

A New Algorithm for Generalized Least Squares Factor Analysis with a Majorization Technique

Kohei Adachi

Graduate School of Human Sciences, Osaka University, Osaka, Japan
Email: adachi@hus.osaka-u.ac.jp

Received 18 January 2015; accepted 22 April 2015; published 27 April 2015

Copyright © 2015 by author and Scientific Research Publishing Inc.
This work is licensed under the Creative Commons Attribution International License (CC BY).
<http://creativecommons.org/licenses/by/4.0/>



Open Access

Abstract

Factor analysis (FA) is a time-honored multivariate analysis procedure for exploring the factors underlying observed variables. In this paper, we propose a new algorithm for the generalized least squares (GLS) estimation in FA. In the algorithm, a majorization step and diagonal steps are alternately iterated until convergence is reached, where Kiers and ten Berge's (1992) majorization technique is used for the former step, and the latter ones are formulated as minimizing simple quadratic functions of diagonal matrices. This procedure is named a majorizing-diagonal (MD) algorithm. In contrast to the existing gradient approaches, differential calculus is not used and only elementary matrix computations are required in the MD algorithm. A simulation study shows that the proposed MD algorithm recovers parameters better than the existing algorithms.

Keywords

Exploratory Factor Analysis, Generalized Least Squares Estimation, Matrix Computations, Majorization

1. Introduction

Using \mathbf{y} for a $p \times 1$ observation vector whose expectation $E(\mathbf{y})$ equals the $p \times 1$ zero vector $\mathbf{0}_p$, the factor analysis (FA) model is expressed as

$$\mathbf{y} = \mathbf{\Lambda}\mathbf{x} + \mathbf{e} \quad (1)$$

with $\mathbf{\Lambda} = (\lambda_{ij})$ a p -variables \times m -factors loading matrix, \mathbf{x} an $m \times 1$ latent factor score vector, \mathbf{e} a $p \times 1$ error vector, and $p > m$. The expectations for \mathbf{x} and \mathbf{e} are assumed to satisfy

$$E(\mathbf{x}) = \mathbf{0}_m, \quad E(\mathbf{e}) = \mathbf{0}_p, \quad E(\mathbf{x}\mathbf{e}') = \mathbf{O}, \quad E(\mathbf{x}\mathbf{x}') = \mathbf{I}_m, \quad E(\mathbf{e}\mathbf{e}') = \mathbf{\Psi} \quad (2)$$

Here, \mathbf{O} is the $m \times p$ matrix of zeros, \mathbf{I}_m is the $m \times m$ identity matrix, and $\mathbf{\Psi}$ is the $p \times p$ diagonal matrix whose diagonal elements are called unique variances. The FA model (1) with the assumptions in (2) imply that the covariance matrix $\mathbf{\Sigma} = E(\mathbf{y}\mathbf{y}')$ is modeled as

$$\mathbf{\Sigma} = \mathbf{\Lambda}\mathbf{\Lambda}' + \mathbf{\Psi} \quad (3)$$

[1] [2]. A main purpose of FA is to estimate the parameter matrices $\mathbf{\Lambda}$ and $\mathbf{\Psi}$ from the inter-variable sample covariance matrix \mathbf{S} ($p \times p$) corresponding to (3). Some authors classify FA as exploratory (EFA) or confirmatory (CFA) [2], where $\mathbf{\Lambda}$ is unconstrained in EFA, while some elements of $\mathbf{\Lambda}$ are constrained in CFA. In this paper, we refer to EFA simply as FA.

Three major approaches for the parameter estimation are least squares (LS), generalized least squares (GLS), and maximum likelihood (ML) procedures [3]. They differ in the definition of the loss function $f(\Theta)$ to be minimized over $\Theta = [\mathbf{\Lambda}, \mathbf{\Psi}]$. The functions for the LS and GLS estimation procedures are defined as $f(\Theta) = \text{tr}(\mathbf{S} - \mathbf{\Sigma})^2$ and

$$f(\Theta) = \text{tr}\{(\mathbf{S} - \mathbf{\Sigma})\mathbf{S}^{-1}\}^2 = \text{tr}\{[\mathbf{S} - (\mathbf{\Lambda}\mathbf{\Lambda}' + \mathbf{\Psi})]\mathbf{S}^{-1}\}^2 \quad (4)$$

respectively, while $f(\Theta)$ is defined as the negative of the log-likelihood derived under the normality assumption for \mathbf{x} and \mathbf{e} in the ML estimation [3] [4].

In all estimation procedures, iterative algorithms are needed for minimizing loss function $f(\Theta)$. They can be roughly classified into gradient and inequality-based algorithms. Here, the gradient ones refer to the algorithms using Newton and related methods [5], in which the partial differentiation of $f(\Theta)$ with respect to Θ is used for updating it. On the other hand, the term ‘‘inequality-based algorithms’’ is not a popular one. We use the term for the algorithms, in which differentiation is not used and the inequality $f(\Theta) \geq f(\Theta_{\text{new}})$ underlies that which guarantees the weakly monotone decrease in the loss function value with updating Θ to Θ_{new} . Similar dichotomization of minimization methodology is also found in [6].

For all of the LS, GLS, and ML estimation, gradient algorithms have been developed: those with the Fletcher-Powell and Newton-Raphson methods have been proposed for the ML estimation [7] [8], while the algorithms using the Newton-Raphson and Gauss-Newton methods have been developed for GLS [9] [10] with the gradient algorithms for GLS also used for LS. On the other hand, inequality-based algorithms have been developed for the LS and ML estimation excluding GLS. Such an algorithm for LS is MINRES [11] in which Θ is partitioned into the subsets of parameters with $\Theta = \{\theta_1, \dots, \theta_q\}$ and the minimization of $f(\Theta)$ over each subset θ_i for $i = 1, \dots, q$ is iterated. The inequality-based one for the ML estimation is the EM algorithm for FA [12] in which $f(\Theta)$ decreases monotonically with the alternate iteration of so-called E- and M-steps [13]. A feature of MINRES and the EM algorithm is that only simple matrix computations such as the inversion of matrices are required and their computer-programs are easily formed. In contrast, the gradient algorithms require more complicated computations such as obtaining or numerically approximating the second derivatives of $f(\Theta)$.

As found in the above discussion, an inequality-based algorithm has not been developed for the GLS estimation in which (4) is minimized over Θ . To propose it is the purpose of this paper. The algorithm to be proposed is also computationally simple as in the existing inequality-based ones: only elementary matrix computations are required such as the inversion and singular value decomposition (SVD) of matrices. A feature of the proposed algorithm to be addressed is using majorization in one of steps. The majorization generally refers to a class of the techniques in which a *majorizing function* $h(\Phi, \Xi)$ is utilized for minimizing a function $g(\Phi)$ over Φ . Here, $h(\Phi, \Xi)$ satisfies $g(\Phi) = h(\Phi, \Phi) \geq h(\Phi^*, \Phi) \geq g(\Phi^*)$, with Φ^* being the minimizer of the majorizing function $h(\Phi^*, \Phi)$ for its latter argument matrix Φ kept fixed [14]. It shows that $g(\Phi)$ decreases with the update of Φ into Φ^* . As described in the next section, the step with a majorization technique and the steps for minimizing the functions of diagonal matrices form the algorithm to be presented. It is thus called majorizing-diagonal (MD) algorithm in this paper.

The MD algorithm is not the first one with majorization in FA. Indeed, the above EM algorithm [12] can be regarded as a majorization procedure with its majorizing function being the full log likelihood derived by supposing that latent factor scores in \mathbf{x} were observed. [15] has also proposed an FA algorithm with a majorization technique. However, in that algorithm, the estimation of a new type [16] [17] is considered, which are different from the LS, GLS, and ML estimation treated as the major procedures in this paper: [15] is beyond the

scope of this paper.

The remaining parts of this paper are organized as follows: the MD algorithm is detailed in the next section, and it is illustrated with a real data set in Section 3. A simulation study for assessing the algorithm is reported in Section 4, which is followed by discussions.

2. Proposed Algorithm

We propose the MD algorithm for minimizing the GLS loss function (4) over the loadings in $\mathbf{\Lambda}$ ($p \times m$) and the unique variances in the diagonal matrix $\mathbf{\Psi}$ ($p \times p$). Here, it is supposed that the sample covariance matrix \mathbf{S} is positive-definite and $\mathbf{\Lambda}$ is of full-column rank, *i.e.*, its rank is m with $p > m$. This supposition and the covariance matrix being modeled as (3) imply that, without loss of generality, we can reparameterize $\mathbf{\Lambda}$ as

$$\mathbf{\Lambda} = \mathbf{L}\mathbf{\Delta}^{1/2} \quad (5)$$

where \mathbf{L} is a $p \times m$ matrix satisfying

$$\mathbf{L}'\mathbf{L} = \mathbf{I}_p \quad (6)$$

and $\mathbf{\Delta}$ is an $m \times m$ positive-definite diagonal matrix. By substituting (5) into the GLS loss function (4), it is rewritten as

$$\begin{aligned} f(\mathbf{L}, \mathbf{\Delta}, \mathbf{\Psi}) &= \text{tr} \left\{ \left[\mathbf{S} - (\mathbf{L}\mathbf{\Delta}\mathbf{L}' + \mathbf{\Psi}) \right] \mathbf{S}^{-1} \right\}^2 = \text{tr} \left\{ \mathbf{I}_p - (\mathbf{L}\mathbf{\Delta}\mathbf{L}' + \mathbf{\Psi}) \mathbf{S}^{-1} \right\}^2 \\ &= p - 2\text{tr} \mathbf{S}^{-1} (\mathbf{L}\mathbf{\Delta}\mathbf{L}' + \mathbf{\Psi}) + \text{tr} \mathbf{S}^{-2} (\mathbf{L}\mathbf{\Delta}\mathbf{L}' + \mathbf{\Psi})^2 \\ &= p - 2\text{tr} \mathbf{S}^{-1} \mathbf{L}\mathbf{\Delta}\mathbf{L}' - 2\text{tr} \mathbf{S}^{-1} \mathbf{\Psi} + \text{tr} \mathbf{S}^{-2} \mathbf{L}\mathbf{\Delta}^2 \mathbf{L}' + 2\text{tr} \mathbf{\Psi} \mathbf{S}^{-2} \mathbf{L}\mathbf{\Delta}\mathbf{L}' + \text{tr} \mathbf{S}^{-2} \mathbf{\Psi}^2 \end{aligned} \quad (7)$$

This function is minimized over \mathbf{L} , $\mathbf{\Delta}$, and $\mathbf{\Psi}$ subject to (6) and the latter two matrices being diagonal ones, by alternately iterating the majorizing and diagonal steps described in the next subsections.

2.1. Majorization Step

Let us consider minimizing (7) over \mathbf{L} subject to (6) while $\mathbf{\Delta}$ and $\mathbf{\Psi}$ are kept fixed. Summarizing the parts irrelevant to \mathbf{L} in (7) into $c = p - 2\text{tr} \mathbf{S}^{-1} \mathbf{\Psi} + \text{tr} \mathbf{S}^{-2} \mathbf{\Psi}^2$, the loss function (7) is rewritten as

$$f(\mathbf{L} | \mathbf{\Delta}, \mathbf{\Psi}) = 2\text{tr} (\mathbf{\Psi} \mathbf{S}^{-2} - \mathbf{S}^{-1}) \mathbf{L}\mathbf{\Delta}\mathbf{L}' + \text{tr} \mathbf{S}^{-2} \mathbf{L}\mathbf{\Delta}^2 \mathbf{L}' + c. \quad (8)$$

Though the optimal \mathbf{L} that minimizes (8) under (6) is not given explicitly, the solution can be obtained using Kiers and ten Berge's [18] majorization technique, whose earlier version is also found in [19]. This technique purposes to minimize a function expressed as the form $\phi(\mathbf{L}) = \text{tr} \mathbf{A}\mathbf{L} + \sum_{j=1}^q \text{tr} \mathbf{B}_j \mathbf{L} \mathbf{C}_j \mathbf{L}' + \text{constant}$. Comparing this with (8), we can find (8) to be a special case of the above $\phi(\mathbf{L})$ with \mathbf{A} being the zero matrix, $\mathbf{B}_1 = 2(\mathbf{\Psi} \mathbf{S}^{-2} - \mathbf{S}^{-1})$, $\mathbf{B}_2 = \mathbf{S}^{-2}$, $\mathbf{C}_1 = \mathbf{\Delta}$, $\mathbf{C}_2 = \mathbf{\Delta}^2$, and $q = 2$. Therefore, the update formula in [18] (pp. 374-375) can be straightforwardly used for (8).

According to the formula, the update of \mathbf{L} by

$$\mathbf{L} = \mathbf{P}\mathbf{Q}' \quad (9)$$

decreases the value of (8) with $f(\mathbf{L} | \mathbf{\Delta}, \mathbf{\Psi}) \leq f(\mathbf{L}_{\text{OLD}} | \mathbf{\Delta}, \mathbf{\Psi})$. Here, \mathbf{L}_{OLD} stands for the matrix \mathbf{L} before the update; \mathbf{P} and \mathbf{Q} are the column-orthonormal matrices that are obtained from the SVD defined as

$$2\{(\lambda_a - \lambda_b) \mathbf{I}_m - \mathbf{S}^{-2} \mathbf{\Psi} - \mathbf{\Psi} \mathbf{S}^{-2} + 2\mathbf{S}^{-1}\} \mathbf{L}_{\text{OLD}} \mathbf{\Delta} + 2(\lambda_c - \mathbf{I}_m - \mathbf{S}^{-2}) \mathbf{L}_{\text{OLD}} \mathbf{\Delta}^2 = \mathbf{P}\mathbf{\Theta}\mathbf{Q}' \quad (10)$$

with $\mathbf{\Theta}$ the diagonal matrix including the singular values of the matrix in the left-hand side, and λ_a , λ_b , and λ_c the largest eigenvalues of $\mathbf{\Psi} \mathbf{S}^{-2} + \mathbf{S}^{-2} \mathbf{\Psi}$, $2\mathbf{S}^{-1}$, and \mathbf{S}^{-2} , respectively.

2.2. Diagonal Steps

In this section, we describe updating each of diagonal matrices $\mathbf{\Delta}$ and $\mathbf{\Psi}$. First, let us consider minimizing the loss function (7) over $\mathbf{\Delta}$ with keeping \mathbf{L} and $\mathbf{\Psi}$ fixed. Since the terms relevant to $\mathbf{\Delta}$ in the loss function (7) are the same as those relevant to \mathbf{L} , the expression (8) into which (7) is rewritten is to be noted again. By taking account of the fact that $\mathbf{\Delta}$ is a diagonal matrix, (8) can be rewritten as

$$f(\Delta|\mathbf{L}, \Psi) = \text{tr}(\mathbf{L}'\mathbf{S}^{-2}\mathbf{L})\Delta^2 - 2\text{tr}\mathbf{L}'(\mathbf{S}^{-1} - \Psi\mathbf{S}^{-2})\mathbf{L}\Delta + c = \text{tr}\mathbf{D}_1\Delta^2 - 2\text{tr}\mathbf{D}_2\Delta + c \quad (11)$$

Here, $\mathbf{D}_1 = \text{diag}(\mathbf{L}'\mathbf{S}^{-2}\mathbf{L})$ and $\mathbf{D}_2 = \text{diag}(\mathbf{L}'(\mathbf{S}^{-1} - \Psi\mathbf{S}^{-2})\mathbf{L})$ with $\text{diag}(\bullet)$ denoting the diagonal matrix whose diagonal elements are those of the parenthesized matrix. Further, we can rewrite (11) as $f(\Delta|\mathbf{L}, \Psi) = \mathbf{D}_1\|\Delta - \mathbf{D}_1^{-1}\mathbf{D}_2\|^2 + \text{tr}\mathbf{D}_1^{-2}\mathbf{D}_2^2 + c$ with $\|\bullet\|$ denoting the Frobenius norm. It shows that the function (11) is minimized for

$$\Delta = \mathbf{D}_1^{-1}\mathbf{D}_2 = \text{diag}(\mathbf{L}'\mathbf{S}^{-2}\mathbf{L})^{-1} \text{diag}(\mathbf{L}'(\mathbf{S}^{-1} - \Psi\mathbf{S}^{-2})\mathbf{L}) \quad (12)$$

for fixed \mathbf{L} and Ψ .

Next, we consider minimizing (7) over Ψ with \mathbf{L} and Δ fixed. Summarizing the parts irrelevant to Ψ in (7) into $c^* = p - 2\text{tr}\mathbf{S}^{-1}\mathbf{L}\Delta\mathbf{L}' + \text{tr}\mathbf{S}^{-2}\mathbf{L}\Delta^2\mathbf{L}'$ and using the fact of Ψ being a diagonal matrix, the loss function (7) can be rewritten as

$$\begin{aligned} f(\Psi|\mathbf{L}, \Delta) &= \text{tr}\mathbf{S}^{-2}\Psi^2 - 2\text{tr}(\mathbf{S}^{-1} - \mathbf{S}^{-2}\mathbf{L}\Delta\mathbf{L}')\Psi + c^* \\ &= \text{tr}\mathbf{D}_3\Psi^2 - 2\text{tr}\mathbf{D}_4\Psi + c^* \\ &= \mathbf{D}_3\|\Psi - \mathbf{D}_3^{-1}\mathbf{D}_4\|^2 + \text{tr}\mathbf{D}_3^{-2}\mathbf{D}_4^2 \end{aligned} \quad (13)$$

with $\mathbf{D}_3 = \text{diag}(\mathbf{S}^{-2})$ and $\mathbf{D}_4 = \text{diag}(\mathbf{S}^{-1} - \mathbf{S}^{-2}\mathbf{L}\Delta\mathbf{L}')$. We can find that (13) is minimized for

$$\Psi = \mathbf{D}_3^{-1}\mathbf{D}_4 = \text{diag}(\mathbf{S}^{-2})^{-1} \text{diag}(\mathbf{S}^{-1} - \mathbf{S}^{-2}\mathbf{L}\Delta\mathbf{L}') \quad (14)$$

for fixed \mathbf{L} and Δ .

2.3. Whole Algorithm

The results in the last two subsections show that the proposed MD algorithm can be listed as follows:

- Step 1. Initialize \mathbf{L} , Δ , and Ψ .
- Step 2. Update \mathbf{L} with (9) l times.
- Step 3. Update Δ with (12).
- Step 4. Update Ψ with (14).
- Step 5. Finish with \mathbf{A} set to (5) if convergence is reached; otherwise, return to Step 2.

It should be noted in Step 2 that the update of \mathbf{L} by (9) does not minimize (7) but only decreases its value, which implies that that update can be replicated (l times) for further decreasing the value of (7). In this paper, we set $l = 5$.

In Step 1, the initialization is performed using the principal component analysis of sample covariance matrix \mathbf{S} . That is, the initial \mathbf{L} and Δ are given by \mathbf{V} and $\mathbf{\Omega}^{1/2}$, respectively, with $\mathbf{\Omega}$ the $m \times m$ diagonal matrix whose diagonal elements are the m largest eigenvalues of \mathbf{S} , and the columns of \mathbf{V} ($p \times m$) being the eigenvectors corresponding to $\mathbf{\Omega}$. The initial Ψ is set to $\text{diag}(\mathbf{S} - \mathbf{V}\mathbf{\Omega}\mathbf{V}')$.

In Step 5, we define the convergence as the decrease in the value of (7) or (4) from the previous round being less than $p^2 \times 0.1^{10}$.

3. Illustration

In this section, we illustrate the performance of the MD algorithm with a 190-person \times 25-item data matrix, which was collected by the author and publicly available at <http://bm.osaka-u.ac.jp/data/big5/>. This data set contains the self-ratings of the persons (university students) for to what extents they are characterized by the personalities described by the 25 items. According to a theory in personality psychology [20], the items can be classified into the five groups shown in the first column of **Table 1**. The 25×25 matrix of the correlation coefficients among those items was obtained from the data set.

We carried out the MD algorithm for the correlation matrix with the number of factors m set to five. **Figure 1** shows the change in the value of loss function (4) until the steps in Section 2.3 were iterated ten times and the change after the tenth iteration. There, we can find that the function value decreased monotonically with iteration, which was finally reached to convergence at the 542 nd iteration.

Table 1. Loadings and unique variances Ψ_{1p} for personality rating data.

Group	Item	Loading matrix					Ψ_{1p}
Neurotic	Worry	0.72	-0.14	-0.12	0.22	0.13	0.32
	Sensitive	0.74	-0.06	-0.04	0.01	-0.08	0.38
	Pessimistic	0.67	-0.29	-0.13	-0.03	0.27	0.35
	Unrest	0.60	0.06	-0.15	-0.08	-0.34	0.33
	Careful	0.62	-0.09	-0.16	0.10	0.32	0.32
Extrovert	Sociable	-0.16	0.81	0.09	0.09	0.12	0.23
	Talkative	0.04	0.82	-0.05	-0.02	-0.09	0.25
	Voluntary	-0.10	0.73	0.17	0.10	0.13	0.29
	Cheerful	-0.21	0.80	0.08	0.16	0.00	0.23
	Showy	0.08	0.66	0.23	-0.04	-0.05	0.38
Open	Creative	-0.09	0.03	0.85	-0.04	-0.07	0.23
	Adventurous	-0.22	0.22	0.67	-0.02	-0.20	0.36
	Progressive	-0.19	0.24	0.64	0.08	0.10	0.40
	Flexible	-0.28	0.31	0.37	0.20	0.03	0.57
	Imaginative	0.19	0.09	0.47	0.11	-0.34	0.44
Agreeable	Mild	-0.15	-0.12	0.15	0.63	0.02	0.38
	Tenderhearted	0.09	0.20	0.13	0.62	-0.01	0.39
	Altruistic	0.13	0.01	0.00	0.71	0.15	0.39
	Cooperative	0.00	0.27	-0.14	0.66	0.08	0.25
	Sympathetic	0.11	0.12	-0.04	0.73	0.21	0.32
Conscious	Deliberate	0.14	-0.04	0.09	0.22	0.60	0.37
	Reliable	-0.13	0.32	0.08	0.26	0.57	0.36
	Diligent	-0.04	0.05	-0.07	0.15	0.77	0.29
	Systematic	0.09	-0.02	-0.09	0.01	0.71	0.40
	Methodical	0.21	0.02	-0.19	0.04	0.74	0.31

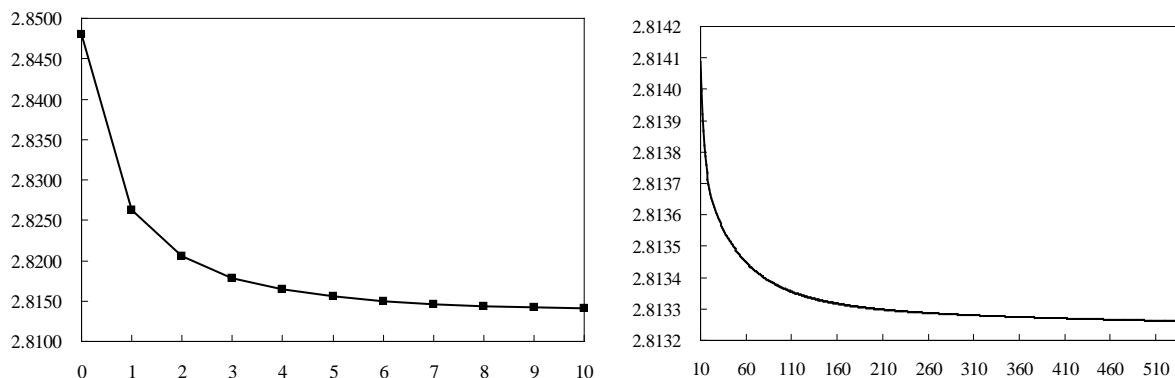


Figure 1. Change in the GLS loss function value as a function of the number of iteration.

As the resulting loading matrix has rotational freedom, that is, the Λ post-multiplied by arbitrary orthonormal matrix satisfies (1) and (2), the loading matrix was rotated by the varimax method [21]. The solution is pre-

sented in **Table 1**. There, bold font is used for the loadings whose absolute values are greater than 0.35. They show that the 25 items are clearly classified into the five groups as predicted by the theory in personality psychology [20], which demonstrates that the MD algorithm provided the reasonable solution.

4. Simulation Study

A simulation study was performed in order to assess how well parameter matrices are recovered by the proposed MD algorithm and compare it with the existing algorithms for the GLS estimation in the goodness of the recovery. We first describe the procedure for synthesizing the data to be analyzed, which is followed by results.

An n -observations \times p -variables data matrix \mathbf{Y} was synthesized according to the matrix versions of the FA model (1) and the assumptions in (2):

$$\mathbf{Y} = \mathbf{X}\mathbf{\Lambda}' + \mathbf{E} = \mathbf{X}\mathbf{\Lambda}' + \mathbf{U}\mathbf{\Psi}^{1/2} = [\mathbf{X}, \mathbf{U}][\mathbf{L}, \mathbf{\Psi}^{1/2}]' = \mathbf{Z}\mathbf{\Gamma}' \quad (15)$$

$$\mathbf{1}'_{p+m} \mathbf{Z} = \mathbf{0}'_{p+m} \quad \text{and} \quad \mathbf{Z}'\mathbf{Z} = \mathbf{I}_{p+m} \quad (16)$$

Here, $\mathbf{1}_{p+m}$ denotes the $(p+m) \times 1$ vector of ones, $\mathbf{Z} = [\mathbf{X}, \mathbf{U}]$ is an $n \times (m+p)$ block matrix whose right $n \times p$ block \mathbf{U} is post-multiplied by $\mathbf{\Psi}^{1/2}$ to give the error matrix $\mathbf{E} = \mathbf{U}\mathbf{\Psi}^{1/2}$, and $\mathbf{\Gamma} = [\mathbf{\Lambda}, \mathbf{\Psi}^{1/2}]$ is an $(m+p) \times p$ block matrix including the loading matrix and the square roots of unique variances. It should be noticed that each row of \mathbf{Y} , \mathbf{X} , and \mathbf{E} corresponds to \mathbf{y}' , \mathbf{x}' , and \mathbf{e}' , respectively, whose transposed vectors appear in (1), and the five equations in (2) can be summarized into the two matrix expressions in (16). The data synthesis procedure follows the next steps:

Step 1. Draw m from $DU(1,5)$, p from $DU(4m,7m)$, and n from $DU(9p,14p)$, with $DU(I,J)$ denoting the discrete uniform distribution defined for the integers within the range $[I,J]$.

Step 2. Draw each loading in $\mathbf{\Lambda}$ from $U(-1,1)$ and each unique variance in $\mathbf{\Psi}$ from $U(0.1,0.7)$ with $U(\alpha,\beta)$ denoting the uniform distribution over the range $[\alpha,\beta]$.

Step 3. Draw each elements of \mathbf{Z} in (15) from $U(-1,1)$ which is followed by centering \mathbf{Z} and post-multiplying it by the matrix that allows the resulting \mathbf{Z} to satisfy (16).

Step 4. Form \mathbf{Y} with (15) and obtain the covariance matrix $\mathbf{S} = n^{-1}\mathbf{Y}'\mathbf{Y}$.

In Step 3 we have used a uniform distribution for \mathbf{Z} , rather than the normal distribution typically used for such a matrix, as a feature of the GLS estimation is that it does not need the normality assumption required in the ML estimation. We replicated the above steps to have 2000 sets of \mathbf{S} . For them, the MD and the existing algorithms were carried out, where the latter are the two gradient algorithms [9] [10], as described in Section 1. We refer to the ones in [9] and [10] as the Newton-Raphson (NR) and Gauss-Newton (GN) algorithms, respectively. In the NR one, we obtained the gradient vector in [9], Equation (32), by pre-multiplying the vector of first derivatives by the Moore Penrose inverse of the corresponding Hessian matrix. Also in the NR and GN algorithms, we used the same initialization and definition of convergence as in Section 2.3.

Let us express the true $[\mathbf{\Lambda}, \mathbf{\Psi}]$ simply as $[\mathbf{\Lambda}, \mathbf{\Psi}]$ and use $[\mathbf{\Lambda}_\#, \mathbf{\Psi}_\#]$ for the solution given by the NR, GN, or MD algorithm. For assessing the recovery of the loading matrix, the averaged absolute difference (AAD) of the elements in $\mathbf{\Lambda}$ to the corresponding estimates, *i.e.*,

$$\text{AAD}_\Lambda(\mathbf{\Lambda}, \mathbf{\Lambda}_\#) = (pm)^{-1} \|\mathbf{\Lambda} - \mathbf{\Lambda}_\#\|_{l_1}, \quad (17)$$

can be used with $\|\cdot\|_{l_1}$ denoting the l_1 norm. Here, it should be noted that $\mathbf{\Lambda}_\#$ has rotational freedom and must be rotated so that the resulting $\mathbf{\Lambda}_\#$ is optimally matched to $\mathbf{\Lambda}$. Such a rotated $\mathbf{\Lambda}_\#$ can be obtained by the orthogonal Procrustes method [22] with $\mathbf{\Lambda}$ a target matrix. The loading matrix $\mathbf{\Lambda}_\#$ in (17) thus stands for the one rotated by the Procrustes method. The recovery of unique variances can also be assessed with the AAD index $\text{AAD}_\Psi(\mathbf{\Psi}, \mathbf{\Psi}_\#) = p^{-1} \|\mathbf{\Psi} - \mathbf{\Psi}_\#\|_{l_1}$, where the unique variances are uniquely determined, thus the additional procedure as for $\mathbf{\Lambda}_\#$ is unnecessary. Smaller values of those AAD indices stand for better recovery.

The statistics of AAD values over 2000 data sets are presented in **Table 2**. There, the averages show that the recovery by the MD algorithm is the best and that for the NR one is the worst. It should be noted that the 50 and 75 percentiles for the NR algorithm are zero, while the maximum and 99 percentile are very large. That is, the recovery by the NR algorithm was perfect for more than 75 percent of the 2000 data sets, but for a few percent of them, recovery was considerably bad, which increased the averages for the NR one. In contrast, the maximum AAD of loadings and unique variances for the MD algorithm are 0.0041 and 0.0013, respectively, which are

Table 2. Statistics for the differences between the true parameter values and their estimated counterparts.

Parameter	Loadings			Unique variances		
	NR	GN	MD	NR	GN	MD
Average	0.0026	0.0020	0.0005	0.0030	0.0007	0.0000
50 percentile	0.0000	0.0013	0.0005	0.0000	0.0001	0.0000
75 percentile	0.0000	0.0031	0.0007	0.0000	0.0005	0.0000
95 percentile	0.0050	0.0061	0.0011	0.0103	0.0035	0.0000
99 percentile	0.0821	0.0095	0.0014	0.0857	0.0105	0.0001
Maximum	0.2199	0.0203	0.0041	0.1406	0.0323	0.0013

small enough to be ignored. That is, the proposed MD algorithm well recovered the true parameter values for all of the 2000 data sets. We can thus conclude that the MD algorithm is superior to the existing ones in the goodness of recovery.

5. Discussion

We proposed the majorizing-diagonal (MD) algorithm for the GLS estimation in FA. In the algorithm, the loading matrix is reparameterized as the product of a column-orthonormal matrix and a diagonal one, and the former one is updated with Kiers and ten Berge's [18] majorization technique, while the latter diagonal matrix and another diagonal one including unique variances are updated so that their quadratic functions are minimized. The iteration of those updates decreases monotonically the GLS loss function. The simulation study demonstrated the exact recovery of loadings and unique variances by the MD algorithm and its superiority to the existing gradient algorithms in the recovery.

One of the tasks remaining for the MD algorithm is to study its mathematical properties as have been done for the algorithms in the other estimation procedures. For example, it has been found that the EM algorithm for the ML estimation [12] can never give an improper solution under a certain condition [23], where the improper solution refers to the one including a negative unique variance. Whether such special features are possessed by the MD algorithm is considered to be found by studying the properties of the matrix update formulas in the algorithm.

References

- [1] Harman, H.H. (1976) *Modern Factor Analysis*. 3rd Edition, The University of Chicago Press, Chicago.
- [2] Mulaik, S.A. (2010) *Foundations of Factor Analysis*. 2nd Edition, CRC Press, Boca Raton.
- [3] Yanai, H. and Ichikawa, M. (2007) *Factor Analysis*. In: Rao, C.R. and Sinharay, S., Eds., *Handbook of Statistics*, Vol. 26: *Psychometrics*, Elsevier, Amsterdam, 257-296.
- [4] Anderson, T.W. and Rubin, H. (1956) *Statistical Inference in Factor Analysis*. In: Neyman, J., Ed., *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability*, Vol. 5, University of California Press, Berkeley, 111-150.
- [5] Lange, K. (2010) *Numerical Analysis for Statisticians*. 2nd Edition, Springer, New York.
- [6] ten Berge, J.M.F. (1993) *Least Squares Optimization in Multivariate Analysis*. DSWO Press, Leiden.
- [7] Jöreskog, K.G. (1967) Some Contributions to Maximum Likelihood Factor Analysis. *Psychometrika*, **32**, 443-482. <http://dx.doi.org/10.1007/BF02289658>
- [8] Jennrich, R.I. and Robinson, S.M. (1969) A Newton-Raphson Algorithm for Maximum Likelihood Factor Analysis. *Psychometrika*, **34**, 111-123. <http://dx.doi.org/10.1007/BF02290176>
- [9] Jöreskog, K.G. and Goldberger, A.S. (1972) Factor Analysis by Generalized Least Squares. *Psychometrika*, **37**, 243-250. <http://dx.doi.org/10.1007/BF02306782>.
- [10] Lee, S.Y. (1978) The Gauss-Newton Algorithm for the Weighted Least Squares Factor Analysis. *Journal of the Royal Statistical Society: Series D (The Statistician)*, **27**, 103-114. <http://dx.doi.org/10.2307/2987906>
- [11] Harman, H.H. and Jones, W.H. (1966) Factor Analysis by Minimizing Residuals (Minres). *Psychometrika*, **31**, 351-369. <http://dx.doi.org/10.1007/BF02289468>

-
- [12] Rubin, D.B. and Thayer, D.T. (1982) EM Algorithms for ML Factor Analysis. *Psychometrika*, **47**, 69-76. <http://dx.doi.org/10.1007/BF02293851>
- [13] Dempster, A.P., Laird, N.M. and Rubin, D.B. (1977) Maximum Likelihood from Incomplete Data via the EM Algorithm. *Journal of the Royal Statistical Society, Series B*, **39**, 1-38.
- [14] Groenen, P.J.F. (1993) The Majorization Approach to Multidimensional Scaling: Some Problems and Extensions. DSWO Press, Leiden.
- [15] Unkel, S. and Trendafilov, N.T. (2010) A Majorization Algorithm for Simultaneous Parameter Estimation in Robust Exploratory Factor Analysis. *Computational Statistics and Data Analysis*, **54**, 3348-3358. <http://dx.doi.org/10.1016/j.csda.2010.02.003>
- [16] Unkel, S. and Trendafilov, N.T. (2010) Simultaneous Parameter Estimation in Exploratory Factor Analysis: An Expository Review. *International Statistical Review*, **78**, 363-382. <http://dx.doi.org/10.1111/j.1751-5823.2010.00120.x>
- [17] Adachi, K. (2012) Some Contributions to Data-Fitting Factor Analysis with Empirical Comparisons to Covariance-Fitting Factor Analysis. *Journal of the Japanese Society of Computational Statistics*, **25**, 25-38. http://dx.doi.org/10.5183/jjcs.1106001_197
- [18] Kiers, H.A.L. and ten Berge, J.M.F. (1992) Minimization of a Class of Matrix Trace Functions by Means of Refined Majorization. *Psychometrika*, **57**, 371-382. <http://dx.doi.org/10.1007/BF02295425>
- [19] Kiers, H.A.L. (1990) Majorization as a Tool for Optimizing a Class of Matrix Functions. *Psychometrika*, **55**, 417-428. <http://dx.doi.org/10.1007/BF02294758>
- [20] Costa, P.T. and McCrae, R.R. (1992) NEO PI-R Professional Manual: Revised NEO Personality Inventory (NEO PI-R) and NEO Five-Factor Inventory (NEO-FFI). Psychological Assessment Resources, Odessa, FL.
- [21] Kaiser, H.F. (1958) The Varimax Criterion for Analytic Rotation in Factor Analysis. *Psychometrika*, **23**, 187-200. <http://dx.doi.org/10.1007/BF02289233>
- [22] Gower, J.C. and Dijksterhuis, G.B. (2004) Procrustes Problems. Oxford University Press, Oxford. <http://dx.doi.org/10.1093/acprof:oso/9780198510581.001.0001>
- [23] Adachi, K. (2013) Factor Analysis with EM Algorithm Never Gives Improper Solutions When Sample Covariance and Initial Parameter Matrices Are Proper. *Psychometrika*, **78**, 380-394. <http://dx.doi.org/10.1007/s11336-012-9299-8>