



Invited Review

Regularized generalized canonical correlation analysis for multiblock or multigroup data analysis

Arthur Tenenhaus^{a,*}, Michel Tenenhaus^{b,1}^aSUPELEC, Plateau de moulon, 3 rue Joliot-Curie, 91192 Gif-sur-Yvette Cedex, France^bHEC Paris, 1 rue de la Libération, 78351 Jouy-en-Josas Cedex, France

ARTICLE INFO

Article history:

Received 2 November 2012

Accepted 3 January 2014

Available online 13 January 2014

Keywords:

Multiblock data analysis

Multigroup data analysis

Regularized generalized canonical correlation analysis

ABSTRACT

This paper presents an overview of methods for the analysis of data structured in blocks of variables or in groups of individuals. More specifically, regularized generalized canonical correlation analysis (RGCCA), which is a unifying approach for multiblock data analysis, is extended to be also a unifying tool for multigroup data analysis. The versatility and usefulness of our approach is illustrated on two real datasets.

© 2014 Elsevier B.V. All rights reserved.

1. Introduction

In this paper, we consider a data matrix \mathbf{X} structured in groups (partition of rows) or in blocks (partition of columns). Rows of \mathbf{X} are related to individuals and columns to variables. Multiblock data analysis concerns the analysis of several sets of variables (blocks) observed on the same set of individuals. Multigroup data analysis concerns the analysis of one set of variables observed on several groups of individuals. Note that there is no established consensus in the literature on the use of the terms “multiblock” and “multigroup”. Therefore these two terms are clearly defined in this paper.

In the *multiblock* framework, a column partition $\mathbf{X} = [\mathbf{X}_1, \dots, \mathbf{X}_j, \dots, \mathbf{X}_J]$ is considered. In this case, each $n \times p_j$ data matrix \mathbf{X}_j is called a block and represents a set of p_j variables observed on n individuals. The number and the nature of the variables usually differ from one block to another but the individuals must be the same across blocks. The main aim is to investigate the relationships between blocks. The data might be preprocessed in order to ensure comparability between variables and blocks. To make variables comparable, standardization is applied (zero mean and unit variance). To make blocks comparable, a possible strategy is to divide each block by $\sqrt{p_j}$ (Wold, Hellberg, Lundstedt, Sjöstrom, & Wold, 1987). This two-step procedure leads to $\text{Trace}(\mathbf{X}_j^t \mathbf{X}_j) = n$ for each block.

* Corresponding author. Address: SUPELEC, Department of Signal Processing and Electronics Systems, Plateau de Moulon, 3 rue Joliot-Curie, 91192 Gif-sur-Yvette Cedex, France. Tel.: +33 (0)1 69 85 14 22; fax: +33 (0)1 69 85 14 29.

E-mail addresses: Arthur.Tenenhaus@supelec.fr (A. Tenenhaus), tenenhaus@hec.fr (M. Tenenhaus).

¹ Tel.: +33 (0)1 39 67 70 00; fax: +33 (0)1 39 67 74 00.

In the *multigroup* framework, a row partition $\mathbf{X} = [\mathbf{X}_1^t, \dots, \mathbf{X}_i^t, \dots, \mathbf{X}_I^t]^t$ is considered. In this framework, the same set of variables is observed on different groups of individuals. Each $n_i \times p$ data matrix \mathbf{X}_i is called a group. The number of individuals of each group can differ from one group to another. The main aim is to investigate the relationships among variables within the various groups. Following the proposal of Kiers and Ten Berge (1994) variables are centered and normalized (i.e. set to unit norm) within each group. This preprocessing leads to $\text{Trace}(\mathbf{X}_i^t \mathbf{X}_i) = p$ for each group.

Many methods exist for multiblock and multigroup data analysis.

Two families of methods have come to the fore in the field of multiblock data analysis. These methods rely on correlation-based or covariance-based criteria. Canonical correlation analysis (Hotelling, 1936) is the seminal paper for the first family and Tucker's inter-battery factor analysis (Tucker, 1958) for the second one. These two methods have been extended to more than two blocks in many ways:

- (1) Main contributions for generalized canonical correlation analysis (GCCA) are found in Horst (1961), Carroll (1968a), Kettenring (1971), Wold (1982, 1985) and Hanafi (2007).
- (2) Main contributions for extending Tucker's method to more than two blocks come from Carroll (1968b), Chessel and Hanafi (1996), Hanafi and Kiers (2006), Hanafi, Kohler, and Qannari (2010, 2011), Hanafi, Mazerolles, Dufour, and Qannari (2006), Krämer (2007), Smilde, Westerhuis, and de Jong (2003), Ten Berge (1988), Van de Geer (1984), Westerhuis, Kourti, and MacGregor (1998), Wold (1982, 1985).

- (3) Carroll (1968b) proposed the “mixed” correlation and covariance criterion. van den Wollenberg (1977) combined correlation and variance for the two-block situation (redundancy analysis). This method is extended to a multiblock situation in this paper.

Regularized generalized canonical correlation analysis (Tenenhaus & Tenenhaus, 2011) includes many of these references as particular cases.

For multigroup data analysis, we may distinguish three families of methods:

(1) Several methods combine the covariance matrices S_i or the correlation matrices R_i related to the various groups.

In a seminal paper, Levin (1966), considering the problem of simultaneous factor analysis, proposed the diagonalization of $\bar{R} = (1/I)\sum_{i=1}^I R_i$. The acronym SUMPCA_c for the Levin method was proposed by Kiers (1991). Kiers and Ten Berge (1994) proposed several simultaneous component analysis methods (see paragraph 3 below): one of them (SCA-P) leads also to the diagonalization of \bar{R} . Krzanowski (1984) proposed to carry out a multigroup PCA (MGPCA) by diagonalizing either $T = \sum_{i=1}^I S_i$ or the within group covariance matrix $S = \sum_{i=1}^I (n_i/n)S_i$.

Krzanowski (1979) proposed to use an approach similar to Carroll’s GCCA (correlation criterion) for comparing group correlation matrices R_1, \dots, R_I . Since GCCA on these matrices yields a trivial solution, Krzanowski replaced each matrix R_i by its k first eigenvectors and then applied GCCA on the obtained matrices.

Flury (1984) proposed a method called common principal component analysis (CPC) for I groups by supposing a special structure on the covariance matrices $\Sigma_1, \dots, \Sigma_I$ defined at the population level. In CPC, the I covariance matrices have the same eigenvectors but the eigenvalues are specific to each group: $\Sigma_i = AA_iA^t$ where A is orthogonal and A_i diagonal. Flury (1987) also proposed a partial common principal component analysis (PCPC) where only q eigenvectors of Σ_i are common to all populations. Flury’s approach is based on normal-theory maximum likelihood and a complicated iterative algorithm is required. Two alternative algorithms have been proposed: (1) Krzanowski (1984) showed empirically that MGPCA and PCPC give very close results; (2) a sequential least squares solution to PCPC can be obtained by using the CCSWA algorithm (Common components and specific weight analysis) described in Hanafi et al. (2006). It is also worth mentioning that CCSWA and HPCA (Hierarchical principal component analysis described in Westerhuis et al. (1998)) are two equivalent methods (Hanafi et al., 2010).

(2) It is always possible to use a multiblock method on multigroup data by considering the transpose of each group. Eslami, Qannari, Kohler, and Bougeard (2013a) proposed to use an approach similar to Carroll’s GCCA (correlation and covariance criteria) on the transpose groups X_1^t, \dots, X_I^t . This approach has later been extended to a multiblock/multigroup situation in Eslami, Qannari, Kohler, and Bougeard (2013b).

(3) Kiers and Ten Berge (1989, 1994) and later Timmerman and Kiers (2003) proposed a generalization of PCA to a multigroup situation under the generic name of Simultaneous Component Analysis (SCA). Each data group X_i of dimension $n_i \times p$ is modeled by a lower rank $n_i \times p$ matrix $\hat{X}_i = X_i W_i P_i^t$ where W_i is a $p \times q$ ($q < p$) weight matrix and P_i a $p \times q$ pattern matrix. A factor matrix $F_i = X_i W_i$ and a loading (or structure) matrix $L_i = R_i W_i$ are also introduced. PCA and SCA methods are about minimizing $\sum_{i=1}^I \|X_i - X_i W_i P_i^t\|^2$ subject to specific constraints on the weight/pattern/structure/factor matrices which are summarized in Table 1.

The reconstructed matrix $\hat{X}_i = X_i W_i P_i^t$ is invariant up to an orthogonal (rotation) matrix A : $X_i W_i P_i^t = X_i W_i A^t A P_i^t$. This invariance can be used to improve interpretation. Various rotation methods are described in Niesing (1997). Moreover, Niesing shows in a comparative study that SCA-P gives the best practical results.

Table 1
PCA and SCA methods.

Methods	Constraints
Separate PCA by group (Pearson (1901))	$P_i = W_i$ and $W_i^t W_i = I \forall i$
SCA-W (Kiers & Ten Berge (1989))	$W_1 = \dots = W_I$
SCA-P (Kiers & Ten Berge (1994))	$P_1 = \dots = P_I$
SCA-S (Kiers & Ten Berge (1994))	$R_i W_i = L_i \forall i^a$
SCA-PF2 (Timmerman & Kiers (2003))	$P_1 = \dots = P_I$ and $F_i^t F_i = D_i \Phi D_i \forall i^b$
SCA-IND (Timmerman & Kiers (2003))	$P_1 = \dots = P_I$ and $F_i^t F_i = D_i^2 \forall i$
SCA-ECP (Timmerman & Kiers (2003))	$P_1 = \dots = P_I$ and $(1/n_i) F_i^t F_i = I \forall i$

^a L is an unknown ($p \times q$) matrix and D_i an unknown ($q \times q$) diagonal matrix.
^b Φ is an unknown ($q \times q$) positive definite matrix with unit diagonal elements.

De Roover, Ceulemans, and Timmerman (2012), De Roover, Ceulemans, Timmerman, and Onghena (2013) and De Roover, Timmerman, Van Mechelen, and Ceulemans (2013) have developed a clusterwise approach to SCA-P, SCA-PF2, SCA-IND and SCA-ECP for tracing structural differences and similarities between different groups of individuals.

Finally, Van Deun, Smilde, van der Werf, Kiers, and Van Mechelen (2009) proposed a simultaneous component analysis framework for multiblock and multigroup data analysis.

In this paper, a modified version of RGCCA, which can be applied to multiblock and multigroup data, is described. The acronym RGCCA will be kept for this more general method. This paper is organized as follows: the general optimization problem for both multiblock and multigroup data analysis is presented in Section 2. A monotonically convergent algorithm is presented in Section 3. An overview of applications of RGCCA for multiblock and multigroup data analysis is given in Sections 4 and 5. The versatility and usefulness of our approach is illustrated on two real datasets in Sections 6 and 7.

2. The optimization problem behind RGCCA for multiblock or multigroup data analysis

RGCCA for multiblock and multigroup data analyses is based on a single optimization problem that we present in this section. We consider I matrices Q_1, \dots, Q_I . Each matrix Q_i is of dimension $m \times p_i$. We also associate to each matrix Q_i a weight column vector w_i of dimension p_i and a symmetric definite positive matrix M_i of dimensions $p_i \times p_i$. Moreover, a design matrix $C = \{c_{i\ell}\}$ is defined with $c_{i\ell} = 1$ if matrices Q_i and Q_ℓ are connected, and $c_{i\ell} = 0$ otherwise. The core optimization problem considered in this paper is defined as follows:

$$\begin{aligned} & \text{Maximize}_{w_1, \dots, w_I} \sum_{i, \ell, i \neq \ell} c_{i\ell} g(\langle Q_i w_i, Q_\ell w_\ell \rangle) \\ & \text{s.c. } w_i^t M_i w_i = 1, \quad i = 1, \dots, I \end{aligned} \tag{1}$$

where $\langle x, y \rangle = x^t y$ is the usual scalar product and g stands for the identity, the absolute value or the square function. By setting $v_i = M_i^{1/2} w_i$ and $P_i = Q_i M_i^{-1/2}$ optimization problem (1) becomes

$$\begin{aligned} & \text{Maximize}_{v_1, \dots, v_I} \sum_{i, \ell, i \neq \ell} c_{i\ell} g(\langle P_i v_i, P_\ell v_\ell \rangle) \\ & \text{s.c. } v_i^t v_i = 1, \quad i = 1, \dots, I \end{aligned} \tag{2}$$

A monotone convergent algorithm can be developed for optimization problem (2). This algorithm will be presented in detail in the next section. It is worth mentioning that all the multiblock and multigroup methods to be presented in this paper are special cases of optimization problem (2).

Download English Version:

<https://daneshyari.com/en/article/478102>

Download Persian Version:

<https://daneshyari.com/article/478102>

[Daneshyari.com](https://daneshyari.com)