

Research article: The *SIMCLAS* model: Simultaneous analysis of coupled binary data matrices with noise heterogeneity between and within data blocks¹

Authors: Wilderjans, T. F.^a, Ceulemans, E.^b, & Van Mechelen, I.^a

Affiliation:

^a Research Group of Quantitative Psychology and Individual Differences, Department of Psychology, Katholieke Universiteit Leuven, Belgium

^b Department of Educational Sciences, Katholieke Universiteit Leuven, Belgium

Contact information:

Tom F. Wilderjans

Research Group of Quantitative Psychology and Individual Differences

Katholieke Universiteit Leuven

Tiensestraat 102, box 3713

3000 Leuven, Belgium

E-mail: tom.wilderjans@psy.kuleuven.be

Telephone: +32. 16. 32.61.23 or Fax: +32. 16. 32.59.93

¹ Requests for reprints should be sent to Tom F. Wilderjans. The first author is a Research Assistant of the Fund for Scientific Research - Flanders (Belgium). The research reported in this paper was partially supported by the Research Council of K.U.Leuven (GOA/2005/04 and EF/2005/07, 'SymBioSys') and by IWT-Flanders (SBO 60045, 'Bioframe'). We would like to thank Gert Storms and his collaborators for providing us with an interesting data set.

The *SIMCLAS* model: Simultaneous analysis of coupled binary data matrices with noise heterogeneity between and within data blocks

Abstract

In many research domains different pieces of information are collected regarding the same set of objects. Each piece of information constitutes a data block, and all these (coupled) blocks have the object mode in common. When analyzing such data, an important aim is to obtain an overall picture of the structure underlying the whole set of coupled data blocks. A further challenge consists of accounting for the differences in information value that exist between and within (i.e., between the objects of a single block) data blocks. To tackle these issues, analysis techniques may be useful in which all available pieces of information are integrated and in which at the same time noise heterogeneity is taken into account. For the case of binary coupled data, however, only methods exist that go for a simultaneous analysis of all data blocks, but that do not account for noise heterogeneity. Therefore, in this paper, the *SIMCLAS* model, being a Hierarchical Classes model for the simultaneous analysis of coupled binary two-way matrices, is presented. In this model, noise heterogeneity between and within the data blocks is accounted for by downweighting entries from noisy blocks/objects within a block. In a simulation study it is shown (1) that the *SIMCLAS* technique recovers the underlying structure of coupled data to a very large extent, and (2) that the

1
2
3
4
5
6
7
8 *SIMCLAS* technique outperforms a Hierarchical Classes technique in which
9
10 all entries contribute equally to the analysis (i.e., noise homogeneity within
11
12 and between blocks). The latter is also demonstrated in an application of both
13
14 techniques to empirical data on categorization of semantic concepts.
15
16
17

18 Key words: data fusion, coupled data, multi-set data, noise heterogeneity,
19
20 simultaneous clusterings, Hierarchical Classes Analysis, overlapping clustering,
21
22 hierarchical relations, multivariate binary data
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Introduction

In many research domains different pieces of information are collected regarding the same set of objects. Each piece of information constitutes a data block, and all these (coupled) blocks have the object mode in common. For example, in psychiatric diagnosis research, such coupled data are encountered when for a set of patients information is available regarding, on the one hand, their diagnosis (i.e., patient by diagnosis data block), and, on the other hand, the symptoms they exhibit (i.e., patient by symptom data block).

To grasp the wealth of information that is present in coupled data, one may uncover the structure underlying each data block, and relate these structures to one another. As such, an overall picture of the structure underlying the whole set of coupled data blocks is acquired. For example, applied to the psychiatric diagnosis case, one may search for the underlying syndromes that may account for both the diagnosis and symptom profiles of the different patients.

A complicating factor herewith may be that the coupled data blocks under study and/or the objects within a single data block may differ regarding their information value. For example, the patient by symptom data block may be more reliable than the patient by diagnosis block, because it is more easy to determine whether or not a symptom is present than to give a full diagnosis (which is a complicated interplay of different symptoms). Moreover, some patients may be more easy to diagnose than others, because of the very typical symptom profile they exhibit. Taking these differences in information value into account may result in stronger inferences with

1
2
3
4
5
6
7
8 respect to the structure underlying the data. However, in most cases no information is
9
10 available regarding the information value of a block/object within a block.

11
12 To handle these tasks, an analysis technique needs to be developed in which all
13
14 available pieces of information are simultaneously integrated and in which noise
15
16 heterogeneity between and within data blocks is taken into account (see, Van Mechelen
17
18 & Smilde, 2009; Van Mechelen & Smilde, 2010). When dealing with real-valued coupled
19
20 data and noise heterogeneity between blocks, one may rely on the *MxLSCA – P*
21
22 approach (Wilderjans, Ceulemans, Van Mechelen, & van den Berg, 2011), a stochastic
23
24 extension of standard simultaneous component analysis (Kiers & Ten Berge, 1994;
25
26 Timmerman & Kiers, 2003; Van Deun, Smilde, van der Werf, Kiers, & Van Mechelen,
27
28 2009). However, in case of binary coupled data, which show up rather frequently in
29
30 psychology, no such technique exists yet. Indeed, although techniques have been
31
32 proposed, like, for example, *concatenated HICLAS*, which can be conceived as a
33
34 Hierarchical Classes counterpart of simultaneous component analysis (Millsap &
35
36 Meredith, 1988; Kiers & Ten Berge, 1989; Ten Berge, Kiers, & Van der Stel, 1992) for
37
38 binary coupled data, these techniques do not account for noise heterogeneity between
39
40 and/or within the blocks.

41
42 To remedy this, the aim of the present paper is to propose a new model and
43
44 associated data-analytic strategy to handle between and within block noise
45
46 heterogeneity in binary coupled data. This model belongs to the family of Hierarchical
47
48 Classes models (De Boeck & Rosenberg, 1988; Van Mechelen, De Boeck, & Rosenberg,
49
50 1995; Leenen, Van Mechelen, De Boeck, & Rosenberg, 1999; Ceulemans, Van Mechelen,
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8 & Leenen, 2003; Ceulemans & Van Mechelen, 2005), and will therefore be called a
9
10 *SIMultaneous HICLAS* model (*SIMCLAS*). Specifically, we will introduce two
11
12 variants of the *SIMCLAS* model. To deal with noise heterogeneity between data
13
14 blocks we will propose *block-homogeneous SIMCLAS*, whereas *block-heterogeneous*
15
16 *SIMCLAS* allows for noise heterogeneity within (and between) data blocks. The
17
18 data-analytic strategy will imply estimating the amount of noise in each data block (in
19
20 the *block-homogeneous* variant) or the amount of noise in each object of each data
21
22 block (in the *block-heterogeneous* variant), and downweighting entries from more noisy
23
24 data blocks/objects in the analysis.
25
26
27
28
29

30
31 The remainder of this paper is organized in five main sections. In the first Section,
32
33 we recapitulate *concatenated HICLAS* analysis and introduce *SIMCLAS* analysis
34
35 (both the *block-homogeneous* and *block-heterogeneous* variant). In the second Section,
36
37 the *SIMCLAS* loss function will be presented, along with an algorithm to fit the
38
39 *SIMCLAS* model to data. In the third Section, by means of an extensive simulation
40
41 study, the optimization and recovery performance of the *SIMCLAS* technique will be
42
43 evaluated and the *SIMCLAS* technique will be compared to the *concatenated*
44
45 *HICLAS* technique. In the fourth Section, both techniques will be applied to empirical
46
47 coupled data on categorization of semantic concepts. The fifth Section, finally, contains
48
49 some concluding remarks.
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

The *SIMCLAS* analysis technique

We assume a data set consisting of N coupled $I \times J_n$ binary data matrices \mathbf{D}^n , with J_n differing across data blocks ($n = 1 \dots N$). Note that, without loss of generality, the rows are the common mode. When concatenating all N data blocks horizontally (Kiers, 2000), an $I \times J^*$ binary matrix \mathbf{D}^* is obtained, with $J^* = \sum_{n=1}^N J_n$. First, we will describe the *concatenated HICLAS* approach. Next, we will introduce the *block-homogeneous* and *block-heterogeneous* variants of the *SIMCLAS* technique.

The concatenated HICLAS approach

In the *concatenated HICLAS* approach for N data matrices¹ (for an introduction to a *HICLAS* analysis for a single data matrix, see De Boeck & Rosenberg, 1988), each binary $I \times J_n$ data matrix \mathbf{D}^n is approximated by a binary model matrix \mathbf{M}^n of the same size:

$$d_{ij_n}^n = m_{ij_n}^n + e_{ij_n}^n, \quad (1)$$

with $d_{ij_n}^n$, $m_{ij_n}^n$, and $e_{ij_n}^n$, denoting the entries of \mathbf{D}^n , \mathbf{M}^n , and the error matrix \mathbf{E}^n ($I \times J_n$), respectively. Further, each model matrix \mathbf{M}^n can be decomposed into a $I \times P$ binary matrix \mathbf{A} and an $J_n \times P$ binary matrix \mathbf{B}^n , where P denotes the rank of the model, and \mathbf{A} is the same for all model matrices \mathbf{M}^n . The P columns of \mathbf{A} and \mathbf{B}^n define P binary variables, which are called bundles, and imply an overlapping clustering of the elements of \mathbf{A} and \mathbf{B}^n .

¹In the remainder of this paper, the terms 'matrix' and 'block' are used interchangeably.

The *concatenated HICLAS* model represents two types of structural relations in \mathbf{M}^n ($n = 1 \dots N$): association and quasi-order. To explain these structural relations, we will use the three hypothetical coupled model matrices \mathbf{M}^1 , \mathbf{M}^2 , and \mathbf{M}^3 in Table 1 as a guiding example. In Table 2, the bundle matrices \mathbf{A} , \mathbf{B}^1 , \mathbf{B}^2 , and \mathbf{B}^3 of the *concatenated HICLAS* model with three bundles for \mathbf{M}^1 , \mathbf{M}^2 , and \mathbf{M}^3 (in Table 1) are presented.

=====
 Insert Table 1 about here
 =====

Association. The association relation is the binary relation between the row and column elements as defined by the 1-entries in \mathbf{M}^n . The *concatenated HICLAS* model represents this association relation by the following decomposition rule:

$$m_{ij_n}^n = \bigoplus_{p=1}^P a_{ip} b_{j_n p}^n, \quad (2)$$

where \bigoplus denotes a Boolean sum (i.e., $1 \bigoplus 1 = 1$), and a_{ip} and $b_{j_n p}^n$ the entries of \mathbf{A} and \mathbf{B}^n , respectively. In Table 1, for example, R_4 and C_5^3 are associated in \mathbf{M}^3 because both elements belong to the second bundle Bu_2 in \mathbf{A} and \mathbf{B}^3 (see Table 2), respectively.

However, R_4 and C_4^3 are not associated in \mathbf{M}^3 (see Table 1) because these elements do not have a bundle in common (see Table 2).

=====
 Insert Table 2 about here
 =====

Quasi-order. On the columns of \mathbf{M}^n , a quasi-order relation \leq is defined as follows:

When for column j the set of associated rows in \mathbf{M}^n is denoted by S^j , then column $j \leq$ column j' in \mathbf{M}^n iff $S^j \subseteq S^{j'}$. The *concatenated HICLAS* model represents the quasi-order relation among the columns in \mathbf{M}^n by subset-superset relations among their bundle patterns in \mathbf{B}^n . For example, in Table 1, one can see that $C_3^2 \leq C_4^2$ in \mathbf{M}^2 because $S^{C_3^2} \subseteq S^{C_4^2}$; therefore, in Table 2 the bundle pattern for C_3^2 is a subset of the bundle pattern for C_4^2 . Also on the (common) rows a quasi-order relation \leq is defined, where for row i S^i consists of the associated column elements from all N matrices \mathbf{M}^n (e.g., $S^{R_3} = \{C_2^1, C_4^1, C_2^2, C_1^3, C_2^3\}$). For example, in Table 1, $R_3 \leq R_4$ because $S^{R_3} \subseteq S^{R_4}$; consequently, in Table 2, the bundle pattern for R_3 is a subset of the bundle pattern for R_4 . Note that the quasi-order relation among the elements of a mode implies a partition of these elements into classes (i.e., elements with an identical bundle pattern), and a hierarchical relation among these classes (for more information, see De Boeck & Rosenberg, 1988).

Note that in a *concatenated HICLAS* analysis the noise in the data is not modeled in an explicit way, as is true for all Hierarchical Classes models. An exception to this, in case of a single binary data block, is the *probabilistic HICLAS* model of Leenen, Van Mechelen, Gelman, and De Knop (2008). This model may be applied to binary coupled data by concatenating all data blocks into a single block, which implies that the noise is modeled in the same way across all data blocks (i.e., noise homogeneity within and between blocks).

Block-homogeneous SIMCLAS for noise heterogeneity between blocks

The *block-homogeneous SIMCLAS* approach allows to model noise heterogeneity between data blocks. Specifically, in a *block-homogeneous SIMCLAS* analysis it is assumed that the absolute values of the noise entries $e_{ij_n}^n$ in (1) follow a Bernoulli distribution with parameter π_n (i.e., $|e_{ij_n}^n| \sim \text{Bern}(\pi_n)$). To account for possible noise heterogeneity, the parameters π_n are allowed to differ between data blocks, implying a different amount of noise for each block; the parameters π_n , restricted to be smaller than .50², equal the probability that a data entry $d_{ij_n}^n$ differs from the corresponding model entry $m_{ij_n}^n$. It is further assumed that all $e_{ij_n}^n$ are sampled independently (i.e., uncorrelated), resulting in all data entries $d_{ij_n}^n$ being locally independent. Note that in case of a single data block (i.e., $N = 1$), the *block-homogeneous SIMCLAS* model reduces to the *probabilistic HICLAS* model of Leenen et al. (2008).

Block-heterogeneous SIMCLAS for noise heterogeneity within and between blocks

To take noise heterogeneity within data blocks into account, the *block-homogeneous SIMCLAS* model may be further extended to the *block-heterogeneous SIMCLAS* model by allowing the parameter π to differ among the rows of each data block (and, as a consequence, also among the data blocks). In particular, it is assumed that the noise entries $e_{ij_n}^n$ are uncorrelated and that $|e_{ij_n}^n| \sim \text{Bern}(\pi_{in})$, with $\pi_{in} \leq .50$ ($i = 1 \dots I$; $n = 1 \dots N$).

²A value of π_n larger than .50 is not realistic, because this implies that model entries may differ from their corresponding data entries with a probability that is above chance.

Data analysis

Aim

Given a set $\tilde{\mathbf{D}}$ of binary coupled data blocks $\mathbf{D}^1, \dots, \mathbf{D}^n$, and a number of bundles P , the aim of a *block-heterogeneous SIMCLAS* analysis is to estimate \mathbf{A} , \mathbf{B}^n , and π_{in} such that the value of the (log)likelihood function

$$l = \sum_{i=1}^I \sum_{n=1}^N \left[\sum_{j_n=1}^{J_n} \left(d_{ij_n}^n - \bigoplus_{p=1}^P a_{ip} b_{j_n p}^n \right)^2 \times \log \left(\frac{\pi_{in}}{1 - \pi_{in}} \right) + J_n \times \log(1 - \pi_{in}) \right] \quad (3)$$

is maximized. In case of a *block-homogeneous SIMCLAS* analysis (3) reduces to:

$$l = \sum_{n=1}^N \left[\sum_{i=1}^I \sum_{j_n=1}^{J_n} \left(d_{ij_n}^n - \bigoplus_{p=1}^P a_{ip} b_{j_n p}^n \right)^2 \times \log \left(\frac{\pi_n}{1 - \pi_n} \right) + I J_n \times \log(1 - \pi_n) \right]. \quad (4)$$

Note that (4) implies that the extent to which $d_{ij_n}^n$ contributes to the likelihood function is inversely related to the amount of noise that is present in the corresponding \mathbf{D}^n ; as such, data entries of noisy data matrices are downweighted in favor of entries from less noisy matrices³. Analogously, (3) implies that data entries belonging to noisy rows of a data matrix are downweighted in favor of entries from less noisy rows.

³When $\pi_1 > \pi_2$, then $c_1 = \log \left(\frac{\pi_1}{1 - \pi_1} \right) > c_2 = \log \left(\frac{\pi_2}{1 - \pi_2} \right)$, with c_1 and c_2 representing the contribution of entries from \mathbf{D}^1 and \mathbf{D}^2 to the likelihood, which has to be maximized, respectively. Note that c_1 and c_2 are negative when $\pi_n \leq .50$, which implies that $|c_1| < |c_2|$. As such, entries from more noisy blocks (i.e., larger π_n and smaller $|c_n|$) imply a smaller decrease in the likelihood than entries from less noisy blocks (i.e., smaller π_n and larger $|c_n|$), resulting in entries from more noisy blocks being downweighted.

The SIMCLAS algorithm

To estimate the *SIMCLAS* model parameters (i.e., the bundle matrices and the noise parameters) that maximize the likelihood function, an iterative algorithm is adopted that consists of the following five steps:

1. Generate initial estimates for the noise parameters π_n ($n = 1 \dots N$), in case of the *block-homogeneous* variant, or π_{in} ($i = 1 \dots I; n = 1 \dots N$) in case of the *block-heterogeneous* variant. For π_n , three types of initial estimates may be obtained (for a discussion of different types of starting procedures, see Ceulemans, Van Mechelen, & Leenen, 2007): (1) rational, (2) random, or (3) smart-random (also called pseudo-rational). To obtain rational starting values for the noise parameters (in each block), first (1) a *HICLAS* analysis (De Boeck & Rosenberg, 1988) on \mathbf{D}^* , or (2) a weighted *HICLAS* analysis on \mathbf{D}^* , with the weight for each data matrix being $\frac{1}{IJ_n}$ (i.e., the inverse of the number of entries in \mathbf{D}^n), implying that all data matrices contribute equally to the analysis, is performed. Next, for each block, the average squared residual between the data and the model entries (see also Step 3 of the *SIMCLAS* algorithm) is computed. Random starting values may be generated by independently sampling values from a $U(0, .5)$ distribution. Smart-random starting values can be obtained by adding randomly sampled numbers to the rational starting values described above, with these numbers being drawn from a $U(-\delta, \delta)$ distribution with δ being equal to the corresponding rational starting value divided by five. Initial starting values for π_{in} can be obtained in a similar way.

2. Estimate \mathbf{A} and \mathbf{B}^n , conditionally upon the noise parameters obtained in Step 1, by means of a simulated annealing procedure (*SA*; for an introduction to *SA*, see Kirkpatrick, Gelatt, & Vecchi, 1983; Aarts, Korst, & van Laarhoven, 1997). Note that *SA* has already been applied successfully for parameter estimation of other Hierarchical Classes models (see, Ceulemans et al., 2007; Wilderjans, Ceulemans, & Van Mechelen, 2008). In the appendix, the principles and the different steps of such an *SA* algorithm are presented, along with the way in which *SA* is implemented in the *SIMCLAS* context. Because a *SA* algorithm may return a suboptimal solution (i.e., local optimum), a multi-start procedure is recommended (see appendix).

3. Re-estimate the noise parameters, conditional upon \mathbf{A} and \mathbf{B}^n (obtained in Step 2), as follows:

$$\hat{\pi}_n = \frac{\sum_{i=1}^I \sum_{j_n=1}^{J_n} (d_{ij_n}^n - \bigoplus_{p=1}^P a_{ip} b_{j_n p}^n)^2}{I J_n},$$

$$\hat{\pi}_{in} = \frac{\sum_{j_n=1}^{J_n} (d_{ij_n}^n - \bigoplus_{p=1}^P a_{ip} b_{j_n p}^n)^2}{J_n}. \quad (5)$$

4. Compute the loss function (4) or (3). When it has increased by a value larger than some tolerance value, repeat Steps 2 and 3, otherwise go to Step 5.

5. Perform a closure operation (Barbut & Monjardet, 1970; Birkhoff, 1940) on \mathbf{A} and \mathbf{B}^n . This is necessary because the bundle matrices obtained at the end of Step 4 do not yet represent the quasi-order relation in \mathbf{M}^n correctly. This closure operation consists of changing each 0-entry in \mathbf{A} and \mathbf{B}^n to 1 iff this modification does not alter \mathbf{M}^n (and consequently does not change the loss function value).

1
2
3
4
5
6
7
8 Because the *SIMCLAS* algorithm may end up in a local rather than the global
9
10 optimum, a multi-start procedure is recommended. In such a procedure, the
11
12 *SIMCLAS* algorithm is run a prespecified number of times, each time from different
13
14 initial estimates for the noise parameters (see Step 1); subsequently, the solution
15
16 yielding the maximum value on the likelihood function is retained. Note the difference
17
18 between the (overall) multi-start procedure for the *SIMCLAS* algorithm as a whole
19
20 (i.e., all five steps) and the multi-start procedure for the *SA* algorithm (i.e., the second
21
22 step of the *SIMCLAS* algorithm).
23
24
25
26
27
28

29 *Model selection*

30
31
32 In general, the (true) number of bundles P of the *SIMCLAS* model is unknown.
33
34 Therefore, in practice, analyses with increasing numbers of bundles are performed and a
35
36 heuristic is used to determine the model with optimal balance between model
37
38 complexity (i.e., the number of bundles), on the one hand, and the fit of the model to
39
40 the data (i.e., the value on the likelihood function), on the other hand (for more
41
42 information, see Ceulemans & Van Mechelen, 2005; Wilderjans et al., 2008). In this
43
44 paper, we will apply for this purpose the generalized scree test, which was proposed by
45
46 Leenen and Van Mechelen (2001) and which has shown good performance as model
47
48 selection heuristic for Hierarchical Classes Analyses (see, Leenen & Van Mechelen, 2001;
49
50 Ceulemans et al., 2003).
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Simulation study

Problem

In this section we present simulation results in which (1) we evaluate the performance of the *SIMCLAS* algorithm, and (2) we compare the *SIMCLAS* modeling technique to the *concatenated HICLAS* approach. At this point, we are interested in two aspects of algorithmic performance: optimization and recovery. With regard to optimization, we will examine the extent to which the *SIMCLAS* algorithm is able to maximize the loss function (4) and (3). Concerning recovery, we will determine the degree to which both approaches succeed in disclosing the true structure underlying a given coupled data set and whether or not the recovery depends on characteristics of the data.

In this section, we will present three simulation studies. In the first simulation study, we will deal with *block-homogeneous SIMCLAS* for the case of two coupled data blocks. Four data characteristics will be taken into account: (1) the size of the common (row) mode, (2) the size of the data blocks, (3) the total amount of noise in the data, and (4) the way the noise is distributed across data blocks. In the following two simulation studies, we try to generalize the results from the first study to case of more than two coupled data blocks, and to *block-heterogeneous SIMCLAS*.

First simulation study: Block-homogeneous SIMCLAS for two data blocks

Design and procedure

Design. We deal with the case of two coupled data matrices \mathbf{D}^1 and \mathbf{D}^2 with noise heterogeneity between data blocks. Four factors were systematically manipulated in a completely randomized factorial design, with all factors considered random. The first two factors, which pertain to the sizes of the data matrices, are:

- (a) the *Size of the common mode*, I , manipulated at three levels: $I = 200$, $I = 100$, $I = 50$;
- (b) the *Array/total ratio*, r :

$$r = \frac{I \times J_1}{(I \times J_1) + (I \times J_2)}, \tag{6}$$

which indicates how much the data blocks differ in size. This factor was manipulated at five levels: .10 (small \mathbf{D}^1 relative to \mathbf{D}^2), .30, .50 (equal size), .70, .90 (large \mathbf{D}^1).

Table 3 presents the sizes of \mathbf{D}^1 and \mathbf{D}^2 that result from an orthogonal combination of the factors Size of the common mode, I , and Array/total ratio, r . Note that the total $I \times (J_1 + J_2)$ was kept constant at a value of 20,000 entries.

=====
 Insert Table 3 about here
 =====

The other two factors, which pertain to the amount of noise in the data matrices, are:

- (c) the *Total amount of noise*, ε^{tot} , in the data, defined as $\varepsilon_1 + \varepsilon_2$, with ε_1 (ε_2) being the expected amount of noise in \mathbf{D}^1 (\mathbf{D}^2). This factor was manipulated at two levels: .20 and .40. Note the difference between ε (i.e., the expected amount of noise) and π (i.e., the estimated amount of noise);
- (d) the *Relative amount of noise*, ε^{rel} , defined as $\frac{\varepsilon_1}{\varepsilon^{tot}}$. This factor was manipulated at five levels: 0, .25, .50, .75, 1. Note that $\varepsilon^{rel} = 1$ ($\varepsilon^{rel} = 0$) implies that only \mathbf{D}^1 (\mathbf{D}^2) is subject to noise, and that $\varepsilon^{rel} = .50$ implies equal amounts of noise for both data blocks.

In Table 4, the values of ε_1 and ε_2 are displayed that are obtained by orthogonally combining the factors Total amount of noise, ε^{tot} , and Relative amount of noise, ε^{rel} .

=====
 Insert Table 4 about here
 =====

Procedure. For each combination of the levels of the four manipulated factors, a set $\tilde{\mathbf{D}}$ of coupled data matrices \mathbf{D}^n ($n = 1, 2$) was constructed as follows: True matrices $\mathbf{A}^{(T)}$ and $\mathbf{B}^{n(T)}$ with four bundles were generated by independently sampling entries from a Bernoulli distribution with a parameter value of .50. Next, the true matrices \mathbf{T}^n were calculated by combining $\mathbf{A}^{(T)}$ and $\mathbf{B}^{n(T)}$ by the *SIMCLAS* decomposition rule in (2). It should be noted that $\mathbf{A}^{(T)}$ and $\mathbf{B}^{n(T)}$ (see above) were generated subject to the constraint that all bundle-specific classes (i.e., a class of elements belonging to one bundle only) were non-empty. This constraint was imposed to ensure that the

$SIMCLAS$ decomposition of \mathbf{T}^n is unique upon a trivial permutation of the bundles (Ceulemans & Van Mechelen, 2003; Van Mechelen et al., 1995). Note that the true number of bundles was not manipulated here, because a pilot study revealed that this factor has no influence on the results. Next, for each \mathbf{T}^n , a data matrix \mathbf{D}^n was constructed by altering the entries of \mathbf{T}^n with a probability ε_n .

Per cell of the design, 20 data sets were generated, resulting in 20 (replications) \times 3 (Size of the common mode) \times 5 (Array/total ratio) \times 2 (Total amount of noise) \times 5 (Relative amount of noise) = 3,000 different data sets. These data sets were analyzed with a $SIMCLAS$ analysis given $P = 4$. A multi-start procedure was implemented using 15 starts (with 100 SA chains in step 2 of the algorithm): Two rational starts, five random starts, and eight smart-random starts (see Section *The SIMCLAS algorithm*). In addition, in order to compare the $SIMCLAS$ technique to the *concatenated HICLAS* approach, a *concatenated HICLAS* analysis with $P = 4$ was performed on each data set, which boils down to a $HICLAS$ analysis on \mathbf{D}^* . To this end, one may apply the second step of the $SIMCLAS$ algorithm to each data set, using 100 starts (i.e., SA chains, see appendix) with the loss function that has to be minimized being equal to the number of discrepancies between the data and the model.

Results

Goodness-of-fit. Studying the degree to which the $SIMCLAS$ algorithm succeeds in maximizing the likelihood function (4) or (3) is problematic because the optimal value for this function is unknown, except in the case of noise-free data. However,

because the true set $\tilde{\mathbf{T}}$ is always a valid *SIMCLAS* model for the coupled data set $\tilde{\mathbf{D}}$, the likelihood for $\tilde{\mathbf{T}}$ is a lower bound for the optimal value of the likelihood function. As a consequence, when the likelihood for $\tilde{\mathbf{T}}$ exceeds the likelihood for $\tilde{\mathbf{M}}$, we are sure that the algorithm got stuck in a local optimum. This was, however, never the case in our simulation.

Goodness-of-recovery. To address the question to which extent the *SIMCLAS* algorithm discloses the true structure underlying coupled data, $\mathbf{A}^{(M)}$ is compared to $\mathbf{A}^{(T)}$ by computing the kappa coefficient⁴ (Cohen, 1960), denoted as $\kappa_{\mathbf{A}}$. Because the bundles of a *SIMCLAS* solution can be freely permuted, the maximal $\kappa_{\mathbf{A}}$ was computed across all possible permutations. The $\kappa_{\mathbf{A}}$ statistic equals one when perfect recovery is encountered and zero when recovery is at chance level.

The mean $\kappa_{\mathbf{A}}$ value, across all simulated data sets, equals .997 ($SD = .01$), implying that the *SIMCLAS* technique recovers the underlying structure very well. To study how the recovery performance varies as a function of the data characteristics, an analysis of variance was performed with $\kappa_{\mathbf{A}}$ as the dependent variable and the data characteristics as independent variables. This analysis, only taking effects with a sizeable intraclass correlation $\hat{\rho}_I$ (Kirk, 1982; Haggard, 1958) into account (i.e., $\hat{\rho}_I \geq .10$), reveals that recovery decreases when the size of the common mode increases, with

⁴The kappa coefficient κ between two dichotomous variables can be computed as follows:

$$\kappa = \frac{(p_{00} + p_{11}) - (p_{0.}p_{.0} + p_{1.}p_{.1})}{1 - (p_{0.}p_{.0} + p_{1.}p_{.1})}, \quad (7)$$

with p_{00} (p_{11}) the proportion of zero-agreements (one-agreements) and $p_{0.}$ and $p_{1.}$ ($p_{.0}$ and $p_{.1}$) the marginal proportion of zeros and ones for the first (second) variable.

1
2
3
4
5
6
7
8 this effect being more pronounced when the amount of noise in the data is large ($\hat{\rho}_I =$
9
10 .15). This interaction, however, is qualified by two three-way interactions ($\hat{\rho}_I$ around
11
12 .13) and a complicated four-way interaction ($\hat{\rho}_I = .16$): The pattern of the two-way
13
14 interaction is more pronounced when the largest data matrix is subject to more noise
15
16 than the other one, or when both matrices contain a considerable (but equal) amount of
17
18 noise. However, perfect recovery is observed when one of both matrices is noise-free.
19
20
21

22
23 For each simulated data set, the $\kappa_{\mathbf{A}}$ statistic was also computed for the
24
25 *concatenated HICLAS* solution. The mean $\kappa_{\mathbf{A}}$ value for *concatenated HICLAS*
26
27 equals .98 ($SD = .06$), implying that, on average, *SIMCLAS* outperforms
28
29 *concatenated HICLAS*. An analysis of variance with the difference in $\kappa_{\mathbf{A}}$ between
30
31 *SIMCLAS* and *concatenated HICLAS* as the dependent variable and the data
32
33 characteristics as the independent variables, shows that *SIMCLAS* especially
34
35 outperforms *concatenated HICLAS* when the largest matrix is most noise-prone and
36
37 when the data contain a large (total) amount of noise ($\hat{\rho}_I = .71$).
38
39
40
41
42
43

44
45 *Second simulation study: Block-homogeneous SIMCLAS for more than two data*
46

47
48 *blocks*
49

50
51 *Design and procedure*
52

53
54 In this second simulation, we extended the first study to the case of ten coupled
55
56 data blocks. Because increasing the number of data blocks lengthens considerably the
57
58 computation time, we fixed (1) the *Size of the common mode, I*, at 50, (2) the
59
60 *Array/total ratio, r*, at .90 by setting the size of large and small data blocks at 50×90
61
62
63
64
65

and 50×10 , respectively, and (3) the *Total amount of noise*, ε_{tot} , which is the sum of ε_{larger} and $\varepsilon_{smaller}$, with ε_{larger} and $\varepsilon_{smaller}$ denoting the noise for the large and the small data blocks, respectively; ε_{tot} is set at .40. However, we also manipulated two factors in a completely randomized design. First, the *Dominance of the large data blocks*, defined as the number of large data blocks (50×90) relative to the total number of data blocks (i.e., 10), is manipulated at three levels: .2, .5, .8. Second, the *Noise level*, in which we manipulated the amount of noise for the larger and smaller data blocks at seven levels: 0 .40, .10 .30, .20 .40, .20 .20, .40 .20, .30 .10, and .40 0.

To generate 20 data sets per cell of the design the same procedure was used as before, resulting in 20 (replications) \times 3 (Number of large data blocks) \times 7 (Noise level) = 420 different data sets. Subsequently, a *SIMCLAS* analysis (with 15 multi-starts overall and 100 *SA* chains in step 2) and a *concatenated HICLAS* analysis (with 100 *SA* chains) with $P = 4$ was applied to each data set.

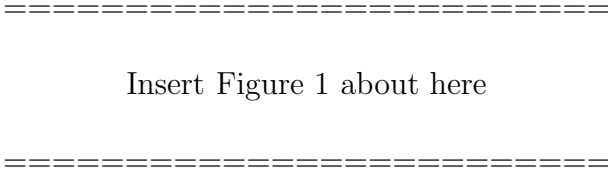
Results

With respect to goodness-of-fit, for all simulated data sets the likelihood for $\tilde{\mathbf{M}}$ exceeds the likelihood for $\tilde{\mathbf{T}}$. To compare *block-homogeneous SIMCLAS* to *concatenated HICLAS* in terms of recovery performance, we computed the $\kappa_{\mathbf{A}}$ value for the solutions obtained with both techniques. It appears that *block-homogeneous SIMCLAS* (mean of .98) outperforms *concatenated HICLAS* (mean of .94) to a larger extent than in the first study. Moreover, an analysis of variance revealed that this effect is more pronounced ($\hat{\rho}_I = .28$) when (1) the large data blocks contain all noise

(i.e., seventh level of *Noise level*), or (2) the large data blocks are more noise-prone and the total amount of noise in the data is large (i.e., fifth level of *Noise level*).

Interestingly, as can be seen in Figure 1, the latter effect becomes larger when the number of large data blocks relative to the total number of blocks increases ($\hat{\rho}_I = .49$).

Finally, when the small blocks are most noise-prone or in case of equal noise, both techniques perform equally well (and almost perfect, i.e., $\kappa_{\mathbf{A}} > .99$).



Third simulation study: Block-heterogeneous SIMCLAS for two data blocks

Design and procedure

In the third simulation study, we extended the first study to the case of two coupled data blocks with noise heterogeneity within and between blocks. In order to limit the number of conditions, the *Size of the common mode*, I , was fixed at 50, and the *Array/total ratio*, r , at .90, resulting in one large and one small block of size 50×90 and 50×10 , respectively. The *Noise level* was manipulated at the same seven levels as in the second study.

To simulate the data we used the same procedure as in the other two simulation studies, except the way in which noise was added to $\tilde{\mathbf{T}}$. Specifically, to obtain heterogeneous noise within data blocks, we selected 25 rows in \mathbf{T}^1 and \mathbf{T}^2 randomly to which a low amount of noise (ε_{low}^n) was added, the other 25 rows were perturbed with a

larger amount of noise (ε_{high}^n); for ε^n levels of .05, .10, .20, .30, and .35, ε_{low}^n and ε_{high}^n were respectively chosen as follows: .025 .075, .050 .150, .100 .300, .150 .450, and .250 .450. As such, conditions were obtained in which noise heterogeneity is present within and between data blocks. For each cell of the design, 20 replications were used, resulting in 20 (replications) \times 7 (Noise levels) = 140 simulated data sets.

Subsequently, a *SIMCLAS* analysis (with 15 multi-starts overall and 100 *SA* chains in step 2) and a *concatenated HICLAS* analysis (with 100 *SA* chains) with $P = 4$ was applied to each simulated data set.

Results

Regarding goodness-of-fit, for all data sets the likelihood for $\tilde{\mathbf{T}}$ does not exceed the likelihood for $\tilde{\mathbf{M}}$. In order to compare the recovery performance of both techniques, we calculated the $\kappa_{\mathbf{A}}$ value for each obtained *block-heterogeneous SIMCLAS* and *concatenated HICLAS* solution. The results show that *SIMCLAS* (mean of .90) clearly outperforms *concatenated HICLAS* (mean of .81) and this to a larger extent than in the first and second study. Further, this effect is more pronounced in these conditions where the largest block is subject to the largest amount of noise. In particular, the average difference in $\kappa_{\mathbf{A}}$ value between *SIMCLAS* and *concatenated HICLAS* in the fifth (i.e., $\varepsilon_1 = .40$ and $\varepsilon_2 = .20$), sixth (i.e., $\varepsilon_1 = .30$ and $\varepsilon_2 = .10$), and seventh (i.e., $\varepsilon_1 = .40$ and $\varepsilon_2 = 0$) condition of *Noise level* amounts to .14, .18, and .32; the mean *SIMCLAS* $\kappa_{\mathbf{A}}$ values in these conditions equals .61, .86, and .95. In the other conditions (i.e., equal amount of noise or the smallest block being most

noise-prone), both strategies perform equally well.

Discussion of the results

To understand why the *SIMCLAS* approach outperforms the *concatenated HICLAS* approach, note that the latter (deterministic) approach is equivalent to a stochastic version of model (2) in which it is assumed that the noise entries $|e_{ij_n}^n| \sim \text{Bern}(\pi)$ ($n = 1 \dots N$). As a consequence, in case coupled data matrices are subject to noise heterogeneity between and/or within data blocks, the concatenation approach, unlike *SIMCLAS*, is based on a misspecified and a too simple model. This also explains why *SIMCLAS* especially outperforms the concatenation approach in those conditions where the difference in noise between and within the data matrices is large, because in those conditions the misspecification of the concatenation approach becomes worse.

Illustrative application

In this section we will illustrate the superior performance of the *block-homogeneous SIMCLAS* technique over the *concatenated HICLAS* approach in terms of disclosing the 'true' structure underlying empirical coupled data. To this end, we will apply both techniques to a data set from the field of categorization of semantic concepts. In this field, one of the main research questions pertains to whether elements (i.e., exemplars) of different categories can be distinguished on the basis of their properties (i.e., features). The data set that we will analyze is a coupled exemplar by feature data set

1
2
3
4
5
6
7
8 from the Animal domain, which belongs to the Leuven Natural Concept Database
9
10 (De Deyne et al., 2008). The first data set was obtained by asking four persons whether
11
12 or not 60 animals (i.e., 30 exemplars of the category of mammals and 30 of the category
13
14 of birds) are characterized by 225 features, where these features resulted from a feature
15
16 generation task in which participants had to list typical features for each category of
17
18 animals. In the second data set, four persons rated the same 60 animals on a different
19
20 set of 764 features that also were obtained by a feature generation task, but now
21
22 participants had to list features that are typical for each animal. As such, the two data
23
24 sets are coupled via the exemplar (i.e., animal) mode. The data were dichotomized
25
26 using a majority rule (i.e., a 1 was recorded when two or more participants judged that
27
28 a feature characterizes an exemplar, and 0 was recorded otherwise).
29
30
31
32
33
34

35
36 Earlier analyses of these data (Ceulemans & Storms, 2010) showed that the two
37
38 categories under study (i.e., mammals and birds) clearly can be distinguished based on
39
40 the feature profiles. Moreover, the classification of the different animals into mammals
41
42 and birds conforms the biological knowledge about the hierarchical classification of the
43
44 animal world. To determine whether *SIMCLAS* or *concatenated HICLAS* best
45
46 unveils this 'true' underlying structure, both techniques were applied with two bundles
47
48 to the coupled data set (*SIMCLAS* with 15 multi-starts and 500 *SA* chains in step 2
49
50 of the algorithm and *HICLAS* with 500 *SA* chains, see Section *The SIMCLAS*
51
52 *algorithm*). The obtained bundle matrix for the animals was compared with the true
53
54 bundle matrix (i.e., partition of the animals) by computing the kappa coefficient
55
56 $\kappa_{animals}$. For both techniques, the $\kappa_{animals}$ value amounts to .97, indicating that the true
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8 structure is recovered almost perfectly. Specifically, except for one mammal that is
9
10 erroneously categorized as a bird, both techniques classify all animals correctly. It can
11
12 be concluded that for this data set *SIMCLAS* does not outperform *concatenated*
13
14 *HICLAS*. This result, however, is not surprising in this case as *concatenated*
15
16 *HICLAS* already recovers the underlying structure almost perfectly, leaving
17
18 *SIMCLAS* no room for improvement (i.e., a ceiling effect). Moreover, the same
19
20 amount of noise appears to be present in both data sets (i.e., the estimated π is .075
21
22 and .068 for the first and second data set, respectively). In the simulation study it was
23
24 shown that in such conditions (i.e., equal noise) both strategies perform equally well.
25
26
27
28
29

30
31 Therefore, we created a new coupled data set as follows: First, we computed for
32
33 each feature, based on a *HICLAS* analysis of the separate data sets, the proportion of
34
35 discrepancies between the data and the model, which may give an indication of the
36
37 amount of noise in the feature scores. Next, we selected in the first data set the features
38
39 with a proportion of discrepancies smaller than .005 (i.e., small amount of noise),
40
41 whereas in the second data set the features with a proportion of discrepancies larger
42
43 than .10 (i.e., large amount of noise) were retained; this resulted in 29 and 167 features
44
45 being selected, respectively. As such, we mimicked the simulation conditions in which
46
47 *SIMCLAS* outperforms *concatenated HICLAS* to a large extent (i.e., largest block
48
49 being most noise-prone). To this new coupled data set, we applied a *SIMCLAS* and a
50
51 *concatenated HICLAS* analysis (with the same number of starts as above) with two
52
53 bundles, and we computed the kappa coefficient for the obtained animal bundle matrix;
54
55 this resulted in *SIMCLAS* ($\kappa_{animals} = .97$) clearly outperforming *concatenated*
56
57
58
59
60
61
62
63
64
65

HICLAS ($\kappa_{animals} = .72$) in terms of disclosing the true partition of the animals.

Specifically, *concatenated HICLAS* misclassifies 9 out of the 60 animals, whereas

SIMCLAS misclassifies only two exemplars.

General discussion

In this paper, the *SIMCLAS* technique, together with an associated data-analytic strategy, was introduced for an integrated analysis of coupled binary data matrices that are subject to noise heterogeneity between and within data blocks. The key idea behind *SIMCLAS* is to downweight in the analysis entries from noisy data matrices (*block-homogeneous SIMCLAS*) or entries from noisy rows within a data matrix (*block-heterogeneous SIMCLAS*). In a simulation study, it was demonstrated that the *SIMCLAS* algorithm succeeds in optimizing the loss function, and recovers the true structure underlying the coupled data well. Moreover, *SIMCLAS* was shown to outperform *concatenated HICLAS*, in which all data entries contribute equally to the analysis. This effect is more pronounced (1) when the largest data matrices are more noise-prone, (2) when the number of large data matrices is large relative to the number of small matrices, and (3) when the noise is heterogeneous within data matrices. It can be expected that in case of more than two coupled data blocks with noise heterogeneity within and between blocks (or with less structured noise within blocks), *SIMCLAS* will outperform *concatenated HICLAS* to a larger extent.

In the remainder of this section, we discuss some generalizations of the *SIMCLAS* technique for, on the one hand, the same type of coupled data, and, on the other hand,

other types of coupled data.

Generalizations for the same type of data. In the *block-homogeneous SIMCLAS* model, the probability that a model entry differs from a data entry is constant within data blocks, whereas, in the *block-heterogeneous SIMCLAS* model, this probability is allowed to differ within blocks, but only in a specific way (i.e., between rows within a block). In some applications, however, other types of noise heterogeneity within data blocks may be encountered. One possibility is that the probability of obtaining a false positive (i.e., $d_{ij_n}^n = 1$ while $m_{ij_n}^n = 0$) may differ from the probability of a false negative (i.e., $d_{ij_n}^n = 0$ while $m_{ij_n}^n = 1$). For instance, in the psychiatric diagnosis example, it may be the case that some psychiatrists easily (or, in contrast, are rather reluctant to) assign a symptom to a patient; this may be because these psychiatrists believe that it is important not to miss relevant symptoms (or, in contrast, that it is harmful to erroneously assign a symptom). To account for this phenomenon within the context of a *HICLAS* analysis of a single two-mode binary data matrix, the *two – error probabilistic HICLAS* model was proposed by Leenen et al. (2008). A similar generalization could be considered for the *SIMCLAS* model, both for the homogeneous and the heterogeneous variant. To implement this, a separate parameter may be introduced for false positives and false negatives, with again these parameters being allowed to differ across data matrices (for the homogeneous variant) or across rows within the same data matrix (in case of *block-heterogeneous SIMCLAS*). However, estimating the associated noise parameters is not a trivial task.

Generalizations for other types of data. The *SIMCLAS* technique may also be

extended to the analysis of two (or more) three-way three-mode coupled data blocks that are subject to noise heterogeneity between and/or within blocks. As an example, take an emotion researcher who wants to study individual differences in situation-specific emotional reactions. To this end, the researcher may, on the one hand, register for a set of persons which physiological reactions they display in a number of situations, and, on the other hand, collect for the same sets of persons and situations information regarding to the appraisals that are activated by the situations. Note that this implies two three-way three-mode data blocks that share two modes (i.e., the persons and the situations). In this case, noise heterogeneity between blocks may occur when, for example, the appraisal data are more noise-prone than the physiological measures, because the former measure is more subjective than the latter. To analyze such data, one may adopt an analysis technique that generalizes the *SIMCLAS* approach to the case of coupled three-way data blocks. In particular, a new model may be introduced, which consists of coupled three-way *HICLAS* models (see Leenen et al., 1999; Ceulemans et al., 2003; Ceulemans & Van Mechelen, 2004, 2005) with a noise parameter for each data block that is allowed to differ across these blocks. Analogous extensions may be considered for the case of noise heterogeneity within data blocks.

**Appendix: Simulated annealing to estimating the bundle matrices,
conditional on the noise parameters**

To estimate, in Step 2 of the *SIMCLAS* algorithm (see Section *The SIMCLAS algorithm*), the binary bundle matrices \mathbf{A} and \mathbf{B}^n that maximize the loss function,

conditionally upon the noise parameters, a simulated annealing procedure is adopted.

Simulated annealing is a local search technique that implies a walk through the solution space. In particular, a chain of solutions, consisting of several subchains, is generated by each time creating a candidate solution based on the current solution. Next, the loss function values of the current and the candidate solution are compared. If the candidate solution has a better loss function value f , it is accepted, which implies that the current solution is replaced by the candidate solution. If the candidate solution, however, has a worse loss function value, it is accepted with a probability that depends on its relative quality (i.e., the difference in loss function value f between the current solution and the candidate one) and the current temperature, a quantity that controls the acceptance probability. At the end of each subchain the temperature is decreased. Subchains are generated until a prespecified subchain stop criterion is met. Finally, the best encountered solution in the chain is retained.

Based on the results of a pilot study and on the *SA* implementations that have been used for other Hierarchical Classes models (see, Ceulemans et al., 2007), we implemented the procedure for generating a single *SA* chain (see Algorithm 1 for pseudo-code) in the *SIMCLAS* algorithm as follows:

1. An initial solution $S_{current}$ and associated initial loss value $L_{current}$ is obtained by replacing the P columns of each bundle matrix by P data vectors sampled at random (i.e., for \mathbf{A} , column vectors are drawn from the different \mathbf{D}^n , whereas for each \mathbf{B}^n , row vectors are chosen from the corresponding \mathbf{D}^n).

2. The initial temperature $T_{initial}$ is obtained by running a subchain of solutions and accepting all solutions; subsequently, the average increase in the likelihood function across those links in which worse solutions are accepted, is divided by $\ln(0.8)$; as such, during the first subchains in which the algorithm is still far from the optimal solution, worse solutions are accepted with a high probability (see Kirkpatrick et al., 1983; Aarts et al., 1997; Ceulemans et al., 2007).
3. A candidate solution S_{trial} and associated loss value L_{trial} is obtained from the current solution $S_{current}$ by altering the value of a randomly chosen cell of a randomly chosen bundle matrix, with each cell of each bundle matrix having the same probability of being changed.
4. A worse candidate solution is accepted if: $p < \exp((L_{trial} - L_{current})/T_{current})$, with p being a number generated from a uniform(0,1) distribution.
5. A subchain stops (1) if the number of generated solutions i_{gen} equals the maximum number of solutions $CL = ((I + \sum_{n=1}^N J_n) \times 2^P) \times 5$, or (2) if the number of accepted solutions i_{acc} equals $CL \times .10$.
6. At the end of each subchain, the temperature is decreased by a factor $\alpha = .90$, implying a smaller acceptance probability for worse solutions: $T_{current} = 0.9 \times T_{current}$.
7. A SA chain stops when (1) the current temperature becomes smaller than $T_{stop} = 0.000001$, or (2) the number of subsequent subchains i_{id} with an identical loss value $L_{current}$ for the last accepted solution in each subchain (i.e., $L_{current} = L_{previous}$)

1
2
3
4
5
6
7
8 equals max_{id} , which is set to five.
9

10 8. The retained solution is the best encountered solution S_{best} across all subchains.
11
12

13
14 To lower the risk of ending in a suboptimal solution (i.e., local optimum), a
15 multi-start procedure may be advised, which consists of running 100 *SA* chains, each
16 time with a different initial solution and initial temperature (see Step 1 and 2), and
17 retaining the best encountered solution across all chains (see, Ceulemans et al., 2007).
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

```

1
2
3
4
5
6
7
8
9
10
11 Initialize( $S_{current}, L_{current}, T_{initial}, CL, L_{best}, L_{previous}$ );
12  $T_{current} := T_{initial}$ ;
13
14 repeat
15      $i_{gen} := 0$ ;
16      $i_{acc} := 0$ ;
17
18     while ( $i_{gen} < CL$ ) and ( $i_{acc} < (.1 \times CL)$ ) do
19          $i_{gen} := i_{gen} + 1$ ;
20         generate  $S_{trial}$  and associated  $L_{trial}$ , based on  $S_{current}$ ;
21         draw  $h$  from  $U(0, 1)$ 
22         if ( $L_{trial} > L_{current}$ ) or ( $h < \exp(\frac{L_{trial} - L_{current}}{T_{current}})$ ) then
23              $S_{current} := S_{trial}$ ;
24              $L_{current} := L_{trial}$ ;
25              $i_{acc} := i_{acc} + 1$ 
26             if  $L_{trial} > L_{best}$  then
27                  $S_{best} := S_{trial}$ ;
28                  $L_{best} := L_{trial}$ 
29             end if
30         end if
31     end while
32
33      $T_{current} := \alpha * T_{current}$ ;
34     if  $L_{current} = L_{previous}$  then
35          $i_{id} := i_{id} + 1$ 
36     else
37          $i_{id} := 1$ ;
38          $L_{previous} := L_{current}$ 
39     end if
40
41 until ( $T_{current} \leq T_{stop}$ ) or ( $i_{id} = max_{i_{id}}$ );
42
43 return  $S_{best}$ ;
44
45
46
47
48
49
50
51
52
53
54
55

```

Algorithm 1: Pseudo-code for generating a single SA chain in Step 2 of the $SIMCLAS$ algorithm

References

- Aarts, E. H. L., Korst, J. H. M., & van Laarhoven, P. J. M. (1997). Simulated annealing. In E. H. L. Aarts & J. K. Lenstra (Eds.), *Local search in combinatorial optimization* (pp. 91–120). Chichester, UK: Wiley.
- Barbut, M., & Monjardet, B. (1970). *Ordre et classification: Algèbre et combinatoire*. Paris: Hachette.
- Birkhoff, G. (1940). *Lattice theory*. Providence: American Mathematical Society.
- Ceulemans, E., & Storms, G. (2010). Detecting intra- and inter-categorical structure in semantic concepts using HICLAS. *Acta Psychologica, 133*, 296–304.
- Ceulemans, E., & Van Mechelen, I. (2003). Uniqueness of n -way n -mode hierarchical classes models. *Journal of Mathematical Psychology, 47*, 259–264.
- Ceulemans, E., & Van Mechelen, I. (2004). Tucker2 hierarchical classes analysis. *Psychometrika, 69*, 375–399.
- Ceulemans, E., & Van Mechelen, I. (2005). Hierarchical classes models for three-way three-mode binary data: Interrelations and model selection. *Psychometrika, 70*, 461–480.
- Ceulemans, E., Van Mechelen, I., & Leenen, I. (2003). Tucker3 hierarchical classes analysis. *Psychometrika, 68*, 413–433.
- Ceulemans, E., Van Mechelen, I., & Leenen, I. (2007). The local minima problem in

1
2
3
4
5
6
7
8 hierarchical classes analysis: An evaluation of a simulated annealing algorithm
9
10 and various multistart procedures. *Psychometrika*, *72*, 377–391.

11
12
13 Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and*
14
15 *Psychological Measurement*, *20*, 37–46.

16
17
18
19 De Boeck, P., & Rosenberg, S. (1988). Hierarchical classes: Model and data analysis.
20
21 *Psychometrika*, *53*, 361–381.

22
23
24
25 De Deyne, S., Verheyen, S., Ameel, E., Vanpaemel, W., Dry, M., Voorspoels, W., &
26
27 Storms, G. (2008). Exemplar by feature applicability matrices and other dutch
28
29 normative data for semantic concepts. *Behavioral Research Methods*, *40*,
30
31 1030–1048.

32
33
34
35
36 Haggard, E. A. (1958). *Intraclass correlation and the analysis of variance*. New York:
37
38 Dryden.

39
40
41
42 Kiers, H. A. L. (2000). Towards a standardized notation and terminology in multiway
43
44 analysis. *Journal of Chemometrics*, *14*, 105–122.

45
46
47
48 Kiers, H. A. L., & Ten Berge, J. M. F. (1989). Alternating least squares algorithms for
49
50 simultaneous components analysis with equal component weight matrices for all
51
52 populations. *Psychometrika*, *54*, 467–473.

53
54
55
56
57 Kiers, H. A. L., & Ten Berge, J. M. F. (1994). Hierarchical relations between methods
58
59 for simultaneous component analysis and a technique for rotation to a simple
60
61
62
63
64
65

- 1
2
3
4
5
6
7 simultaneous structure. *British Journal of Mathematical and Statistical*
8
9
10 *Psychology*, 47, 109–126.
11
12
13 Kirkpatrick, S., Gelatt, C. D. J., & Vecchi, M. P. (1983). Optimization by simulated
14
15 annealing. *Science*, 220, 671–680.
16
17
18
19 Kirk, R. E. (1982). *Experimental design: Procedures for the behavioral sciences* (2nd
20
21 ed.). Belmont, CA: Brooks/Cole.
22
23
24
25 Leenen, I., & Van Mechelen, I. (2001). An evaluation of two algorithms for hierarchical
26
27 classes analysis. *Journal of Classification*, 18, 57–80.
28
29
30
31 Leenen, I., Van Mechelen, I., De Boeck, P., & Rosenberg, S. (1999). INDCLAS: A
32
33 three-way hierarchical classes model. *Psychometrika*, 64, 9–24.
34
35
36
37 Leenen, I., Van Mechelen, I., Gelman, A., & De Knop, S. (2008). Bayesian hierarchical
38
39 classes analysis. *Psychometrika*, 73, 39–64.
40
41
42
43 Millsap, R. E., & Meredith, W. (1988). Component analysis in cross-sectional and
44
45 longitudinal data. *Psychometrika*, 53, 123–134.
46
47
48
49 Ten Berge, J. M. F., Kiers, H. A. L., & Van der Stel, V. (1992). Simultaneous
50
51 components analysis. *Statistica Applicata*, 4, 377–392.
52
53
54
55 Timmerman, M. E., & Kiers, H. A. L. (2003). Four simultaneous component models for
56
57 the analysis of multivariate time series from more than one subject to model
58
59 intraindividual and interindividual differences. *Psychometrika*, 68, 105–121.
60
61
62
63
64
65

- 1
2
3
4
5
6
7
8 Van Deun, K., Smilde, A. K., van der Werf, M. J., Kiers, H. A. L., & Van Mechelen, I.
9
10 (2009). A structured overview of simultaneous component based data integration.
11
12 *BMC Bioinformatics*, *10*, 246.
13
14
15
16 Van Mechelen, I., De Boeck, P., & Rosenberg, S. (1995). The conjunctive model of
17
18 hierarchical classes. *Psychometrika*, *60*, 505–521.
19
20
21
22 Van Mechelen, I., & Smilde, A. K. (2009). *A generic model for data fusion* (Paper
23
24 presented at The Sixth Edition of the ThRee-way methods in Chemistry and
25
26 Psychology -TRICAP- meeting).
27
28
29
30 Van Mechelen, I., & Smilde, A. K. (2010). A generic linked-mode decomposition model
31
32 for data fusion. *Chemometrics and Intelligent Laboratory Systems*, *104*, 83–94.
33
34
35
36
37 Wilderjans, T. F., Ceulemans, E., & Van Mechelen, I. (2008). The CHIC model: A
38
39 global model for coupled binary data. *Psychometrika*, *73*, 729–751.
40
41
42
43 Wilderjans, T. F., Ceulemans, E., Van Mechelen, I., & van den Berg, R. A. (2011).
44
45 Simultaneous analysis of coupled data matrices subject to different amounts of
46
47 noise. *British Journal of Mathematical and Statistical Psychology*, *64*, 277–290.
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

TABLE 1.

Hypothetical two-way two-mode coupled model matrices \mathbf{M}^1 , \mathbf{M}^2 , and \mathbf{M}^3

	\mathbf{M}^1				\mathbf{M}^2					\mathbf{M}^3					
	C_1^1	C_2^1	C_3^1	C_4^1	C_1^2	C_2^2	C_3^2	C_4^2	C_5^2	C_1^3	C_2^3	C_3^3	C_4^3	C_5^3	C_6^3
R_1	1	0	0	0	0	0	1	1	0	0	1	0	0	1	1
R_2	1	0	0	0	0	0	1	1	0	0	1	0	0	1	1
R_3	0	1	0	1	0	1	0	0	0	1	1	0	0	0	0
R_4	1	1	0	1	0	1	1	1	0	1	1	0	0	1	1
R_5	1	1	1	0	1	1	1	1	1	0	1	1	0	1	1
R_6	0	1	1	0	1	1	0	1	1	0	1	1	0	0	0

TABLE 2.

Concatenated HICLAS model with tree bundles for the coupled model matrices \mathbf{M}^1 , \mathbf{M}^2 , and \mathbf{M}^3 in

Table 1

	A			B¹			B²			B³					
	Bu ₁	Bu ₂	Bu ₃	Bu ₁	Bu ₂	Bu ₃	Bu ₁	Bu ₂	Bu ₃	Bu ₁	Bu ₂	Bu ₃			
R ₁	0	1	0	C ₁ ¹	0	1	0	C ₁ ²	0	0	1	C ₁ ³	1	0	0
R ₂	0	1	0	C ₂ ¹	1	0	1	C ₂ ²	1	0	1	C ₂ ³	1	1	1
R ₃	1	0	0	C ₃ ¹	0	0	1	C ₃ ²	0	1	0	C ₃ ³	0	0	1
R ₄	1	1	0	C ₄ ¹	1	0	0	C ₄ ²	0	1	1	C ₄ ³	0	0	0
R ₅	0	1	1					C ₅ ²	0	0	1	C ₅ ³	0	1	0
R ₆	0	0	1									C ₆ ³	0	1	0

TABLE 3.

The sizes of $\mathbf{D}^1 (I \times J_1)$ and $\mathbf{D}^2 (I \times J_2)$ that are obtained by an orthogonal combination of the factors
Size of the common mode, I , and Array/total ratio, r

I	r	size of $\mathbf{D}^1 (I \times J_1)$	size of $\mathbf{D}^2 (I \times J_2)$
50	.10	(50 × 40)	(50 × 360)
50	.30	(50 × 120)	(50 × 280)
50	.50	(50 × 200)	(50 × 200)
50	.70	(50 × 280)	(50 × 120)
50	.90	(50 × 360)	(50 × 40)
100	.10	(100 × 20)	(100 × 180)
100	.30	(100 × 60)	(100 × 140)
100	.50	(100 × 100)	(100 × 100)
100	.70	(100 × 140)	(100 × 60)
100	.90	(100 × 180)	(100 × 20)
200	.10	(200 × 10)	(200 × 90)
200	.30	(200 × 30)	(200 × 70)
200	.50	(200 × 50)	(200 × 50)
200	.70	(200 × 70)	(200 × 30)
200	.90	(200 × 90)	(200 × 10)

TABLE 4.

The values of ε_1 and ε_2 resulting from an orthogonal combination of the factors Total amount of noise, ε^{tot} , and Relative amount of noise, ε^{rel}

ε^{tot}	ε^{rel}	ε_1	ε_2
.20	0	0	.20
.20	.25	.05	.15
.20	.50	.10	.10
.20	.75	.15	.05
.20	1	.20	0
.40	0	0	.40
.40	.25	.10	.30
.40	.50	.20	.20
.40	.75	.30	.10
.40	1	.40	0

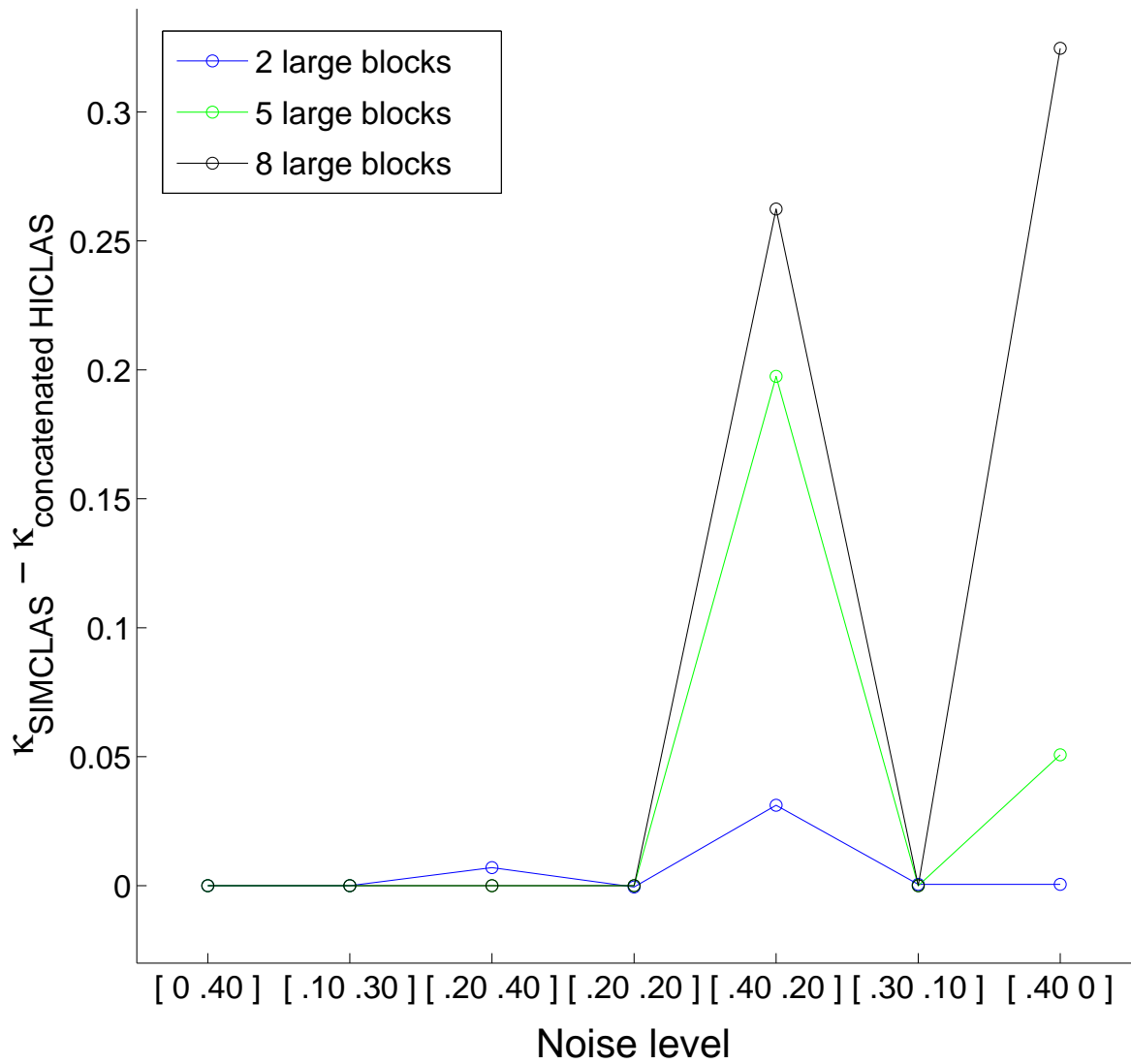


FIGURE 1.

Mean difference in κ between *SIMCLAS* and *concatenated HICLAS* as a function of the Noise level and the Dominance of the large data blocks