear multivariate analysis*. Department of Data Theory. University of Leiden.
(A revised version will be published by DSWO-Press in 1988).
- Gifi, A. (1985). *PRINCALS*, UG-85-03. Department of Data Theory. University of Leiden.
- Gittens, R. (1984). *Canonical analysis-A review with applications in ecology*. Berlin: Springer Verlag.
- Hotelling, H. (1936). Relations between two sets of variables. *Biometrika*, 28, 321-377.
- Kettenring, J.R. (1971). Canonical analysis of several sets of variables. *Biometrika*, 58, 433-460.
- De Leeuw, J. (1983). *Nonlinear joint bivariate analysis*. Paper presented at the European Meeting of the Psychometric Society. Jouy-en-Josas. France.
- De Leeuw, J. (1984). *Canonical Analysis of Categorical Data*. (Doctoral dissertation, University of Leiden, 1973). Leiden: DSWO-Press.
- De Leeuw, J. (1986). *Multivariate Analysis with Optimal Scaling*. Paper presented at the International Conference on Advances in Multivariate Statistical Analysis. Calcutta. India.
- Meulman, J.J. (1986). *A distance approach to nonlinear multivariate analysis*. Leiden: DSWO-Press.
- Styan, G.P. (1973). Hadamard products and multivariate statistical analysis. *Linear Algebra and its applications*, 6, 217-240.
- Thomson, G.H. (1947). The maximum correlation of two weighted batteries. *British Journal of Psychology - Statistical Section*, 1, 27-34.
- Tijssen, R.J.W. (1985). *A new approach to nonlinear canonical correlation analysis*. Leiden Psychological Reports - Psychometrics and Research methodology, RR 85-01. University of Leiden.
- Van der Burg, E. and De Leeuw, J. (1983). Non-linear canonical correlation. *British Journal of Mathematical and Statistical Psychology*, 36, 54-80.
- Van der Burg, E., De Leeuw, J. and Verdegaal, R. (1988). Nonlinear canonical correlation with m sets of variables. *Psychometrika*, 2, 171-197.
- Van de Geer, J.P. (1984). Linear relations between k sets of variables. *Psychometrika*, 49, 79-94.
- Van Putten, C.M. (1987). *Leerlingen in het individueel beroepsonderwijs nader beschouwd*. (Doctoral dissertation - in Dutch). University of Leiden.
- Wilkinson, J.H. (1965). *The algebraic Eigenvalue Problem*. Oxford: Clarendon Press.
- Young, F.W., De Leeuw, J. and Takane, Y. (1976). Regression with qualitative and quantitative variables: an alternating least squares method with optimal scaling features. *Psychometrika*, 43, 279-281.
- Young, F.W. (1981). Quantitative Analysis of Qualitative Data. *Psychometrika*, 46, 357-388.