

Tilburg University

Clusterwise Simultaneous Component Analysis for Analyzing Structural Differences in Multivariate Multiblock Data

De Roover, Kim; Ceulemans, Eva; Timmerman, Marieke E.; Vansteelandt, Kristof; Stouten, Jeroen; Onghena, Patrick

Published in:
Psychological Methods

DOI:
[10.1037/a0025385](https://doi.org/10.1037/a0025385)

Publication date:
2012

Document Version
Peer reviewed version

[Link to publication in Tilburg University Research Portal](#)

Citation for published version (APA):

De Roover, K., Ceulemans, E., Timmerman, M. E., Vansteelandt, K., Stouten, J., & Onghena, P. (2012). Clusterwise Simultaneous Component Analysis for Analyzing Structural Differences in Multivariate Multiblock Data. *Psychological Methods*, 17(1), 100-119. <https://doi.org/10.1037/a0025385>

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Clusterwise Simultaneous Component Analysis for Analyzing Structural Differences in
Multivariate Multiblock Data

Kim De Roover

Katholieke Universiteit Leuven

Eva Ceulemans

Katholieke Universiteit Leuven

Marieke E. Timmerman

University of Groningen

Kristof Vansteelandt

Universitair psychiatrisch centrum, Katholieke Universiteit Leuven

Jeroen Stouten

Katholieke Universiteit Leuven

Patrick Onghena

Katholieke Universiteit Leuven

Author Notes:

The research reported in this paper was partially supported by the fund for Scientific Research-Flanders (Belgium), Project No. G.0477.09 awarded to Eva Ceulemans, Marieke Timmerman and Patrick Onghena and by the Research Council of K.U.Leuven (GOA/10/02). Correspondence concerning this paper should be addressed to Kim De Roover, Department of Educational Sciences, Andreas Vesaliusstraat 2, B-3000 Leuven, Belgium. E-mail: Kim.DeRoover@ped.kuleuven.be.

Abstract

Many studies yield multivariate multiblock data, that is, multiple data blocks that all involve the same set of variables (e.g., the scores of different groups of subjects on the same set of variables). The question then rises whether or not the same processes underlie the different data blocks. To explore the structure of such multivariate multiblock data, component analysis can be very useful. Specifically, two approaches are often applied: principal component analysis (PCA) on each data block separately and different variants of simultaneous component analysis (SCA) on all data blocks simultaneously. The PCA approach yields a different loading matrix for each data block and is thus not useful for discovering structural similarities. The SCA approach may fail to yield insight into structural differences, since the obtained loading matrix is identical for all data blocks. We introduce a new generic modeling strategy, called Clusterwise SCA, that comprises the separate PCA approach and SCA as special cases. The key idea behind Clusterwise SCA is that the data blocks form a few clusters, where data blocks that belong to the same cluster are modeled with SCA and thus have the same structure, and different clusters have different underlying structures. In this paper, we use the SCA-ECP variant of SCA. An algorithm for fitting Clusterwise SCA-ECP solutions is proposed and evaluated in a simulation study. Finally, the usefulness of Clusterwise SCA is illustrated by empirical examples from eating disorder research and social psychology.

Keywords: multivariate data, multigroup data, multilevel data, principal component analysis, simultaneous component analysis, clustering

Introduction

In the behavioral sciences, many studies yield multivariate multiblock data, that is, multiple data blocks that all involve the same set of variables. For an example, one can think of data from different groups of subjects that are measured on the same variables or of data from multiple subjects that have scored the same variables on multiple measurement occasions (the latter type of data are sometimes called multioccasion-multisubject data; see Kroonenberg, 2008). In the first case the different groups constitute the separate data blocks, in the second case the different subjects.

The question then rises whether or not the same structure underlies each data block. For example, in emotion psychology, there has been a long-lasting debate about the structure of emotions (e.g., Ekman, 1999; Fontaine, Scherer, Roesch, & Ellsworth, 2007; Russel & Barrett, 1999). In this debate, many cross-cultural psychologists argue that the structure may differ between cultures (Eid & Diener, 2001; Fontaine, Poortinga, Setiadi, & Markam, 2002; MacKinnon & Keating, 1989; Rodriguez & Church, 2003). In a similar vein, between subjects, structural differences in the time-varying experience of emotions can be expected. For example, studies have indicated that subjects differ with respect to so-called emotional granularity (Barrett, 1998; Tugade, Fredrickson, & Barret, 2004). For subjects with a low emotional granularity negative or positive emotional states are highly correlated across time. Subjects with high emotional granularity experience emotions in a more differentiated manner, distinguishing between a variety of positive and negative emotions, rather than feeling overall positive or negative.

To explore the similarities and differences in the structure of multivariate multiblock data, two approaches have been proposed within the component analysis literature. A first

approach is to perform a separate principal component analysis (PCA; Jolliffe, 1986; Meredith & Millsap, 1985; Pearson, 1901) on each data block, which summarizes the information in the data block by reducing the variables to a few components. For instance, when analyzing cross-cultural emotion data, where each data block holds the emotion scores of inhabitants of a particular country, a PCA is conducted for each country separately (e.g., Fontaine et al., 2002). In this analysis, each data block \mathbf{X}_i is decomposed into a component score matrix \mathbf{F}_i , containing the scores of the inhabitants on the components, and a loading matrix \mathbf{B}_i , indicating the extent to which the scores on the emotion variables are determined by the respective components. Because a separate loading matrix is obtained for each data block, this approach leaves plenty of freedom to trace differences in the underlying structure of the different data blocks. To gain further insight, van de Vijver and Leung (1997) proposed to rotate the loading matrices of the different data blocks towards each other and calculate Tucker congruence coefficients (Tucker, 1951) between them. However, when the loadings differ substantially, such coefficients do not reveal potential similarities across data blocks. Alternatively, one may compare the component loadings of the variables directly. Because of the large amount of information this is often not very insightful either. Especially when the number of data blocks becomes large, it is practically infeasible to trace differences and similarities.

A second approach is simultaneous component analysis (SCA; Millsap & Meredith, 1988; Kiers, 1990; Kiers & ten Berge, 1994; Timmerman & Kiers, 2003; Van Deun, Smilde, van der Werf, Kiers, & Van Mechelen, 2009). In an SCA analysis all data blocks are reduced simultaneously, based on the assumption that the same components underlie the different data blocks and thus that the same loading matrix can be used to reconstruct these data. Applying SCA to our example from cross-cultural emotion research, one common loading matrix

would be obtained for all cultures. Since a common loading matrix is imposed, this approach is much more parsimonious than performing a PCA on each data block separately. However, this approach often leaves little room to find structural differences between data blocks.

One may conclude that both PCA and SCA may fail to yield insight, because in the separate PCA approach structural similarities are hard to trace and in SCA structural differences may be difficult to detect. Therefore, we introduce a new generic modeling strategy called Clusterwise SCA, which comprises the separate PCA approach and the SCA approach as special cases. Clusterwise SCA is based on the same principle as clusterwise linear regression, namely, that the data contain a few clusters, where each cluster has a different underlying structure (Brusco, Cradit, Steinley, & Fox, 2008; DeSarbo, Oliver, & Rangaswamy, 1989; Späth, 1979, 1982). More specifically, the key idea behind Clusterwise SCA is that the different data blocks form a limited number of mutually exclusive clusters, where data blocks that belong to the same cluster can be reconstructed by means of the same loadings, whereas data blocks that belong to different clusters have a different underlying component structure and thus imply a different loading matrix. In this paper, the data blocks within each cluster are modeled using the SCA variant that imposes Equal average Cross-Products constraints (SCA-ECP; Timmerman & Kiers, 2003). Therefore, the proposed modeling strategy will be called Clusterwise SCA-ECP.

The remainder of this paper is organized as follows: In the following section, the separate PCA and SCA-ECP approaches are recapitulated, followed by the introduction of the new Clusterwise SCA-ECP model and a comparison of the latter model to other approaches for detecting structural differences and similarities in multivariate multiblock data. Next, the Data Analysis section describes the aim of and an algorithm for Clusterwise SCA-ECP analysis and proposes a model selection procedure. In the fourth section, the

performance of the algorithm is evaluated in three simulation studies. Then, the model is illustrated with an application to time series data from eating disorder patients, and an application to data from social psychology, with subjects being nested in experimental conditions. Finally, directions for future research are discussed.

Model

Data Structure and Preprocessing

Clusterwise SCA-ECP can be applied to all kinds of multivariate multiblock data, where ‘multivariate’ indicates that multiple variables are involved and ‘multiblock’ implies that the data can be divided in separate data blocks according to the hierarchical structure of the data. Different configurations are possible: for instance, the data blocks may represent different groups (e.g., cultures) where the observations within the data blocks stem from different individuals that belong to these groups, or the data blocks may represent individuals with different time points constituting the observations within the blocks (i.e., multioccasion-multisubject data). More formally, Clusterwise SCA-ECP requires I data blocks \mathbf{X}_i ($N_i \times J$) that contain scores on J variables, where the number of observations N_i ($i = 1, \dots, I$) in each data block may differ between data blocks, subject to the restriction that N_i is larger than the number of components to be fitted (but preferably larger than J , for the sake of stable model estimates). These I data blocks can be concatenated into a N (observations) \times J (variables) data matrix \mathbf{X} , where $N = \sum_{i=1}^I N_i$. A graphical presentation of the data structure is given in

Figure 1.

[Insert Figure 1 about here]

Note that three-way three-mode data (for an introduction, see Kroonenberg, 2008) are a special case of multivariate multiblock data, in which all the groups consist of the same subjects, or in which all the subjects are measured on the same time points. Moreover, note that a distinction can be made between ‘multigroup’ and ‘multilevel’ multiblock data, based on whether the data blocks are considered fixed or random respectively (Timmerman, Kiers, Smilde, Ceulemans, & Stouten, 2009). For instance, if each data block contains the scores of a single subject, the data blocks are fixed when one is only interested in the subjects in the study, and random when one wants to generalize the conclusions towards a larger populations of subjects. As Clusterwise SCA-ECP is a deterministic method (i.e., no distributional assumptions are made), it is applicable to both multigroup and multilevel data.

Clusterwise SCA-ECP is designed for tracing similarities and differences in the within structure of the data (i.e., the correlational structure within each of the data blocks). Therefore, the differences between the data blocks in the means of the variables, which is often called the between structure, should be removed from the data. This is achieved by centering the variables per data block. Moreover, arbitrary differences between the variables in measurement scale are usually eliminated in component analysis by scaling the data. In the case of multivariate multiblock data, two scaling options have been advocated: The first option is autoscaling (Kiers & ten Berge, 1994), which implies that each variable is rescaled per data block, such that the sum of squares per variable is equal to the number of observations N_i for the data block in question. This leads to a variance of one per variable for each data block. Note that, when the centering per data block is combined with autoscaling, the preprocessing is equivalent to calculating z-scores within each data block. The second option is to rescale across the data blocks (Timmerman & Kiers, 2003), implying a variance of one per variable across all data blocks, so that differences between the data blocks in

variability are preserved. In this paper we have chosen the first option, to focus on the correlational differences between the data blocks rather than the differences in variances. In what follows, we assume each data block \mathbf{X}_i to be centered and autoscaled.

Principal Component Analysis on Each of the Data Blocks Separately

Applying PCA to each data block \mathbf{X}_i ($i = 1, \dots, I$) separately, implies that the underlying structure is allowed to differ across data blocks in that the data blocks may be characterized by different loading matrices \mathbf{B}_i . Formally, this model can be written as follows:

$$\mathbf{X}_i = \mathbf{F}_i \mathbf{B}_i' + \mathbf{E}_i \quad (1)$$

where \mathbf{F}_i ($N_i \times Q$) denotes the component score matrix containing the Q component scores for observation 1, ..., N_i where the number of components Q is assumed to be the same across data blocks, \mathbf{B}_i ($J \times Q$) denotes the loading matrix for the i -th data block, and \mathbf{E}_i ($N_i \times J$) denotes the matrix of residuals. Note that, since the variables are standardized per data block, the loadings in \mathbf{B}_i equal the correlations between the respective components and variables, in case of orthogonal components.

The loading matrices \mathbf{B}_i ($i = 1, \dots, I$) of the PCA solutions may be orthogonally or obliquely rotated, provided that such a transformation is compensated for in the component score matrices \mathbf{F}_i . Therefore, standard rotational procedures (e.g., Varimax; Kaiser, 1958) can be applied to obtain solutions which are easier to interpret. When an oblique rotation is used, the components become correlated to some extent. Consequently, the loadings may not be interpreted as correlations, but as weights that indicate the extent to which each variable is influenced by the respective components.

Simultaneous Component Analysis

SCA differs from the separate PCA approach in that the underlying components are assumed to be the same across the data blocks, implying that the loadings of the variables are identical. Specifically, the SCA-ECP model (Kiers & ten Berge, 1994; Timmerman & Kiers, 2003) is given by

$$\mathbf{X}_i = \mathbf{F}_i \mathbf{B}' + \mathbf{E}_i \quad (2)$$

where \mathbf{F}_i ($N_i \times Q$) denotes the component score matrix of the i -th data block, \mathbf{B} ($J \times Q$) denotes the loading matrix which is identical for all data blocks and therefore does not have an index i , and \mathbf{E}_i ($N_i \times J$) denotes the matrix of residuals. In the SCA-ECP model, the component score matrices \mathbf{F}_i are constrained in that the variances and the correlations of the component scores in \mathbf{F}_i are restricted to be equal across data blocks (with all the component correlations being zero if the components are orthogonal), i.e., $N_i^{-1} \mathbf{F}_i' \mathbf{F}_i = \mathbf{\Phi}$. To partly identify the solution, the variances of the component scores are fixed at 1, implying that the diagonal elements of $\mathbf{\Phi}$ are set to 1.

Note that Timmerman and Kiers (2003) also described three less restrictive variants of the SCA model, for which the variances and/or correlations of the component scores in \mathbf{F}_i are allowed to vary across data blocks. Differences in the component variances and/or correlations may reveal some structural differences between the data blocks. For instance, when some of the data blocks have hardly any variance on a certain component, this may indicate that the component concerned does not underlie those data blocks. In that case, a rotation can be applied to distinguish such 'distinctive' components from the common components more clearly (see DISCO-SCA; Schouteden, Van Deun, Van Mechelen &

Pattyn, 2011). However, in practice, differences in component variances (or correlations) are usually gradual rather than distinct, making it hard to deduce structural differences between the data blocks. On top of that, SCA-ECP imposes equality of the component variances and correlations across data blocks and thus rules out tracing structural differences. In the next section, we will motivate why we selected the SCA-ECP variant for the Clusterwise SCA.

As is the case for PCA solutions, the components of an SCA-ECP solution can be freely rotated without altering the fit of the solution. Specifically, the loading matrix \mathbf{B} can be transformed by multiplying it by any rotation matrix, provided that such a transformation is compensated for in the component score matrices \mathbf{F}_i ($i = 1, \dots, D$).

Clusterwise SCA-ECP

As stated in the Introduction, the key idea behind Clusterwise SCA-ECP is that the different data blocks fall apart into K mutually exclusive clusters. The data blocks that belong to the same cluster are driven by the same processes, whereas data blocks that belong to different clusters come about through different mechanisms. Clusterwise SCA-ECP deals with these differences in underlying structure by partitioning the I data blocks into K clusters and modeling the data blocks within each cluster by an SCA-ECP model, assuming the number of components Q to be the same across the clusters.

In this paper, the goal is to capture the structural differences between the data blocks with the clustering. To ensure that data blocks with a different within structure will be allocated to different clusters, the SCA-ECP model is imposed for each cluster. Any alternative SCA model would leave room for structural differences between the data blocks within a cluster, like, for example, differences in variances of the component scores.

Formally, the model equation of Clusterwise SCA-ECP is given by

$$\mathbf{X}_i = \sum_{k=1}^K p_{ik} \mathbf{F}_i^{(k)} \mathbf{B}^{(k)'} + \mathbf{E}_i = \sum_{k=1}^K p_{ik} \mathbf{F}_i \mathbf{B}^{(k)'} + \mathbf{E}_i \quad (3)$$

where K is the number of clusters, p_{ik} denotes the entries of the binary partition matrix \mathbf{P} ($I \times K$) which equal 1 when data block i is assigned to cluster k and 0 otherwise, $\mathbf{F}_i^{(k)}$ ($N_i \times Q$) denotes the component score matrix of data block i when assigned to cluster k , $\mathbf{B}^{(k)}$ ($J \times Q$) denotes the loading matrix of cluster k ($k = 1, \dots, K$) and \mathbf{E}_i ($N_i \times J$) denotes the matrix of residuals. The loading matrices $\mathbf{B}^{(k)}$ are shared by all data blocks that belong to a particular cluster. As can be seen in the left part of Equation 3, the index k in $\mathbf{F}_i^{(k)}$ is mostly omitted in the remainder of this paper, because each data block \mathbf{X}_i is assigned to one cluster only. Note that since the parameter estimates of an SCA-ECP solution have rotational freedom, the components of a Clusterwise SCA-ECP solution can also be freely rotated within each cluster without altering the fit of the solution.

[Insert Table 1 about here]

To illustrate the Clusterwise SCA-ECP model, we make use of the hypothetical data matrix \mathbf{X} in Table 1. \mathbf{X} contains the scores of 4 subjects on 6 variables pertaining to emotions and physical activity at 8, 9, 7, and 10 measurement occasions respectively. This data matrix can be perfectly reconstructed by the Clusterwise SCA-ECP solution with two clusters and two components, of which the partition matrix, component loading matrices and component score matrices are presented in Tables 2, 3 and 4, respectively.

[Insert Table 2, Table 3 and Table 4 about here]

From Table 2, it can be read that the 4 subjects fall apart into two clusters, where subjects 1 and 4 belong to the first cluster and subjects 2 and 3 to the second. To trace the

structural differences between these two clusters, we inspect the cluster loading matrices $\mathbf{B}^{(1)}$ and $\mathbf{B}^{(2)}$ in Table 3. For both clusters, a positive affect component and a negative affect component are found. The difference between the two clusters lies in how positive and negative emotions are related to physical activity. Specifically, $\mathbf{B}^{(1)}$ shows that for the subjects in cluster 1 there is a relation between physical activity and positive affect with both being unrelated to negative affect, whereas $\mathbf{B}^{(2)}$ reveals that for the subjects of cluster 2 being physically active is related to negative affect and not to positive affect. From the component score matrices \mathbf{F}_i in Table 4, it can be derived how positive or negative the subjects feel on the different measurement occasions. Note that the scores in Table 1 can be reconstructed by combining the matrices in Tables 2, 3, and 4 according to Equation 3. For example, from Table 1 it be read that subject 1 has a score of -1.4 on the variable ‘happy’. By means of Equation 3, we can reconstruct this score as follows: $[-1.4 \quad -1.1] * [1 \quad 0]'$, where * denotes the inner product of the vectors, $[-1.4 \quad -1.1]$ are the component scores of the first observation of the subject in question and $[1 \quad 0]$ are the loadings of ‘happy’ for the first cluster, to which the subject belongs.

Finally, from Equation 3 it is clear that the Clusterwise SCA-ECP model is a generic modeling strategy for tracing differences in underlying structure, in that it comprises the separate PCA and SCA-ECP approaches as special cases. Specifically, when the number of clusters K is set to one, the model reduces to a regular SCA-ECP model. In case the number of clusters K equals the number of data blocks I , the model boils down to a separate PCA on each data block.

Relations to Existing Models

In the past decades, several models and associated algorithms have been developed for tracing differences in the underlying structure of multivariate data. In this section, we focus on methods for multiblock data, discarding a number of techniques that combine clustering and dimension reduction in the context of single block data – reduced K-means (Bock, 1987; de Soete & Carroll, 1994), factorial K-means (Vichi & Kiers, 2001; Timmerman, Ceulemans, Kiers, & Vichi, 2010), mixtures of factor analyzers (McLachlan & Peel, 2000), mixture SEM (Dolan & van der Maas, 1998; Jedidi, Jagpal, & DeSarbo, 1997; Yung, 1997), high-dimensional data clustering (Bouveyron, Girard, & Schmid, 2007) – and three-way three-mode data – Tucker3 clustering (Rocci & Vichi, 2005), three-way factorial K-means (Vichi, Rocci & Kiers, 2007). All existing techniques differ from Clusterwise SCA-ECP in at least one of the following respects: Is the method exploratory or confirmatory? Is the model deterministic or stochastic, i.e. are assumptions made about the underlying probability distributions? What kind of data is required: two-mode data or one-mode data? What type of data reduction is performed?

Within the family of deterministic models, Clusterwise SCA-ECP is related to points of view analysis (PVA; Tucker & Messick, 1963). This method was proposed to handle multiple one-mode (dis)similarity data blocks. The data blocks are clustered and for each cluster multidimensional scaling (MDS) is performed (Kruskal & Wish, 1978). It can be concluded that Clusterwise SCA-ECP differs from this method with respect to the kind of data that is dealt with and the type of data reduction, as Clusterwise SCA-ECP deals with two-mode instead of one-mode data blocks and uses SCA-ECP instead of MDS for the data reduction.

Within the stochastic framework, Clusterwise SCA-ECP is mainly related to multigroup structural equation modeling (multigroup SEM; Jöreskog, 1971; Kline, 2004;

Sörbom, 1974) and multigroup exploratory factor analysis (multigroup EFA; Dolan, Oort, Stoel, & Wicherts, 2009; Hessen, Dolan, & Wicherts, 2006) on the one hand, and multilevel latent class analysis (Vermunt, 2003, 2008a) and multilevel mixture factor analysis (Varriale & Vermunt, 2010; Vermunt, 2008b) on the other hand.

Multigroup SEM is a confirmatory method for multiblock data which imposes a structural equation model (Haavelmo, 1943; Kline, 2004) for each data block. Similarly, multigroup EFA is an exploratory method, in which an exploratory common factor model (Lawley & Maxwell, 1962) is estimated for each data block. The factor loading structure can be constrained to be the same for each data block (like in SCA) or different for each data block (like in separate PCA). To determine whether or not the factor loading structure is invariant over the blocks, the fit of the constrained and unconstrained model are compared. Thus, multigroup SEM and multigroup EFA differ from Clusterwise SCA-ECP in that these methods do not perform a clustering of the data blocks. Moreover, in multigroup SEM and multigroup EFA, it is assumed that the observations within a particular data block are independent. Therefore, modeling multioccasion-multisubject data with multigroup SEM or EFA would not be correct, due to the serial dependencies between successive measurements. Thus, Clusterwise SCA-ECP has as major advantage over multigroup SEM and EFA that no assumptions about within-block dependencies are made. When these dependencies exist, they may be captured in the component scores.

Multilevel latent class analysis (MLCA) as well as multilevel mixture factor analysis (MMFA) imply a clustering of the data blocks. However, a crucial difference between the latter techniques and Clusterwise SCA-ECP is that in MLCA and MMFA the clustering is based on between-block differences in variable means, thus capturing the between structure,

while in Clusterwise SCA-ECP the clustering of the data blocks is performed according to between-block differences in within-block structure.

From this overview of related models, it can be concluded that no other model is available for two-mode multivariate multiblock data that clusters the data blocks based on the within-block structure.

Data Analysis

Aim of Clusterwise SCA-ECP Analysis

For a given number of clusters K and components Q , the aim of a Clusterwise SCA-ECP analysis is to find the partition matrix \mathbf{P} , the component score matrices \mathbf{F}_i and the loading matrices $\mathbf{B}^{(k)}$ that minimize the loss function:

$$L = \sum_{i=1}^I \sum_{k=1}^K p_{ik} \|\mathbf{X}_i - \mathbf{F}_i \mathbf{B}^{(k)'}\|^2 \quad (4)$$

Note that on the basis of the loss function value L , the percentage of variance accounted for (VAF) can be computed as follows:

$$\text{VAF}(\%) = \frac{\|\mathbf{X}\|^2 - L}{\|\mathbf{X}\|^2} \times 100. \quad (5)$$

Algorithm

Ideally, one would wish to develop a Clusterwise SCA-ECP algorithm that returns a globally optimal solution. To this end, the algorithm should sieve through all possible partitions of the data blocks in a smart way. However, for the related K-means clustering

problem (MacQueen, 1967), this becomes very time consuming if the data set contains a large number of observations (Aloise, Hansen, & Liberti, in press; Brusco, 2006), because the number of possible partitions grows very large. Our problem requires additional computations since data reduction is performed within each cluster, thus a search of all possible partitions seems not feasible for Clusterwise SCA-ECP. Therefore, we developed a fast relocation algorithm that can handle large numbers of data blocks, but may end in a local minimum.

More specifically, to obtain a (K, Q) Clusterwise SCA-ECP solution that minimizes the loss function (Equation 4), an alternating least squares (ALS) procedure is used, that consists of four steps:

1. Randomly initialize partition matrix \mathbf{P} : Initialize the partition matrix \mathbf{P} by randomly assigning the I data blocks to one of the K clusters, where each data block has an equal probability of being assigned to each cluster. If one of the clusters is empty, repeat this procedure until all clusters contain at least one element.
2. Estimate the SCA-ECP model for each cluster: Estimate the \mathbf{F}_i and $\mathbf{B}^{(k)}$ matrices for each cluster $k = 1, \dots, K$ by performing a rationally started SCA-ECP analysis (Timmerman & Kiers, 2003) on the \mathbf{X}_i data blocks assigned to the k -th cluster. To this end, the loading matrix $\mathbf{B}^{(k)}$ is rationally initialized, based on the singular value decomposition of the vertical concatenation of the data blocks within cluster k , denoted by $\mathbf{X}^{(k)}$. Next, an ALS procedure is performed, in which $\mathbf{F}^{(k)}$ and $\mathbf{B}^{(k)}$ are iteratively re-estimated, where $\mathbf{F}^{(k)}$ is the vertical concatenation of the component scores of all data blocks within cluster k . Specifically, $\mathbf{F}^{(k)}$ is (re-)estimated by performing a singular value decomposition for all data blocks \mathbf{X}_i that belong to cluster k : $\mathbf{X}_i \mathbf{B}^{(k)}$ is decomposed into \mathbf{U}_i , \mathbf{S}_i and \mathbf{V}_i with $\mathbf{X}_i \mathbf{B}^{(k)} = \mathbf{U}_i \mathbf{S}_i \mathbf{V}_i'$; a least squares

estimate of $\mathbf{F}_i^{(k)}$ is given by $\mathbf{F}_i^{(k)} = \sqrt{N_i} \mathbf{U}_i \mathbf{V}_i'$ (ten Berge, 1993). $\mathbf{B}^{(k)}$ is updated by $\mathbf{B}^{(k)} = ((\mathbf{F}^{(k)'} \mathbf{F}^{(k)})^{-1} \mathbf{F}^{(k)'} \mathbf{X}^{(k)})'$. This way, cluster loading matrices $\mathbf{B}^{(k)}$ are obtained that 'average' the underlying structure of the data blocks within a cluster. Note that, due to the SCA-ECP restrictions on the component variances and correlations of each $\mathbf{F}_i^{(k)}$, the $\mathbf{F}^{(k)}$ and $\mathbf{B}^{(k)}$ matrices for each cluster k cannot be obtained directly by means of a singular value decomposition of $\mathbf{X}^{(k)}$.

3. Re-estimate the partition matrix \mathbf{P} : For each data block, a component score matrix $\tilde{\mathbf{F}}_i^{(k)}$ is computed for each cluster k , by means of the same computations that are used for updating the component score matrices within SCA-ECP (see step 2). The extent to which the data block fits in the different clusters is quantified by computing the following block and cluster specific partition criterion: $L_{ik} = \left\| \mathbf{X}_i - \tilde{\mathbf{F}}_i^{(k)} \mathbf{B}^{(k)'} \right\|^2$. Each data block is assigned to the cluster k for which L_{ik} is minimal. When one of the K clusters is empty after this procedure, the data block with the worst fit in its current cluster, is moved to the empty cluster.
4. Steps 2 and 3 are repeated until convergence is reached, i.e., until the decrease of the loss function value L (Equation 4) for the current iteration is smaller than the convergence criterion of 1e-6.

To reduce the probability to end up in a local minimum, it is advised to use a multistart procedure with different random initializations of the partition matrix \mathbf{P} . In the algorithm, a rational start is used for the SCA-ECP analysis within each cluster (see step 2 of the algorithm). The results of a pilot study reveal that the use of multiple random starts rather than a single rational start does not improve the performance of the Clusterwise SCA-ECP

algorithm, and therefore we deem the use of multiple random starts to be unnecessary. The described algorithm has been implemented in Matlab R2010a and can be obtained freely from the first author. Moreover, user-friendly software for applying Clusterwise SCA-ECP has already been developed (De Roover, Ceulemans, & Timmerman, in press) and is available for potential users at: <http://ppw.kuleuven.be/okp/software/MBCA/>. It is also possible to run the software as a stand-alone application (i.e., without Matlab).

Model Selection

When applying Clusterwise SCA-ECP analysis, the underlying number of clusters K and components Q is usually unknown. To deal with this, Clusterwise SCA-ECP solutions are estimated using several values for K and Q . Subsequently, a solution is selected on the basis of formal model selection techniques, interpretability and stability of the solutions. As a formal model selection technique, one may consider to use a generalization of the well-known scree test (Cattell, 1966) which aims at selecting a model with an optimal balance between fit and parsimony. Specifically, one may plot the percentage of variance accounted for (Equation 5) of the different solutions against the number of components for each value of K . The number of clusters K is established, by examining the general increase in fit that is obtained by adding a cluster and choosing the number of clusters after which this general increase in fit levels off. Finally, considering the solutions with K clusters only, the number of components Q is determined for which it holds that adding more components does not ‘significantly’ increase the fit to the data. The use of this model selection procedure will be further illustrated below when applying it to empirical data sets. To evaluate the stability of the retained solution, one can examine how often (specific aspects of) the obtained partition show up across the different random starts in the multistart procedure.

Simulation Studies

In this section, we first present a large simulation study in which the performance of the Clusterwise SCA-ECP algorithm is evaluated when the correct number of clusters and components is known. Next, two smaller simulation experiments are discussed in which we investigate how the performance is affected by less favorable analysis conditions. Specifically, we examine the effect of using an incorrect number of clusters or components and the effect of the presence of masking variables in the data (i.e., variables in which the cluster structure is not reflected; see e.g., Brusco & Cradit, 2001; Vichi & Kiers, 2001).

Simulation Study 1

Problem.

The first simulation study is an extensive study in which the Clusterwise SCA-ECP algorithm is evaluated with respect to goodness of fit, sensitivity to local minima and goodness of recovery, under optimal conditions, i.e., when the data to be analyzed are generated from a Clusterwise SCA-ECP model with a known number of clusters K and components Q . Furthermore, we examine whether the performance of the algorithm is influenced by seven factors: (1) the number of data blocks, (2) the number of observations per data block, (3) the number of underlying clusters and (4) components, (5) the cluster size, (6) the amount of error on the data and (7) the structure of the loading matrices of the different clusters. These factors were chosen because their influence is often investigated in simulation studies on clustering algorithms and/or because they are theoretically interesting. Factors 1 and 2 pertain to sample size. Based on previous research (Brusco & Cradit, 2005;

Hands & Everitt, 1987), we expect that the Clusterwise SCA-ECP algorithm will perform better when more information is available (i.e., more data blocks and/or more observations per data block). With respect to Factors 3 and 4, which refer to the complexity of the underlying model, we hypothesize that the goodness of fit and recovery will decrease with increasing complexity (Brusco & Cradit, 2005; Milligan, Soon, & Sokol, 1983; Timmerman et al., 2010). Regarding Factor 5, cluster size, we conjecture that better results will be obtained when the clusters are of equal size (Brusco & Cradit, 2001; Milligan et al., 1983; Steinley, 2003). For Factor 6, amount of error, we expect that performance will deteriorate when the data contain more error (Brusco & Cradit, 2005). Finally, Factor 7, structure of the cluster loading matrices, was manipulated to study whether the performance deteriorates when the structures underlying the different clusters overlap more, where overlap is defined in terms of congruence between the loadings of the different clusters. Moreover, we expect ordinal interactions between (some of) these factors: Specifically, we conjecture that the worst results will be obtained with small data sets that contain, apart from a lot of error, many clusters of different sizes, of which the cluster loading matrices are highly congruent.

Design and procedure.

In this simulation study, the number of variables J was fixed at 12. Furthermore, the seven factors that were introduced above, were systematically varied in a complete factorial design:

1. the *number of data blocks* I at 2 levels: 20, 40;

2. the *number of observations per data block* N_i at 3 levels: $N_i \sim U[15;20]$, $N_i \sim U[30;70]$, $N_i \sim U[80;120]$, with U indicating a discrete uniform distribution between the given numbers;
3. the *number of clusters* K at 3 levels: 2, 3, 4;
4. the *number of components* Q at 3 levels: 2, 3, 4;
5. the *cluster size*, at 3 levels (see Milligan et al., 1983): equal (equal number of data blocks in each cluster); unequal with minority (10% of the data blocks in one cluster and the remaining data blocks distributed equally over the other clusters); unequal with majority (60% of the data blocks in one cluster and the remaining data blocks distributed equally over the other clusters);
6. the *error level* e , which is the expected proportion of error variance in the data blocks \mathbf{X}_i : .00, .20, .40.
7. the *structure of the loading matrices* $\mathbf{B}^{(k)}$ of the different clusters at 3 levels: simple structure, random with low congruence and random with high congruence.

With respect to the latter factor, at the simple structure level, the cluster loading matrices all had simple structure, with the set of variables constituting the different components varying across the different clusters. For instance, the cluster loading matrices for K and Q equal to four were constructed as follows:

$$\mathbf{B}^{(1)} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad \mathbf{B}^{(2)} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad \mathbf{B}^{(3)} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad \mathbf{B}^{(4)} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \end{bmatrix}$$

Random cluster loading matrices with low congruence were obtained by sampling the loadings uniformly between -1 and 1. Finally, random cluster loading matrices with high congruence were constructed as follows: First, a common base matrix was uniformly sampled between -1 and 1. Subsequently, the rows of this matrix were rescaled to have a sum of squares equal to .7. Next, for each cluster a new random matrix, uniformly sampled between -1 and 1 and rescaled to have a rowwise sum of squares equal to .3, was added to the base matrix to form the cluster loading matrix.

To evaluate how much the resulting cluster loading matrices differ in each level of factor 7, they were orthogonally procrustes rotated to each other (i.e., for each pair of cluster loading matrices, one was chosen to be the target matrix and the other was rotated towards the target matrix) and a congruence coefficient φ^1 (Tucker, 1951) was computed, for each pair of corresponding components in all pairs of $\mathbf{B}^{(k)}$ matrices. Subsequently, a grand mean of

¹ The congruence coefficient (Tucker, 1951) between two column vectors \mathbf{x} and \mathbf{y} is defined as their normalized inner product: $\varphi_{\mathbf{xy}} = \frac{\mathbf{x}'\mathbf{y}}{\sqrt{\mathbf{x}'\mathbf{x}}\sqrt{\mathbf{y}'\mathbf{y}}}$.

the obtained φ -values was calculated, over the components and cluster pairs. The resulting φ -values led to the conclusion that the differences between the cluster loading matrices are large at the random, low congruence level (average φ across data sets = .41, $SD = 0.09$), rather small at the simple structure level (average $\varphi = .71$, $SD = 0.07$), and very small at the random, high congruence level (average $\varphi = .93$, $SD = 0.02$). Therefore, the levels of this factor can also be labeled ‘low congruence’, ‘medium congruence’ and ‘high congruence’, respectively.

For each cell of the factorial design, 50 data matrices \mathbf{X} were generated, consisting of $I \mathbf{X}_i$ data blocks. These \mathbf{X}_i data blocks were constructed as follows:

$$\mathbf{X}_i = \mathbf{F}_i \mathbf{B}^{(k)'} + \mathbf{E}_i \quad (6)$$

where the component score matrix \mathbf{F}_i and error matrix \mathbf{E}_i were generated by randomly sampling entries from a standard normal distribution. The partition matrix \mathbf{P} was generated by first computing the size of the different clusters and then randomly assigning the correct number of data blocks to the clusters. The cluster loading matrices $\mathbf{B}^{(k)}$ were generated as described above (factor 7). Subsequently, the error matrices \mathbf{E}_i and the cluster loading matrices $\mathbf{B}^{(k)}$ were rescaled to obtain data that contain a proportion e of error variance (factor 6). Finally, the resulting \mathbf{X}_i matrices were standardized columnwise, and were vertically concatenated into the matrix \mathbf{X} .

In total, 2 (number of data blocks) \times 3 (number of observations per data block) \times 3 (number of clusters) \times 3 (number of components) \times 3 (cluster size) \times 3 (error level) \times 3 (structure of cluster loading matrices) \times 50 (replicates) = 72,900 simulated data matrices were generated. Each data matrix \mathbf{X} was analyzed with the Clusterwise SCA-ECP algorithm, using

the correct number of clusters K and components Q . The algorithm was run 25 times, each time using a different random start, and the best solution was retained.

Results.

Goodness of fit and sensitivity to local minima.

Goodness of fit.

On average the VAF (Equation 5) equals 99% ($SD = 1.18$), 84% ($SD = 1.54$) and 69% ($SD = 2.74$) when 0%, 20% and 40% error variance is present in the data (factor 6), respectively. Note that in the 0% error condition, the means and standard deviations of the simulated component score matrices will differ from zero and one, due to sampling fluctuations, implying that the simulated data do not perfectly comply with the SCA-ECP assumptions. This explains why these data sets cannot be fitted perfectly. The results for the 20% and 40% error conditions indicate that part of the error variance is fitted and thus that overfitting occurs. To gain more insight into when overfitting occurs, the amount of overfit was computed as the difference between the VAF by the estimated model and the VAF by the true model:

$$\text{overfit (\%)} = \text{VAF}_{\mathbf{M}} - \text{VAF}_{\mathbf{T}} = \frac{\left(1 - \sum_{n=1}^N \sum_{j=1}^J \|\mathbf{X} - \mathbf{M}\|^2\right) - \left(1 - \sum_{n=1}^N \sum_{j=1}^J \|\mathbf{X} - \mathbf{T}\|^2\right)}{N \times J} \times 100 \quad (7)$$

where \mathbf{T} is the true data matrix (i.e., the data before error was added) and \mathbf{M} contains the reconstructed scores of the best solution out of the 25 random runs. Note that \mathbf{X} always differs from \mathbf{T} , even in those cases where no error was added, because \mathbf{X} , unlike \mathbf{T} , is a preprocessed matrix. For all 72,900 data sets the overfit value was larger than zero, implying

that overfitting always occurs. Subsequently, an analysis of variance was performed with overfit as the dependent variable and the seven factors as independent variables. Finally, to examine which main and interaction effects have a large effect size, omega squared proportions of total variance $\hat{\omega}^2$ (Hays, 1963; Olejnik & Algina, 2000) were computed. A large omega squared was found for the main effect of the amount of error variance ($\hat{\omega}^2 = .68$), which implies that the error variance accounts for 68% of the variance in the overfit value. As expected, the overfit increases when more error variance is present in the data (Figure 2a). The amount of error further interacts with the number of components ($\hat{\omega}^2 = .09$): the more components, the more overfitting (see Figure 2a). Finally, a main effect of the number of observations per data block ($\hat{\omega}^2 = .14$) was revealed, implying that the amount of overfit is larger when less observations are available (Figure 2b). Other effects are not discussed, because they account for less than 7.5% of the variance of the dependent variable.

[Insert Figure 2 about here]

Sensitivity to local minima.

The sensitivity of the Clusterwise SCA-ECP algorithm to local minima should be evaluated by comparing the loss function value of the retained solution to that of the global minimum. However, the global minimum of the Clusterwise SCA-ECP analysis of a data set \mathbf{X} is unknown, because the simulated data do not perfectly comply with the Clusterwise SCA-ECP assumptions, as was explained before.

Alternatively, we first evaluated whether the best fitting solution out of the 25 solutions from the multistart procedure yielded a higher loss function than the solution

resulting from seeding the algorithm with the true \mathbf{F}_i , $\mathbf{B}^{(k)}$ and \mathbf{P} . If this were the case, the multistart solution is a local minimum for sure. Such a local minimum was found for 1,230 out of the 72,900 simulated data matrices (1.69%). Out of these 1,230 data sets 1,206 belong to the condition in which the cluster loading matrices are random with high congruence. Moreover, the majority (1,199) of these 1,206 data sets contain a lot of error variance (i.e., $e = .40$) and/or have a low number of observations per data block (i.e., between 15 and 20).

To investigate further the issue of local minima, we examined the proportion of random runs that had a loss function value that was equal to that of a proxy of the global minimum. This proportion will be called ‘global minimum proportion’ in the remainder of the paper. Note that the proxy is the best fitting solution resulting from seeding the algorithm with the true \mathbf{F}_i , $\mathbf{B}^{(k)}$ and \mathbf{P} on the one hand and from the 25 random runs on the other hand. On average, the global minimum proportion equals .81 ($SD = 0.26$). An analysis of variance was performed with the global minimum proportion as the dependent variable, which revealed main effects of the number of observations per data block, of the number of clusters K and of the type of cluster loading matrices, and an interaction effect of the number of observations per data block and the type of cluster loading matrices. The effect of the number of observations per data block ($\hat{\omega}^2 = .11$) implies that having less observations per data block leads to a lower global minimum proportion (Figure 3b). The effect of the number of clusters ($\hat{\omega}^2 = .17$) implies that the more clusters are present in the data, the lower the global minimum proportion (Figure 3a). With respect to the effect of the amount of congruence between the cluster loading matrices ($\hat{\omega}^2 = .19$), the global minimum proportion is lower for the random loading matrices with a high congruence (Figure 3b). The interaction of the latter effect with the number of observations per data blocks ($\hat{\omega}^2 = .09$) implies that the effect of

the congruence between the cluster loading matrices is stronger when the number of observations per data block is lower, and vice versa (Figure 3b).

[Insert Figure 3 about here]

Goodness of recovery.

The goodness of recovery will be evaluated with respect to (a) the clustering of the data blocks and (b) the cluster loading matrices.

Recovery of the clustering of the data blocks.

To examine the goodness of recovery of the cluster membership of the data blocks, the *Adjusted Rand Index (ARI)* (Hubert & Arabie, 1985) is computed between the true partition of the data blocks and the estimated partition. The *ARI* equals 1 if the two partitions are identical, and equals 0 when the overlap between the two partitions is at chance level.

With an overall mean *ARI* of .98 ($SD = 0.11$) the Clusterwise SCA-ECP algorithm appears to recover the clustering of the data blocks very well. An analysis of variance was performed with *ARI* as the dependent variable and the seven factors as independent variables. An interaction effect between the congruence of the cluster loading matrices and the number of observations per data block ($\hat{\omega}^2 = .12$) was found, which implies that the *ARI* is lower when the congruence between the cluster loading matrices is high, given that the number of observations per data block is low (i.e., between 15 and 20; Figure 4).

[Insert Figure 4 about here]

Recovery of the cluster loading matrices.

To evaluate the recovery of the cluster loading matrices, we obtained a goodness-of-cluster-loading-recovery statistic (*GOCL*) by computing congruence coefficients φ (Tucker, 1951) between the components of the true and estimated loading matrices and averaging across components and clusters as follows:

$$GOCL = \frac{\sum_{k=1}^K \sum_{q=1}^Q \varphi(\mathbf{B}_q^{(k)T}, \mathbf{B}_q^{(k)M})}{KQ} \quad (8)$$

with φ being the Tucker phi coefficient and $\mathbf{B}_q^{(k)T}$ and $\mathbf{B}_q^{(k)M}$ indicating the q -th component of the true and estimated cluster loading matrices, respectively. The rotational freedom of the Clusterwise SCA-ECP model was dealt with by an orthogonal procrustes rotation of the estimated loading matrices towards the true loading matrices. To take the permutational freedom of the clusters into account, the permutation was chosen that maximizes the *GOCL* value. The *GOCL* statistic takes values between 0 (no recovery at all) and 1 (perfect recovery).

In our simulation study, the *GOCL*-value varies between .86 and 1.00, with a mean value of .9969 ($SD = 0.01$), indicating an excellent recovery of the $\mathbf{B}^{(k)}$ matrices. An analysis of variance was performed with *GOCL* as the dependent variable, which revealed a main effect of the number of observations per data block ($\hat{\omega}^2 = .20$) – the lower the number of observations per data block, the worse the cluster loading matrices are recovered (Figure 5) – and a main effect of the error level ($\hat{\omega}^2 = .14$) – the cluster loading matrices are more difficult to recover when the data contain more error (Figure 5). The main effects of error and

the number of observations per data block are qualified by the interaction of these factors ($\hat{\omega}^2 = .08$): The effect of the number of observations per data block is stronger when the data contain more error (Figure 5).

[Insert Figure 5 about here]

Simulation Study 2

To investigate how the goodness of recovery is affected by performing Clusterwise SCA-ECP with an incorrect number of clusters or components, we simulated data according to the seven factor design from Simulation Study 1, with five replications per cell of the design. Each of the resulting 7,290 data matrices \mathbf{X} was analyzed four times, each time using 25 random starts: 1) with K clusters and $Q-1$ components per cluster, 2) with K clusters and $Q+1$ components per cluster, 3) with $K-1$ clusters and Q components per cluster, and 4) with $K+1$ clusters and Q components per cluster – where K denotes the correct number of clusters and Q the correct number of components that is underlying the data. When using incorrect values of Q (analysis 1 and 2), the overall *ARI* equals .79 ($SD = 0.34$) and .96 ($SD = 0.16$) when $Q-1$ and $Q+1$ components are extracted for each cluster, respectively. Thus, extracting too many components (‘overextraction’) appears to be less problematic for the recovery of the clustering than extracting too few components (‘underextraction’). This makes sense as for some data sets all Q components will be needed to properly distinguish between the clusters, whereas overextraction will most often yield all the information needed to disentangle the different clusters. Moreover, these results are in line with literature on the effects of overextraction and underextraction in component analysis (Fava & Velicer, 1992;

Wood, Tataryn, & Gorsuch, 1996), i.e., it is generally found that overextraction introduces less error to the estimated loading structure than underextraction.

When using incorrect values of K , the most interesting question is how this affects the recovery of the clustering. When estimating too few clusters, we hypothesize that either two underlying clusters will be fused or that the data blocks of one true cluster will be spread across two or more of the estimated clusters. Vice versa, when estimating too many clusters, we expect that either one true cluster will be split up or that an additional cluster will be composed in the analysis out of an amalgam of data blocks that belong to different true clusters. To evaluate which of these phenomena occurs, we computed the *ARI* between the estimated clustering and the true clustering in which two clusters were merged (too low K used) or the true clustering splitting one of the true clusters into two clusters (too high K used); note that all possible splits and fusions were considered, retaining the split or fusion that yielded the highest *ARI* value. If these *ARI* values equal 1, this implies that the analysis deals with misspecification of K by cleanly fusing two true clusters or splitting one true cluster. Indeed, on average, *ARI* amounts to .99 ($SD = 0.07$) for the analyses with one cluster too few (note that the data sets with two true clusters were not taken into account) and .98 ($SD = 0.08$) for the analyses with one cluster too many.

Simulation Study 3

An additional small simulation study was performed to evaluate the goodness of recovery of Clusterwise SCA-ECP analyses when the data contain variables that do not reflect the cluster structure of interest. These variables can mask the clustering and are thus named ‘masking variables’ (see e.g., Brusco & Cradit, 2001; Vichi & Kiers, 2001). The

design of the simulation study consists of factors 3, 4, 6, and 7 of Simulation Study 1, with the number of data blocks (factor 1) fixed at 40, the number of observations per data block (factor 2) between 80 and 120, and equal cluster sizes (factor 5). Each cell of the design was replicated 10 times. To each of the thus obtained 810 data matrices, consisting of 12 variables, 6 masking variables were added. To test the influence of two different kinds of masking variables, the six masking variables could either be summarized by the same two components across all datablocks (810 data sets in the no clustering condition) or by two components that differed across the data blocks and thus reflected another clustering of the data blocks than the clustering of interest (810 data sets in conflicting clustering condition). The loadings of these additional masking components are sampled uniformly between -1 and 1 (as in level 1 of factor 7 in Simulation Study 1). Each simulated data matrix is analyzed by Clusterwise SCA-ECP using the correct number of clusters and components.

On average, the *ARI* between the true clustering and the estimated clustering amounts to .97 ($SD = 0.13$) and .71 ($SD = 0.43$) in the no clustering and conflicting clustering conditions, respectively. In the conflicting clustering condition, the *ARI* largely depends on the level of congruence of the cluster loading matrices for the 12 original variables ($\hat{\omega}^2 = .50$): the mean *ARI* is .99, .86 and .29 for cluster loading matrices with low, medium and high congruence respectively.

Conclusion

Based on the results of Simulation Study 1 we can conclude that, for simulated data sets, the Clusterwise SCA-ECP algorithm is not prone to end in a local minimum when using a multistart procedure with 25 random starts, given that the correct number of clusters and

components are used. Furthermore, the algorithm appears to recover both the partition matrix and the relative sizes of the cluster loadings very well. The latter result puts the tendency of the algorithm for overfitting in perspective, as it suggests that the overfitting primarily boils down to a slight overestimation of the loadings, which is a well-known result in component analysis and which does not affect the interpretation of the components (see e.g., Velicer, Peacock, & Jackson, 1982). On the basis of Simulation Study 2, we can infer that the performance of Clusterwise SCA-ECP is affected more by using too few clusters or components than by using too many. Typically, when using one cluster too many, one of the underlying clusters is split up in the estimated model. Finally, Simulation Study 3 has shown that masking variables have little influence on the goodness of recovery, except when the masking variables reflect a conflicting clustering and the clustering of interest is less distinct due to high congruence of the corresponding cluster loading matrices.

The results of the simulation studies can be translated into some guidelines about the use of Clusterwise SCA-ECP analysis in practice. First, Simulation Study 1 indicated that 25 random starts are sufficient under ideal conditions. Therefore, when analyzing empirical data, one may adopt the following strategy: When exploring different numbers of clusters and components, 25 starts can be used, but afterwards, when one retains some solutions (i.e., some particular combinations of K and Q) for further investigation, it may be wise to repeat the estimation of these solutions with a larger number of random starts (say, 50 or 100) to reduce the probability of obtaining a local minimum. Also, Simulation Study 1 indicates that one should be cautious when the number of observations per data block is low, the number of clusters and components may be large, the data may contain a lot of error, and when the cluster loadings may be highly congruent. Among these data characteristics, the number of observations per data block is the only one that can be observed; therefore, we emphasize that

it is important to have enough observations within each data block (preferably more observations than variables). Second, the results of Simulation Study 2 imply that, when in doubt about the underlying number of clusters or components, it is advisable to consider high enough numbers. Third, based on Simulation Study 3, we hypothesize that when the data contain masking variables, Clusterwise SCA-ECP will still be able to recover the clustering, unless the structures of the clusters are highly congruent or the number of masking variables is too large.

Illustrative Application 1

In this section, we present an empirical example from eating disorder research. In this research domain, it has been observed that a substantial proportion of patients with eating disorders, like Anorexia Nervosa (AN) or Bulimia Nervosa (BN), engage in high levels of physical activity (Beumont, Arthur, Russel, & Touyz, 1994; Davis, 1997; Solenberger, 2001). Partly due to the absence of a clear operational definition of excessive exercise, few studies have provided estimates of the prevalence of this high-level exercising among eating disordered patients. However, there is some consensus that approximately 80% of the AN and 55% of the BN patients engage in excessive exercising (Davis, Kaptein, Kaplan, Olmsted, & Woodside, 1998; Epling, Pierce, & Stefan, 1983; Peñas-Lledó, Vaz Leal, & Waller, 2002). To explain this excessive physical activity at the within-patient level, two hypotheses have been put forward about the underlying psychological processes. The drive for thinness hypothesis states that eating disorder patients engage in physical activity because they are actively trying to lose weight by burning calories (Davis, 1997; Davis, Kennedy, Ravelski, & Dionne, 1994; Heatherton & Baumeister, 1991). The affect regulation hypothesis reads that physical activity is a way of coping with chronically negative affect (Davis,

Katzman, & Kirsh, 1999; Holtkamp, Hebebrand, & Herpertz-Dahlmann, 2004; Thome & Espelage, 2004). In the subgroup of patients who do not display excessive levels of physical activity, physical activity at the within-patient level is not expected to be related to burning calories but to increased positive affect just like in normal control subjects (Gauvin, Rejeski, & Norris, 1996; Kelsey et al., 2006).

To trace whether there are interindividual differences in how physical activity is related to drive for thinness and affect, Vansteelandt, Rijmen, Pieters, Probst and Vanderlinden (2007) studied 32 female patients of the specialized inpatient eating disorder unit of the University Psychiatric Centre in Leuven, Belgium. Of these patients, 19 suffered from AN and 13 from BN. During one week, at nine randomly selected times a day the patients were signalled by an electronic device to fill out a questionnaire. The questionnaire consisted of 22 items that aimed at measuring the momentary drive for thinness, positive and negative emotional states, urge to be physically active and physical activity. Inevitably, patients typically missed some of the 63 signals (7 days \times 9 signals) resulting in a different number of assessments for different patients. On average 45.60 observations were obtained per patient ($SD = 10.78$). This study results in a 1,459 observations by 22 variables data matrix \mathbf{X} , consisting of 32 submatrices \mathbf{X}_i , one for each patient, where each \mathbf{X}_i was centered and autoscaled. Vansteelandt et al. (2007) also asked all patients to fill in the Eating Disorder Inventory (EDI), the Commitment to Exercise Scale (CES), the Symptom Checklist (SCL-90) and the Beck Depression Inventory (BDI) as between-patient measures.

To trace the most important interindividual differences in within-patient relations between drive for thinness, affect and physical activity, we analyzed the data set with Clusterwise SCA-ECP, using 25 random starts and with K and Q varying from 1 to 8. The Clusterwise SCA-ECP solutions with K equal to one are SCA-ECP solutions (see Model

section). Separate PCA's were also performed. In Figure 6 the percentage of explained variance of the different Clusterwise SCA-ECP solutions is plotted against the number of components for each value of K .

[Insert Figure 6 about here]

To select a model out of the obtained solutions, we used the model selection procedure that was described in the Data Analysis section. First, the number of clusters K is determined. Since going from one (i.e., SCA-ECP) to two clusters gives an increase in fit of about 2.5%, while adding more clusters hardly improves the percentage of explained variance, we decided to choose between the solutions with two clusters. Second, out of these solutions, a model with two components is chosen, since the scree line of the Clusterwise SCA-ECP solutions with two clusters shows an elbow at that point. On top of that, the solution with two components and two clusters appeared to be the best in terms of interpretability of the cluster loading matrices and stability of the subject partition. For the eating disorder data, the best fitting partition allocates 25 subjects in cluster 1 and 7 subjects in cluster 2. This partition occurs once within the 25 runs. Out of the remaining runs, however, 22 yield a solution with a very similar partition in which only one or two subjects are assigned to another cluster. Reallocating these subjects does not alter the structure of the cluster loading matrices. Since the partition with 7 subjects in cluster 2, which is the global optimum to the best of our knowledge, is retrieved only once in the 25 randomly started runs of the algorithm, it may as well be missed in another 25 runs. So, in this case, performing a higher number of runs (e.g., 100), would be useful to decrease the probability of landing in a local minimum.

[Insert Table 5 about here]

The Normalized Varimax rotated component loadings² for the two clusters of this solution are given in Table 5. In this table, it can be seen that in cluster 1 all emotions load high on the first component, which is therefore labeled positive (PA) versus negative affect (NA). Feeling fat and ugly also load high on this component and are related to negative affect. The second component of this cluster shows high loadings for the items measuring urge to be physically active (UPhA), physical activity (PhA) and drive for thinness (DT). This indicates that urge to be physically active, physical activity and drive for thinness (especially drive for burning calories and losing weight) covary, which is in line with the drive for thinness hypothesis. In cluster 2 a different structure shows up. The first component has highly negative loadings for the negative emotions and highly positive loadings for the emotions ‘pleased’ and ‘happy’. Therefore this first component is labeled negative affect versus being pleased. On the second component all positive emotional states are loading high, together with two of the physical activity items. The second component thus indicates that positive affect and physical activity are related for the subjects in cluster 2.

Summarizing, the patients in the first cluster tend to engage in physical activity when they feel the momentary urge to be physically active and burn calories (drive for thinness). For the subjects in cluster 2, however, physical activity is not driven by a drive for thinness or by affect regulation. Instead, physical activity is related to positive affect, so it seems that they tend to move more when they are feeling happier (or vice versa). As a result, it appears

² Some eating disorder patients (especially in cluster 2) displayed no variance on one or more of the variables. As the invariant scores cannot be rescaled to have a variance of one, they were replaced by zeros. Therefore, in this case, even though the components are orthogonal, the loadings should be read as weights rather than correlations.

that the Clusterwise SCA-ECP made a distinction between the cluster of patients engaging in excessive activity to burn calories and the cluster of patients who do not engage in excessive physical activity and for whom activity is associated with positive affect like in normal controls. When we use the prevalence estimates of excessive exercising in eating disorders mentioned in the introduction of this section, we expect about 22 of the subjects ($.80 \times 25$ AN patients + $.55 \times 7$ BN patients) to belong to cluster 1 and 10 to cluster 2 which is in line with the results (25 belong to cluster 1 and 7 to cluster 2).

These between-cluster differences in the mechanisms behind physical activity are further related to the severity of the eating disorder. Specifically, comparing the scores of the patients in clusters 1 and 2 on the Beck Depression Inventory ($t = 4.36, p < .001$), the Commitment to Exercise Scale ($t = 3.48, p = .002$) and the Eating Disorder Inventory ($t = 2.66, p = .02$) reveals that the patients in cluster 1 are more depressed, are more committed to physical exercise and have higher body dissatisfaction. These results are in line with the original study of Vansteelandt et al. (2007) who also found that the within-patient relation between drive for thinness and urge to be physically active/physical activity – the relation found in cluster 1 – was stronger for patients scoring higher on depression in particular and pathology in general.

This application illustrates the value of Clusterwise SCA-ECP as a generic modeling strategy, since it shows that, using this strategy, a solution with a good balance between fit and model complexity can be obtained. Specifically, when two components are extracted, performing separate PCA's explains roughly 10% more variance than the standard (i.e., non-clusterwise) SCA-ECP approach. Yet, by classifying the 32 subjects in two clusters only, Clusterwise SCA-ECP was already able to explain about 2.5% more variation than standard SCA-ECP. Since fitting more than two clusters leads to considerably smaller increases in fit,

the additional clusters seem to chart more idiosyncratic structural differences. We can conclude that the selected Clusterwise SCA-ECP solution is only slightly more complex than the corresponding SCA-ECP solution, but yields much more information on the underlying mechanisms of physical activity for the eating disorder patients. More specifically, in line with earlier findings in eating disorder research, for the eating disorder patients that display high levels of (urge for) physical activity (cluster 1), being active is associated with drive for thinness (Davis, 1997), whereas for the patients that do not engage in excessive physical activity (cluster 2), activity is related to positive affect (Gauvin et al., 1996; Kelsey et al., 2006).

Illustrative Application 2

The second empirical example stems from social psychology and addresses how people react emotionally to unequal (and hence, unfair) outcomes dependent on the personal need of another person. That is, people in need may arouse empathy and this is likely to have an effect on the emotional reaction. This was examined using a public good dilemma game in which the participant and an opposing player contribute to a public good. Note that a good becomes a public good, when individuals can benefit from the good even if they did not contribute to its provision. Yet, when the combined investment of both players is sufficient, the public good can be provided and both obtain a reward. An important principle in public good dilemmas is equality (Van Dijk & Wilke, 1995, 2000): When both players contribute equally, they want to be rewarded equally. Violation of equality is perceived as unjust and elicits anger (Stouten, De Cremer, & Van Dijk, 2005). In case of inequality, the subject's negative emotional reactions may be mollified if the advantaged person is in personal need,

because feelings of sympathy and empathy can be aroused (Stouten, Ceulemans, Timmerman, & Van Hiel, 2011).

To study the effects of equality violation and the opponent being in need on a broad range of emotions, Stouten et al. (2011) performed an experimental study in which 282 participants played a public good dilemma with an alleged opponent, in which it was said that sufficient contributions were made to provide the public good. Hence, both players earned a reward. In the experiment, two factors were fully crossed, yielding 6 conditions. Firstly, the reward was equal for both players (equal conditions) or was higher for the opponent (unequal conditions). Secondly, the amount of empathy was manipulated at three levels. In the high empathy conditions, the participants received a message about a negative personal event that happened to the opponent and they are asked to imagine how the opponent felt. In the low empathy conditions, participants also received the above message but are asked to take an objective perspective. In the control conditions, participants did not receive a message. After the game, participants were asked how they felt using 21 emotion terms, to be rated on a 7-point Likert scale. This study resulted in a 282 participants by 21 emotions data matrix, where each of the participants is nested within a particular equality/empathy condition.

To investigate how the structure of the experienced emotions within a condition possibly differs across the 6 conditions, we applied Clusterwise SCA-ECP to the data. We used 25 random starts, and varied K from 1 to 6 and Q from 1 to 6. Clusterwise SCA-ECP solutions with K equal to one are SCA-ECP solutions, and Clusterwise SCA-ECP with K equal to six boils down to the separate PCA approach since there are six conditions (see Model section). In Figure 7, the percentages of explained variance of the different Clusterwise SCA-ECP solutions is plotted against the numbers of components for each value of K . Upon inspection of Figure 7, we decided to select a solution in which the conditions are

grouped into two clusters, because going from one (SCA-ECP) to two clusters gives a large increase in fit, while adding more clusters gives substantially smaller improvements in fit. Subsequently, we selected the model with two components, since the scree line of the Clusterwise SCA-ECP solutions with two clusters showed a clear elbow at two components.

[Insert Figure 7 about here]

In the selected solution, cluster 1 contains all equal conditions (equal/high empathy, equal/low empathy, equal/control) and cluster 2 all unequal conditions (unequal/high empathy, unequal/low empathy, unequal/control). Thus, the equality factor seems to have the largest effect on the underlying structure of the emotion patterns in the different conditions.

[Insert Table 6 about here]

The Normalized Varimax rotated component loadings for the two clusters of the selected solution, are given in Table 6. In this table, it can be seen that the first component of the cluster with the equal conditions is a positive affect versus negative affect (PA vs. NA) component. The second component is characterized by high loadings for empathic emotions such as ‘compassionate’ and ‘tenderhearted’ and thus can be labeled empathy. The components of the unequal conditions are rather similar to the ones of the equal conditions, and can be labeled PA vs. NA and empathy as well. However, a very different loading pattern across the two clusters is found for three emotions – fearful, surprised and sympathetic.

‘Fearful’ loads highly on the PA vs. NA component in the equal conditions cluster, while in the cluster with unequal conditions it is related to the empathy emotions of the second component. A tentative interpretation of this reads that within the unequal conditions, it may be that people who react empathically (for instance, based on the message received) to the unequal treatment also feel fearful, because they do not know what is going to happen

next. Additionally, participants may feel powerless, which is also argued to be one of the core components of fear (Ellsworth & Smith, 1988). Hence, due to the uncertainty of the situation and the fact that participants have to accept the given situation they also feel fearful.

‘Surprised’ loads high on the empathy component in the equal conditions, and on the PA vs. NA component in the unequal conditions. Thus, in the equal conditions the feelings of surprise go together with an empathic reaction, whereas in the unequal conditions surprise is associated with anger about the unequal treatment of the players of the game. This negative surprise is elicited as part of the unexpected result of receiving unequal outcomes (e.g., feeling stunned, seeking understanding).

Concerning ‘sympathetic’, it can be seen that in the unequal conditions this emotion loads higher on the empathy component than in the equal conditions. Although participants received an unequal (and unfair) outcome, the fact that the other person was in need elicits understanding, and therefore the feeling of sympathy. This implies that, in the unequal conditions, sympathy is associated with an empathetic reaction towards the opponent, while this is less outspoken in the equal conditions.

This empirical example with subjects nested within groups shows that Clusterwise SCA-ECP can also be a useful tool to explore structural differences between groups of subjects. Moreover, this application demonstrates that Clusterwise SCA-ECP is able to detect interpretable differences in component structure, even when these differences are small.

Discussion

In scientific research, many questions concern differences and similarities in the underlying structure of multivariate multiblock data. To answer such research questions, two commonly used strategies are separate principal component analyses and simultaneous component analysis, with the first one often being too flexible to provide insight and the second one too parsimonious to detect differences. Therefore, the Clusterwise SCA-ECP model was introduced as a generic modeling strategy for detecting structural differences and similarities in multiblock multivariate data, where each data block can consist of multiple measurements for a certain subject or of multiple subjects that belong to a particular group. Specifically, the most important structural differences between the data blocks are captured by clustering the data blocks: Data blocks with similar component loadings are grouped in the same cluster, whereas data blocks from different clusters have different component loadings. Thus, the structural differences between the data blocks can be derived from the comparison of the cluster loading matrices. Clusterwise SCA-ECP comprises the separate PCA (number of clusters equal to number of data blocks) and SCA-ECP (number of clusters equal to one) approaches as special cases, as well as all the intermediate solutions (number of clusters between one and number of data blocks). Therefore, it allows to look for a solution with a good balance between fit and model complexity.

We see various directions for future research, at the level of the data, the level of the model, and the level of the data analysis.

Level of the data. The data blocks were autoscaled, implying that the information on the between-block differences in variability is lost. For certain research questions, these differences can be very interesting; for example, to study whether or not emotions vary more over time for some people than for others and whether these differences in variability are related to personality traits (Oravecz, Tuerlinckx, & Vandekerckhove, 2009). Therefore, it

may be worthwhile to consider standardizing the variables over the data blocks. This implies severe adjustments to the modeling strategy, however.

Level of the model. Using SCA-ECP within each cluster, all between-block differences in correlations between component scores must go into the clustering. Thus, data blocks with the same component loadings but different correlations between the components, would be assigned to separate clusters. For example, when personality traits are measured for subjects from a number of countries, it is possible that although in some countries the same personality traits are found, they are differentially correlated. In that case, Clusterwise SCA-ECP would assign these countries to separate clusters, while putting them in the same cluster would be more parsimonious. A solution would be to use a variant of SCA within the clusters that allows for differences in correlations between the components.

Level of the data analysis. In this paper, the number of components was fixed across the clusters. This is often not very realistic. For example, in studies on emotional granularity, subjects with a high granularity are expected to have more components in their affect space than subjects with a low emotional granularity. Thus, in future research it would be useful to allow the number of components to vary between clusters.

To conclude, we presented Clusterwise SCA-ECP, a versatile and generic modeling strategy for detecting and describing structural differences and similarities in multivariate multiblock data. User-friendly Clusterwise SCA-ECP software is available at <http://ppw.kuleuven.be/okp/software/MBCA/>. Based on simulation studies, guidelines were formulated for the use of Clusterwise SCA-ECP analysis in practice.

References

- Aloise, D., Hansen, P., Liberti, L. (in press). An improved column generation algorithm for minimum sum-of-squares clustering. *Mathematical Programming*.
- Barrett, L. F. (1998). Discrete emotions or dimensions? The role of valence focus and arousal focus. *Cognition and Emotion*, *12*, 579–599.
- Beumont, P. J. V., Arthur, B., Russell, J. D., & Touyz, S. W. (1994). Excessive physical activity in dieting disorder patients: Proposals for a supervised exercise program. *International Journal of Eating Disorders*, *15*, 21–36.
- Bock, H. H. (1987). On the interface between cluster analysis, principal component analysis, and multidimensional scaling. In H. Bozdogan & A. K. Gupta (Eds.), *Multivariate statistical modeling and data analysis* (pp. 17–34). Dordrecht, The Netherlands: Reidel Publishing.
- Bouveyron, C., Girard, S., & Schmid, C. (2007). High-Dimensional Data Clustering. *Computational Statistics and Data Analysis*, *52*, 502–519.
- Brusco, M. J. (2006). A repetitive branch-and-bound algorithm for minimum within-cluster sums of squares partitioning. *Psychometrika*, *71*, 347–363.
- Brusco, M. J., & Cradit, J. D. (2001). A variable selection heuristic for K -means clustering. *Psychometrika*, *66*, 249–270.
- Brusco, M. J., & Cradit, J. D. (2005). ConPar: A method for identifying groups of concordant subject proximity matrices for subsequent multidimensional scaling analyses. *Journal of Mathematical Psychology*, *49*, 142–154.

- Brusco, M. J., Cradit, J. D., Steinley, D., & Fox, G. (2008). Cautionary remarks on the use of clusterwise regression. *Multivariate Behavioral Research, 43*, 29–49.
- Cattell, R. B. (1966). The scree test for the number of factors. *Multivariate Behavioral Research, 1*, 245–276.
- Davis, C. (1997). Eating disorders and hyperactivity: A psychobiological perspective. *Canadian Journal of Psychiatry, 42*, 168–175.
- Davis, C., Kaptein, S., Kaplan, A. S., Olmsted, M. P., & Woodside, D. B. (1998). Obsessionality in anorexia nervosa: The moderating influence of exercise. *Psychosomatic Medicine, 60*, 192–197.
- Davis, C., Katzman, D. K., & Kirsh, C. (1999). Compulsive physical activity in adolescents with anorexia nervosa: A psychobehavioral spiral of pathology. *Journal of Nervous and Mental Disease, 187*, 336–342.
- Davis, C., Kennedy, S. H., Ravelski, E., & Dionne, M. (1994). The role of physical activity in the development and maintenance of eating disorders. *Psychological Medicine, 24*, 957–967.
- DeSarbo, W. S., Oliver, R. L., & Rangaswamy, A. (1989). A simulated annealing methodology for clusterwise linear regression. *Psychometrika, 54*, 707–736.
- de Soete, G., & Carrol, J. D. (1994). K-means clustering in a low-dimensional Euclidean space. In E. Diday, Y. Léchevallier, M. Schader, P. Bertrand, & B. Burtschy (Eds.), *New approaches in classification and data analysis* (pp. 212–219). Berlin, Germany: Springer.

- De Roover, K., Ceulemans, E., & Timmerman, M. E. (in press). How to perform multiblock component analysis in practice. *Behavior Research Methods*.
- Dolan, C. V., Oort, F. J., Stoel, R. D., & Wicherts, J. M. (2009). Testing measurement invariance in the target rotated multigroup exploratory factor model. *Structural Equation Modeling, 16*, 295–314.
- Dolan, C. V., & van der Maas, H. L. J. (1998). Fitting multivariate normal finite mixtures subject to structural equation modeling. *Psychometrika, 63*, 227–253.
- Eid, M., & Diener, E. (2001). Norms for experiencing emotions in different cultures: Inter- and intranational differences. *Journal of Personality and Social Psychology, 81*, 869–885.
- Ekman, P. (1999). Basic emotions. In T. Dalgleish & M. J. Power (Eds.), *Handbook of cognition and emotion* (pp. 45–60). Chichester, UK: Wiley.
- Ellsworth, P. C., & Smith, C. A. (1988). From appraisal to emotion: Differences among unpleasant feelings. *Motivation and Emotion, 12*, 271–302.
- Epling, W. F., Pierce, W. D., & Stefan, L. (1983). A theory of activity-based anorexia. *International Journal of Eating Disorders, 7*, 27–46.
- Fava, J. L., & Velicer, W. F. (1992). The effects of overextraction on factor and component analysis. *Multivariate Behavioral Research, 27*, 387–415.
- Fontaine, J. R. J., Poortinga, Y. H., Setiadi, B., & Markam, S. S. (2002). Cognitive structure of emotion terms in Indonesia and the Netherlands. *Cognition & Emotion, 16*, 61–86.

- Fontaine, J. R. J., Scherer, K. R., Roesch, E. B., & Ellsworth, P. C. (2007). The world of emotions is not two-dimensional. *Psychological Science, 18*, 1050–1057.
- Gauvin, L., Rejeski, W. J., & Norris, J. L. (1996). A naturalistic study of the impact of acute physical activity on feeling states and affect in women. *Health Psychology, 15*, 391–397.
- Haavelmo, T. (1943). The statistical implications of a system of simultaneous equations. *Econometrica, 11*, 1–2.
- Hands, S., & Everitt, B. (1987). A monte carlo study of the recovery of cluster structure in binary data by hierarchical clustering techniques. *Multivariate Behavioral Research, 22*, 235–243.
- Hays, W. L. (1963). *Statistics for psychologists*. New York: Holt, Rinehart & Winston.
- Heatherton, T. F., & Baumeister, R. F. (1991). Binge eating as escape from self-awareness. *Psychological Bulletin, 110*, 86–108.
- Hessen, D. J., Dolan, C. V, & Wicherts, J. M. (2006). Multi-group exploratory factor analysis and the power to detect uniform bias. *Applied Psychological Research, 30*, 233–246.
- Holtkamp, K., Hebebrand, J., & Herpertz-Dahlmann, B. (2004). The contribution of anxiety and food restriction on physical activity levels in acute anorexia nervosa. *International Journal of Eating Disorders, 36*, 163–171.
- Hubert, L., & Arabie, P. (1985). Comparing partitions. *Journal of Classification, 2*, 193–218.
- Jedidi, K., Jagpal, H. S., & DeSarbo, W. S. (1997). STEMM: A general finite mixture structural equation model. *Journal of Classification, 14*, 23–50.

Jolliffe, I. T. (1986). *Principal component analysis*. New York: Springer.

Jöreskog, K. G. (1971). Simultaneous factor analysis in several populations. *Psychometrika*, 36, 409–426.

Kaiser, H. F. (1958). The Varimax criterion for analytic rotation in factor analysis. *Psychometrika*, 23, 187–200.

Kelsey, K. S., DeVellis, B. M., Begum, M., Belton, L., Hooten, E. G., & Campbell, M. K. (2006). Positive affect, exercise and self-reported health in blue-collar women. *American Journal of Health Behavior*, 30, 199–207.

Kiers, H. A. L. (1990). *SCA. A program for simultaneous components analysis of variables measured in two or more populations*. Groningen, The Netherlands: iec ProGAMMA.

Kiers, H. A. L., & ten Berge, J. M. F. (1994). Hierarchical relations between methods for Simultaneous Components Analysis and a technique for rotation to a simple simultaneous structure. *British Journal of Mathematical and Statistical Psychology*, 47, 109–126.

Kline, R. B. (2004). *Principles and practice of structural equation modeling* (2nd ed.). New York: Guilford Press.

Kroonenberg, P. M. (2008). *Applied multiway data analysis*. Hoboken, NJ: Wiley.

Kruskal, J. B., & Wish, M. (1978). *Multidimensional scaling*. Newberry Park, California: Sage.

Lawley, D. N., & Maxwell, A. E. (1962). Factor analysis as a statistical method. *The Statistician*, 12, 209–229.

- MacKinnon, N. J., & Keating, L. J. (1989). The structure of emotions: Canada-United States comparisons. *Social Psychology Quarterly*, *52*, 70–83.
- MacQueen, J. B. (1967). Some methods for classification and analysis of multivariate observations. *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability* (pp. 281–297). Berkeley, CA: University of California Press.
- McLachlan, G. J., & Peel, D. (2000). *Finite mixture models*. New York: Wiley.
- Meredith, W., & Millsap, R. E. (1985). On component analyses. *Psychometrika*, *50*, 495–507.
- Milligan, G. W., Soon, S. C., & Sokol, L. M. (1983). The effect of cluster size, dimensionality, and the number of clusters on recovery of true cluster structure. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, *5*, 40–47.
- Millsap, R. E., & Meredith, W. (1988). Component analysis in cross-sectional and longitudinal data. *Psychometrika*, *53*, 123–134.
- Olejnik, S., & Algina, J. (2000). Measures of effect size for comparative studies: Applications, interpretations, and limitations. *Contemporary Educational Psychology*, *25*, 241–286.
- Oravecz, Z., Tuerlinckx, F., & Vandekerckhove, J. (2009). A hierarchical Ornstein-Uhlenbeck model for continuous repeated measurement data. *Psychometrika*, *74*, 395–418.

- Pearson, K. (1901). On lines and planes of closest fit to systems of points in space. *Philosophical Magazine*, 2, 559–572.
- Peñas-Lledó, E., Vaz Leal, F. J., & Waller, G. (2002). Excessive exercise in anorexia nervosa and bulimia nervosa: Relation to eating characteristics and general psychopathology. *International Journal of Eating Disorders*, 31, 370–375.
- Rocci, R., & Vichi, M. (2005). Three-mode component analysis with crisp or fuzzy partition of units. *Psychometrika*, 70, 715–736.
- Rodriguez, C., & Church, A. T. (2003). The structure and personality correlates of affect in Mexico: Evidence of cross-cultural comparability using the Spanish language. *Journal of Cross Cultural Psychology*, 34, 211–230.
- Russell, J. A., & Barrett, L. F. (1999). Core affect, prototypical emotional episodes, and other things called emotion: Dissecting the elephant. *Journal of Personality and Social Psychology*, 76, 805–819.
- Schouteden, M., Van Deun, K., Van Mechelen, I., & Pattyn, S. (2010). SCA and rotation to distinguish common and specific information in coupled data. Manuscript in preparation.
- Solenberger, S. E. (2001). Exercise and eating disorders: A 3-year inpatient hospital record analysis. *Eating Behaviors*, 2, 151–168.
- Sörbom, D. (1974). A general method for studying differences in factor means and factor structure between groups. *British Journal of Mathematical and Statistical Psychology*, 27, 229–239.

- Späth, H. (1979). Algorithm 39: Clusterwise linear regression. *Computing*, 22, 367–373.
- Späth, H. (1982). Algorithm 48: A fast algorithm for clusterwise linear regression. *Computing*, 29, 175–181.
- Steinley, D. (2003). Local optima in K-means clustering: What you don't know may hurt you. *Psychological Methods*, 8, 294–304.
- Stouten, J., Ceulemans, E., Timmerman, M. E., & Van Hiel, A. (2011). Tolerance of justice violations: The effects of need on emotional reactions after violating equality in social dilemmas. *Journal of Applied Social Psychology*, 41, 357–380.
- Stouten, J., De Cremer, D., & Van Dijk, E. (2005). All is well that ends well, at least for proselves: Emotional reactions to equality violation as a function of social value orientation. *European Journal of Social Psychology*, 35, 767–783.
- ten Berge, J. M. F. (1993). *Least squares optimization in multivariate analysis*. Leiden: DSWO press.
- Thome, J., & Espelage, D. L. (2004). Relations among exercise, coping, disordered eating, and psychological health among college students. *Eating Behavior*, 5, 337–351.
- Timmerman, M. E., Ceulemans, E., Kiers, H. A. L., & Vichi, M. (2010). Factorial and reduced K-means reconsidered. *Computational Statistics & Data Analysis*, 54, 1858–1871.
- Timmerman, M. E., & Kiers, H. A. L. (2003). Four simultaneous component models of multivariate time series from more than one subject to model intraindividual and interindividual differences. *Psychometrika*, 68, 105–122.

- Timmerman, M. E., Kiers, H. A. L., Smilde, A. K., Ceulemans, E., & Stouten, J. (2009). Bootstrap confidence intervals in multi-level simultaneous component analysis. *British Journal of Mathematical & Statistical Psychology*, *62*, 299–318.
- Tucker, L. R. (1951). *A method for synthesis of factor analysis studies* (Personnel Research section Rep. No. 984). Washington, DC: Department of the Army.
- Tucker, L. R., & Messick, S. (1963). An individual differences model for multidimensional scaling. *Psychometrika*, *28*, 333–367.
- Tugade, M. M., Fredrickson, B. L., & Barrett, L. F. (2004). Psychological resilience and positive emotional granularity: Examining the benefits of positive emotions on coping and health. *Journal of Personality*, *72*, 1161–1190.
- van de Vijver, F. J. R., & Leung, K. (1997). *Methods and data analysis for cross-cultural research*. Thousand Oaks, California: Sage.
- Van Deun, K., Smilde, A. K., van der Werf, M. J., Kiers, H. A. L., & Van Mechelen, I. (2009). A structured overview of simultaneous component based data integration. *BMC Bioinformatics*, *10*, 246.
- Van Dijk, E., & Wilke, H. (1995). Coordination rules in asymmetric social dilemmas: A comparison between public good dilemmas and resource dilemmas. *Journal of Experimental Social Psychology*, *31*, 1–27.
- Van Dijk, E., & Wilke, H. (2000). Decision-induced focusing in social dilemmas: Give-some, keep-some, take-some, and leave-some dilemmas. *Journal of Personality and Social Psychology*, *78*, 92–104.

- Vansteelandt, K., Rijmen, F., Pieters, G., Probst, M., & Vanderlinden, J. (2007). Drive for thinness, affect regulation and physical activity in eating disorders: A daily life study. *Behaviour Research and Therapy*, *45*, 1717–1734.
- Varriale, R., & Vermunt, J. K. (2010). Multilevel mixture factor models. Manuscript submitted for publication.
- Velicer, W. F., Peacock, A. C., & Jackson, D. N. (1982). A comparison of component and factor patterns: A monte carlo approach. *Multivariate Behavioral Research*, *17*, 371–388.
- Vermunt, J. K. (2003). Multilevel latent class models. *Sociological Methodology*, *33*, 213–239.
- Vermunt, J. K. (2008a). Latent class and finite mixture models for multilevel data sets. *Statistical Methods in Medical Research*, *17*, 33–51.
- Vermunt, J. K. (2008b). Multilevel latent variable modeling: An application in education testing. *Austrian Journal of Statistics*, *37*, 285–299.
- Vichi, M., & Kiers, H. A. L. (2001). Factorial k-means analysis for two-way data. *Computational Statistics & Data Analysis*, *37*, 49–64.
- Vichi, M., Rocci, R., & Kiers, H. A. L. (2007). Simultaneous component and clustering models for three-way data: Within and between approaches. *Journal of Classification*, *24*, 71–98.
- Wood, J. M., Tataryn, D. J., & Gorsuch, R. L. (1996). Effects of under- and overextraction on principal axis factor analysis with varimax rotation. *Psychological Methods*, *1*, 354–365.

Yung, Y. F. (1997). Finite mixtures in confirmatory factor-analysis models. *Psychometrika*, 62, 297–330.

Table 1

Hypothetical data matrix \mathbf{X} with the (rounded off) scores of 4 subjects on 6 variables measuring emotions and physical activity.

Person	Observation	Variables					
		Happy	Pleased	Sad	Ashamed	Moving	Sporting
1	1	-1.4	-1.4	-1.1	-1.1	-1.4	-1.4
	2	0.5	0.5	2.4	2.4	0.5	0.5
	3	-0.2	-0.2	-0.6	-0.6	-0.2	-0.2
	4	-0.2	-0.2	0.7	0.7	-0.2	-0.2
	5	1.4	1.4	-0.2	-0.2	1.4	1.4
	6	0.3	0.3	0.9	0.9	0.3	0.3
	7	0.2	0.2	-0.8	-0.8	0.2	0.2
	8	1.6	1.6	-1.4	-1.4	1.6	1.6
2	1	2.6	2.6	-0.8	-0.8	-0.8	-0.8
	2	-0.7	-0.7	0.7	0.7	0.7	0.7
	3	0.2	0.2	0.8	0.8	0.8	0.8
	4	-0.1	-0.1	-0.2	-0.2	-0.2	-0.2
	5	-1.9	-1.9	0.2	0.2	0.2	0.2
	6	-0.4	-0.4	-1.2	-1.2	-1.2	-1.2
	7	-1.8	-1.8	-1.1	-1.1	-1.1	-1.1
	8	0.8	0.8	0.1	0.1	0.1	0.1
	9	-0.9	-0.9	0.7	0.7	0.7	0.7
3	1	-0.2	-0.2	0.1	0.1	0.1	0.1
	2	-2.1	-2.1	-0.5	-0.5	-0.5	-0.5
	3	-0.8	-0.8	0.3	0.3	0.3	0.3
	4	1.4	1.4	-0.6	-0.6	-0.6	-0.6
	5	-1.1	-1.1	0.5	0.5	0.5	0.5
	6	1.0	1.0	0.7	0.7	0.7	0.7
	7	0.1	0.1	1.7	1.7	1.7	1.7
4	1	1.1	1.1	1.4	1.4	1.1	1.1
	2	-0.3	-0.3	-2.0	-2.0	-0.3	-0.3
	3	0.7	0.7	-0.2	-0.2	0.7	0.7
	4	-2.1	-2.1	-1.2	-1.2	-2.1	-2.1
	5	-0.4	-0.4	2.9	2.9	-0.4	-0.4
	6	-0.8	-0.8	0.8	0.8	-0.8	-0.8
	7	-1.6	-1.6	1.4	1.4	-1.6	-1.6
	8	0.5	0.5	-1.1	-1.1	0.5	0.5
	9	0.3	0.3	-0.5	-0.5	0.3	0.3
	10	0.0	0.0	-0.3	-0.3	0.0	0.0

Table 2

Partition matrix \mathbf{P} of the Clusterwise SCA-ECP decomposition with two clusters and two components of \mathbf{X} in Table 1. The ones indicate to which cluster each person is assigned.

Person	Clusters	
	Cluster 1	Cluster 2
Person 1	1	0
Person 2	0	1
Person 3	0	1
Person 4	1	0

Table 3

Component loading matrices $\mathbf{B}^{(1)}$ and $\mathbf{B}^{(2)}$ of the Clusterwise SCA-ECP decomposition with two clusters and two components of \mathbf{X} in Table 1.

	Cluster 1		Cluster 2	
	Positive affect & physical activity	Negative affect	Positive affect	Negative affect & physical activity
Happy	1	0	1	0
Pleased	1	0	1	0
Sad	0	1	0	1
Ashamed	0	1	0	1
Moving	1	0	0	1
Sporting	1	0	0	1

Table 4

Component score matrices \mathbf{F}_i of the Clusterwise SCA-ECP decomposition with two clusters and two components of \mathbf{X} in Table 1.

Person	Observation	Components			
		Cluster 1		Cluster 2	
		Positive affect & physical activity	Negative affect	Positive affect	Negative affect &
1	1	-1.4	-1.1		
	2	0.5	2.4		
	3	-0.2	-0.6		
	4	-0.2	0.7		
	5	1.4	-0.2		
	6	0.3	0.9		
	7	0.2	-0.8		
	8	1.6	-1.4		
2	1			2.6	-0.8
	2			-0.7	0.7
	3			0.2	0.8
	4			-0.1	-0.2
	5			-1.9	0.2
	6			-0.4	-1.2
	7			-1.8	-1.1
	8			0.8	0.1
	9			-0.9	0.7
3	1			-0.2	0.1
	2			-2.1	-0.5
	3			-0.8	0.3
	4			1.4	-0.6
	5			-1.1	0.5
	6			1.0	0.7
	7			0.1	1.7
4	1	1.1	1.4		
	2	-0.3	-2.0		
	3	0.7	-0.2		
	4	-2.1	-1.2		
	5	-0.4	2.9		
	6	-0.8	0.8		
	7	-1.6	1.4		
	8	0.5	-1.1		
	9	0.3	-0.5		
	10	0.0	-0.3		

Table 5

Normalized Varimax rotated loadings for the two clusters of the eating disorder data. Loadings greater than +/- .30 are highlighted in bold face. PA is positive affect, NA is negative affect. UPhA is urge to be physically active, PhA is physical activity and DT is drive for thinness.

	Cluster 1		Cluster 2	
	PA vs. NA	Drive for thinness hypothesis	Pleased vs. NA	PA & physical activity
Pleased (PA)	.71	.09	.47	.45
Happy (PA)	.72	.07	.34	.53
Appreciated (PA)	.48	.08	-.08	.73
Love (PA)	.41	.12	.09	.69
Sad (NA)	-.76	-.05	-.75	.01
Angry (NA)	-.69	.00	-.32	.03
Lonely (NA)	-.65	-.01	-.68	-.04
Anxious (NA)	-.59	.06	-.30	-.06
Tense (NA)	-.66	.10	-.51	-.22
Guilty (NA)	-.55	.04	-.31	.02
Irritated (NA)	-.52	.06	-.36	-.08
Ashamed (NA)	-.53	.01	-.09	.05
Feel fat (DT)	-.43	.28	-.15	-.19
Feel ugly (DT)	-.46	.18	-.23	.01
Want to burn calories (DT)	-.19	.61	.00	-.04
Want to lose weight (DT)	-.30	.34	-.08	.01
Want to be active (UPhA)	-.07	.78	.11	-.12
Want to move (UPhA)	-.07	.77	.12	-.11
Want to sport (UPhA)	-.07	.73	-.02	-.09
Am active (PhA)	.15	.61	-.03	.46
Am moving (PhA)	.11	.61	-.06	.51
Am sporting (PhA)	.16	.48	-.02	.17

Table 6

Normalized Varimax rotated loadings for the two clusters of the dilemma game data.

Loadings greater than +/- .35 are highlighted in bold face. PA is positive affect, NA is negative affect.

	Equal conditions		Unequal conditions	
	PA vs. NA	Empathy	PA vs. NA	Empathy
Angry	-.84	-.06	-.76	.00
Hurt	-.77	.01	-.47	.20
Irritated	-.84	.14	-.79	-.02
Annoyed	-.85	.10	-.87	-.10
Frustrated	-.85	.15	-.83	-.10
Dissappointed	-.84	.09	-.75	-.00
Indignant	-.81	.15	-.76	-.20
Enraged	-.80	.10	-.74	-.05
Hostile	-.75	-.04	-.61	-.01
Fearful	-.77	.16	-.17	.49
Surprised	-.29	.46	-.50	.01
Relieved	.28	.54	.43	.51
Elated	.57	.39	.27	.62
Warm	.44	.47	.30	.51
Happy	.72	.28	.64	.35
Satisfied	.77	.23	.63	.31
Sympathetic	.16	.25	.03	.62
Compassionate	-.34	.57	.02	.61
Tenderhearted	-.21	.72	-.07	.74
Concerned	-.37	.52	-.14	.55
Tender	.09	.76	.22	.42

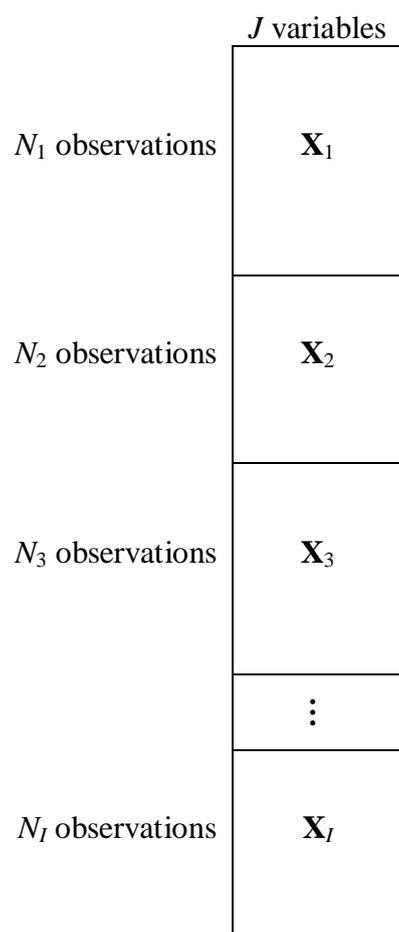


Figure 1. Graphical presentation of multivariate multiblock data.

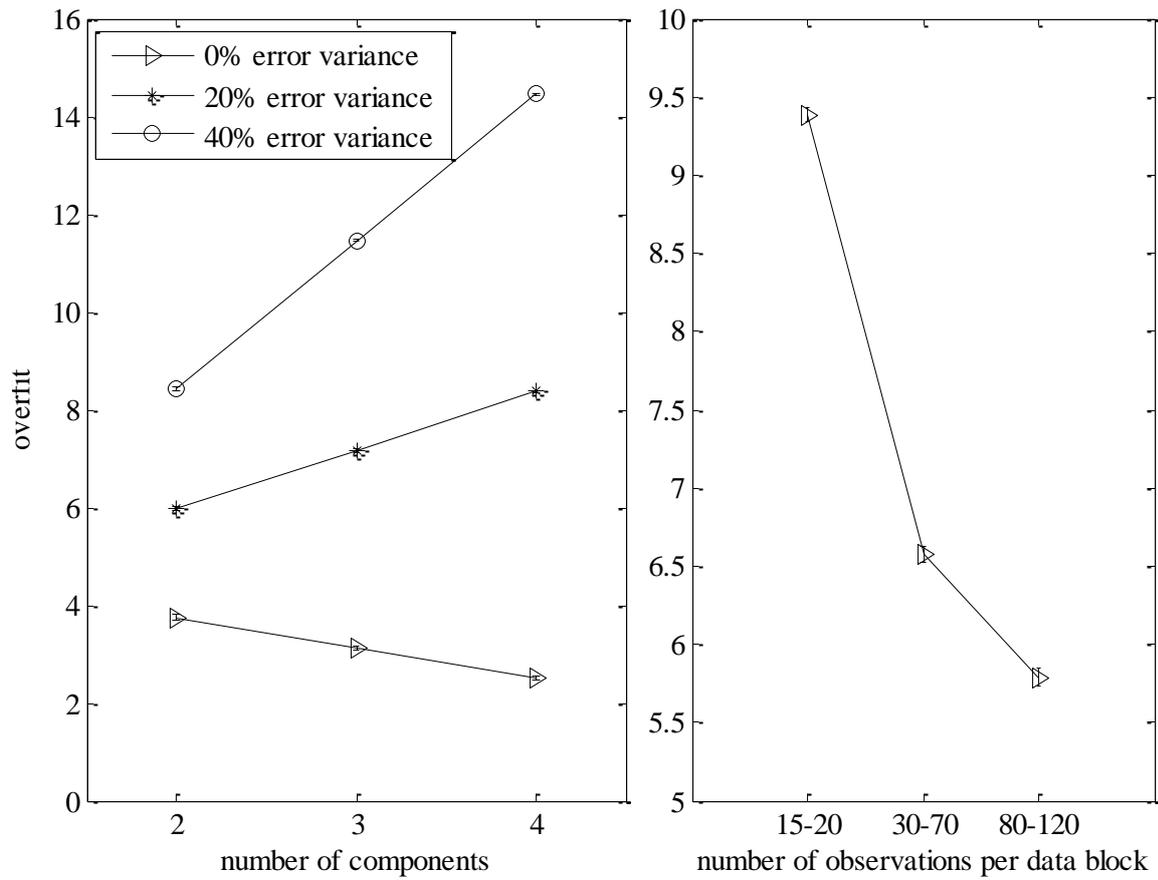


Figure 2. Mean overfit values and associated 95% confidence intervals as a function of the number of components and the amount of error variance ($e \times 100\%$) (a) and as a function of the number of observations per data block (b).

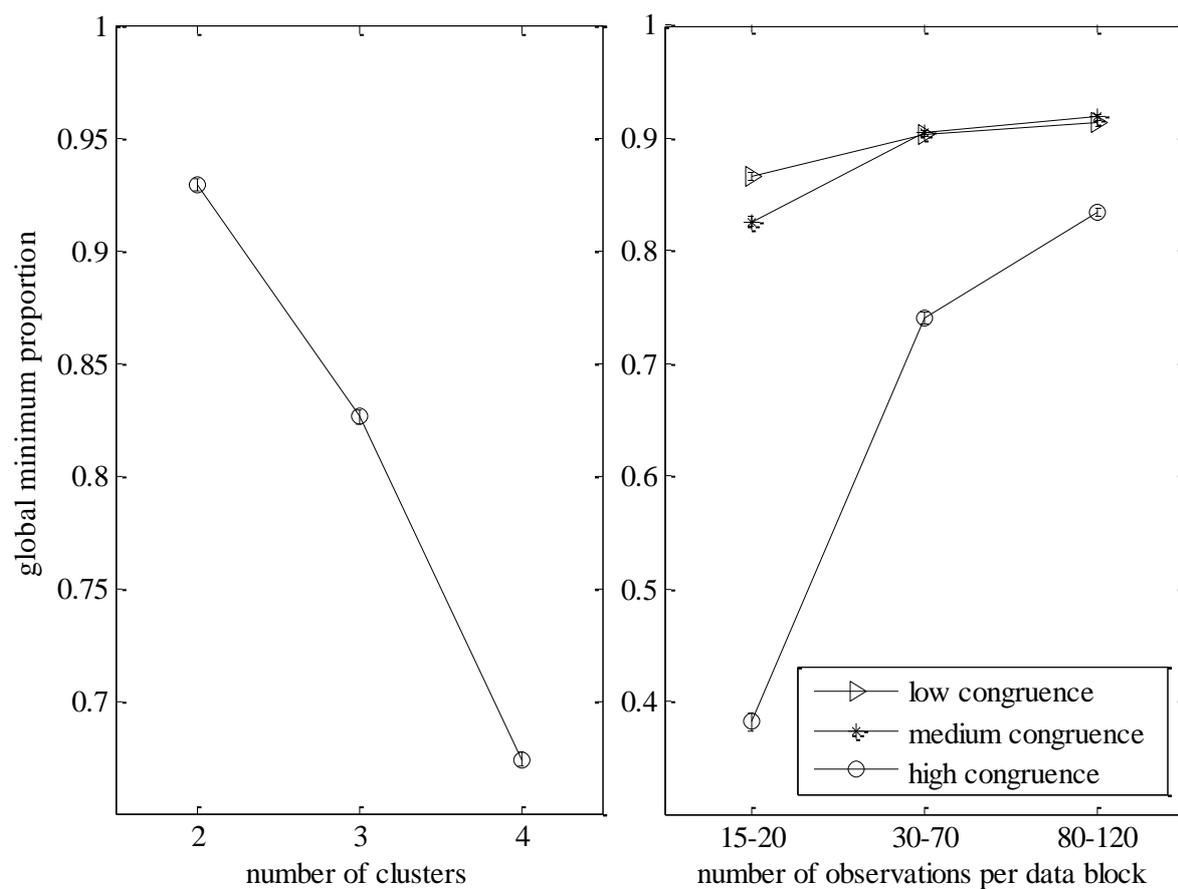


Figure 3. Mean values and associated 95% confidence intervals of the proportion of random runs with a loss function value equal to that of the proxy of the global minimum ('global minimum proportion') as a function of the number of clusters (a) and as a function of the number of observations per data block and the congruence of the cluster loading matrices (b).

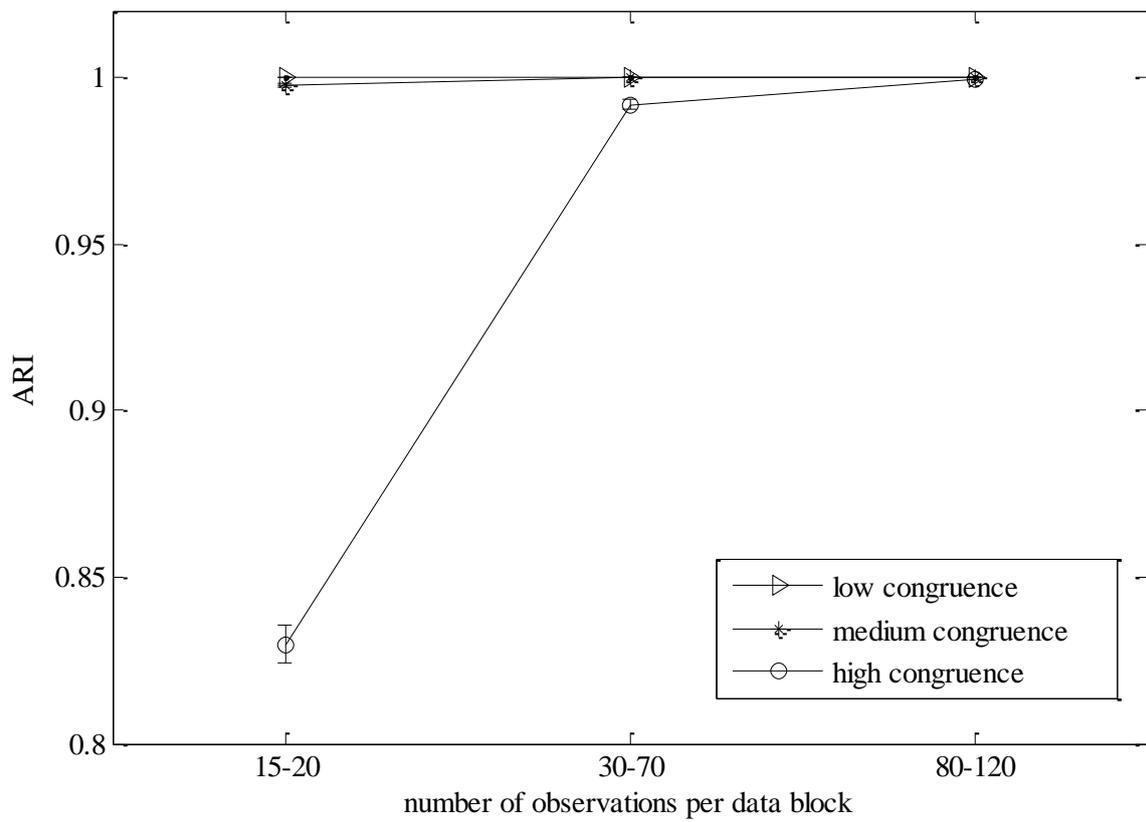


Figure 4. Mean ARI and associated 95% confidence intervals as a function of the number of observations per data block and the congruence of the cluster loading matrices (low, medium or high).

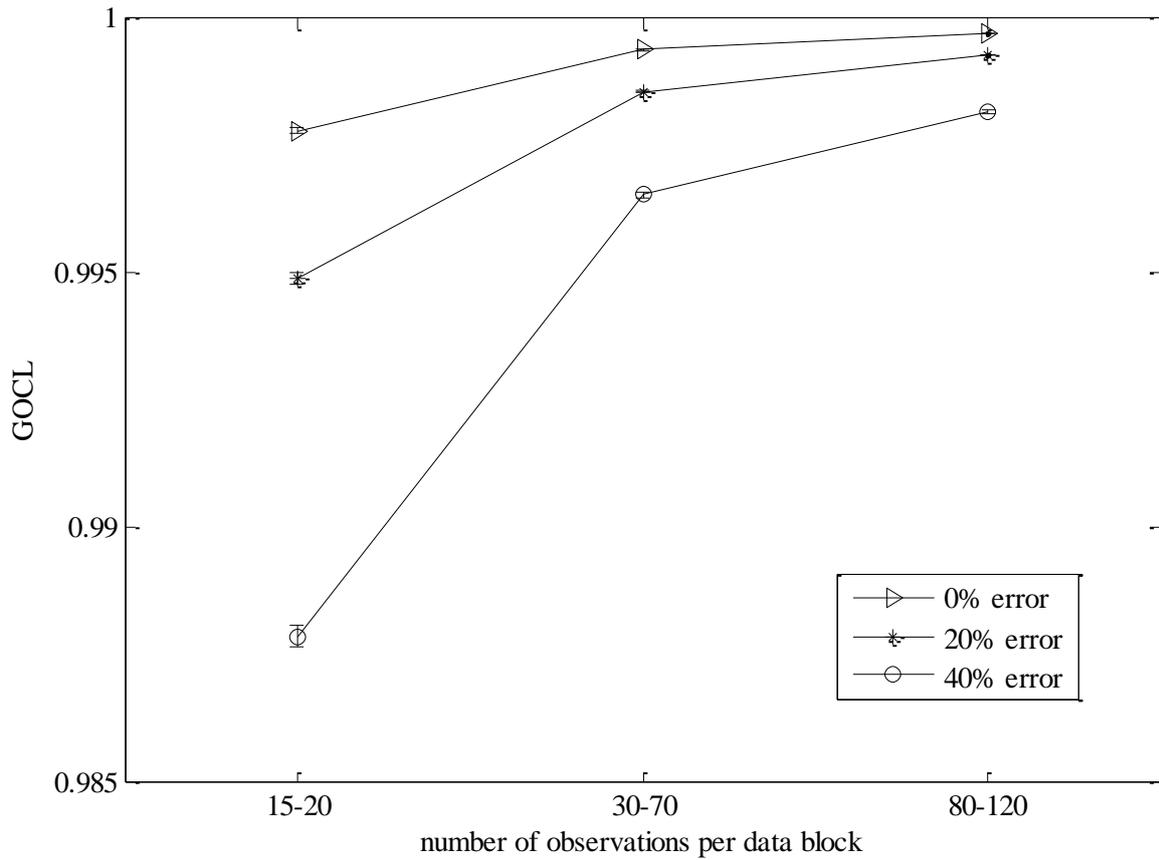


Figure 5. Mean *GOCL* and associated 95% confidence intervals as a function of the number of observations per data block and the amount of error variance ($e \times 100\%$).

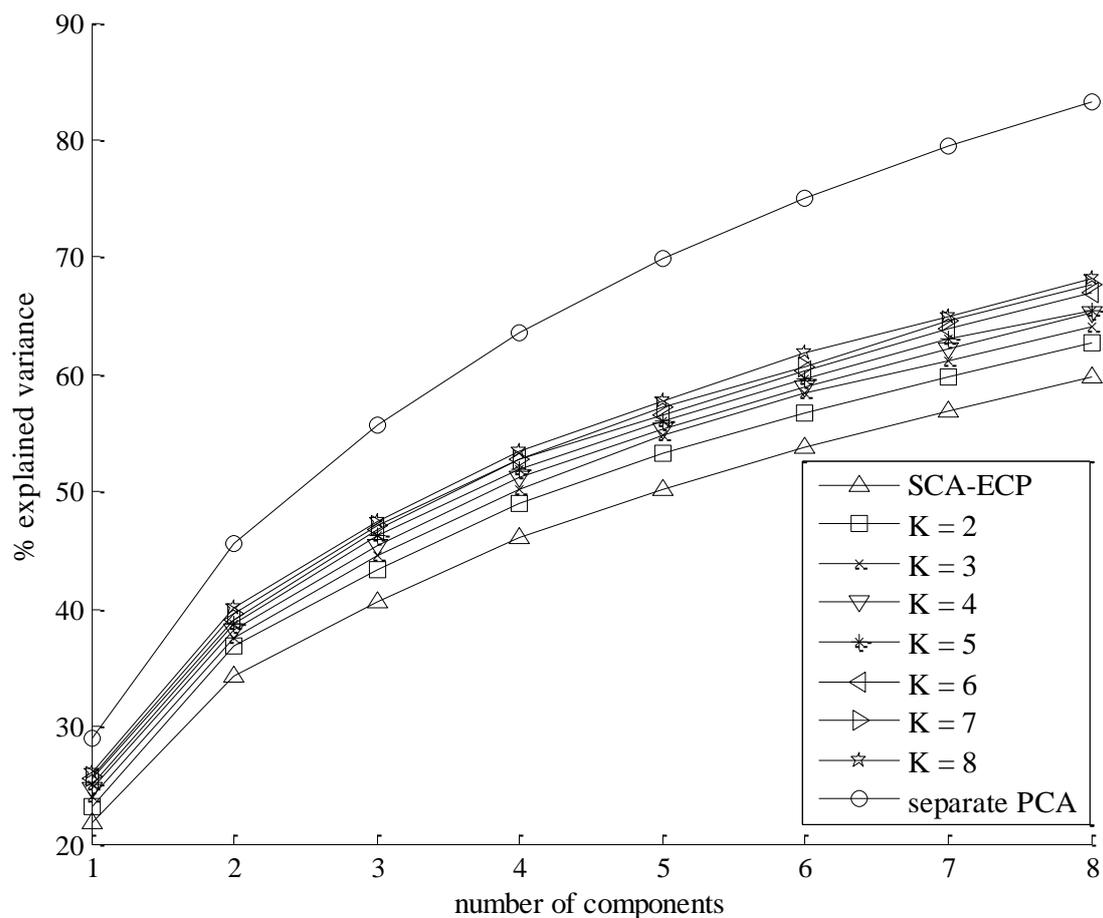


Figure 6. Percentage of explained variance for separate PCA, SCA-ECP and Clusterwise SCA-ECP solutions for the eating disorder data, with the number of components varying from 1 to 8, and the number of clusters for Clusterwise SCA-ECP varying from 2 to 8.

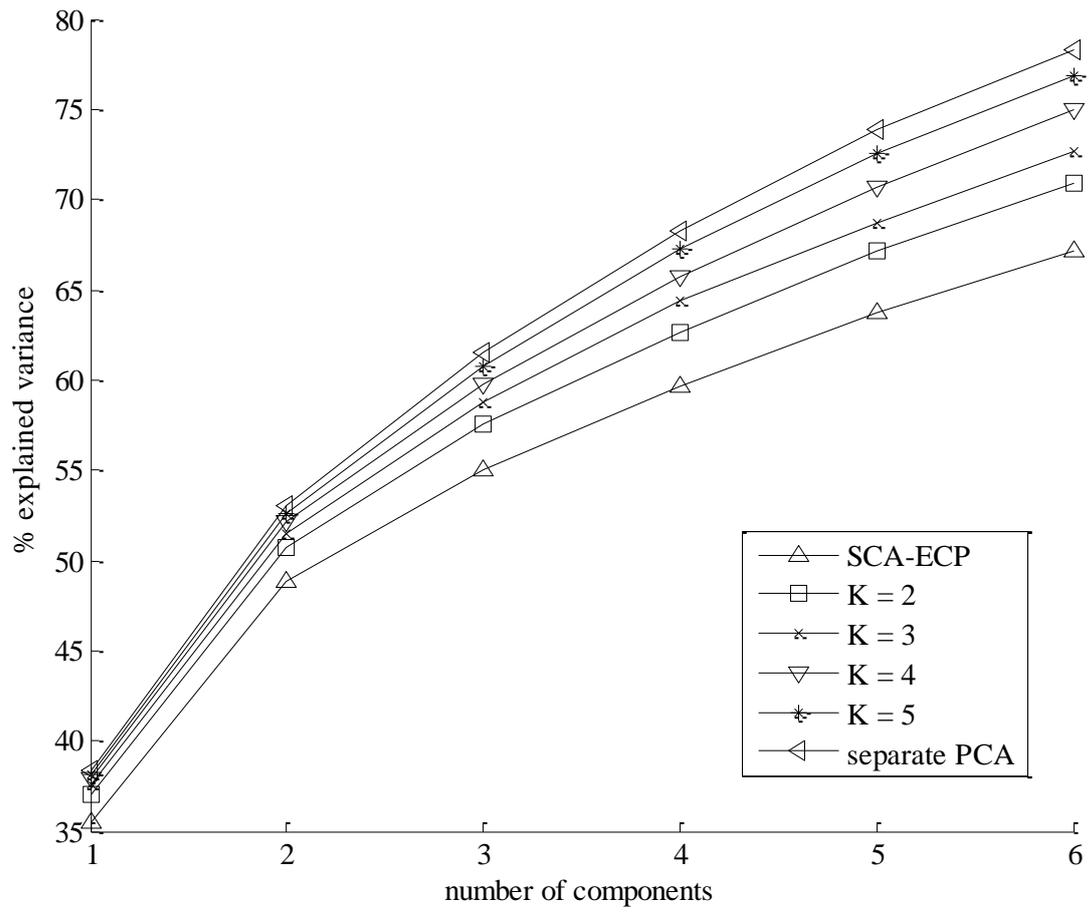


Figure 7. Percentage of explained variance for separate PCA, SCA-ECP and Clusterwise SCA-ECP solutions for the dilemma game data, with the number of components varying from 1 to 6, and the number of clusters for Clusterwise SCA-ECP varying from 2 to 5.