

University of Groningen

Estimating the variance of estimator of the latent factor linear mixed model using supplemented expectation-maximization algorithm

Angraini, Yenni; Notodiputro, Khairil Anwar; Folmer, Henk; Saefuddin, Asep; Toharudin, Toni

Published in:
Symmetry

DOI:
[10.3390/sym13071286](https://doi.org/10.3390/sym13071286)

IMPORTANT NOTE: You are advised to consult the publisher's version (publisher's PDF) if you wish to cite from it. Please check the document version below.

Document Version
Publisher's PDF, also known as Version of record

Publication date:
2021

[Link to publication in University of Groningen/UMCG research database](#)

Citation for published version (APA):

Angraini, Y., Notodiputro, K. A., Folmer, H., Saefuddin, A., & Toharudin, T. (2021). Estimating the variance of estimator of the latent factor linear mixed model using supplemented expectation-maximization algorithm. *Symmetry*, 13(7), [1286]. <https://doi.org/10.3390/sym13071286>

Copyright

Other than for strictly personal use, it is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license (like Creative Commons).

The publication may also be distributed here under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license. More information can be found on the University of Groningen website: <https://www.rug.nl/library/open-access/self-archiving-pure/taverne-amendment>.

Take-down policy

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

Downloaded from the University of Groningen/UMCG research database (Pure): <http://www.rug.nl/research/portal>. For technical reasons the number of authors shown on this cover page is limited to 10 maximum.

Article

Estimating the Variance of Estimator of the Latent Factor Linear Mixed Model Using Supplemented Expectation-Maximization Algorithm

Yenni Angraini ^{1,2,*} , Khairil Anwar Notodiputro ^{1,*}, Henk Folmer ², Asep Saefuddin ¹ and Toni Toharudin ³ ¹ Department of Statistics, IPB University, Bogor 16680, Indonesia; asaefuddin@apps.ipb.ac.id² Faculty of Spatial Sciences, University of Groningen, 9747 Groningen, The Netherlands; h.folmer@rug.nl³ Department of Statistics, University of Padjadjaran, Bandung 16426, Indonesia; toni.toharudin@unpad.ac.id

* Correspondence: y_angraini@apps.ipb.ac.id (Y.A.); khairil@apps.ipb.ac.id (K.A.N.)

Abstract: This paper deals with symmetrical data that can be modelled based on Gaussian distribution, such as linear mixed models for longitudinal data. The latent factor linear mixed model (LFLMM) is a method generally used for analysing changes in high-dimensional longitudinal data. It is usual that the model estimates are based on the expectation-maximization (EM) algorithm, but unfortunately, the algorithm does not produce the standard errors of the regression coefficients, which then hampers testing procedures. To fill in the gap, the Supplemented EM (SEM) algorithm for the case of fixed variables is proposed in this paper. The computational aspects of the SEM algorithm have been investigated by means of simulation. We also calculate the variance matrix of beta using the second moment as a benchmark to compare with the asymptotic variance matrix of beta of SEM. Both the second moment and SEM produce symmetrical results, the variance estimates of beta are getting smaller when number of subjects in the simulation increases. In addition, the practical usefulness of this work was illustrated using real data on political attitudes and behaviour in Flanders-Belgium.

Keywords: latent factor linear mixed model (LFLMM); expectation-maximization (EM) algorithm; supplemented EM algorithm; longitudinal data analysis



Citation: Angraini, Y.; Notodiputro, K.A.; Folmer, H.; Saefuddin, A.; Toharudin, T. Estimating the Variance of Estimator of the Latent Factor Linear Mixed Model Using Supplemented Expectation-Maximization Algorithm. *Symmetry* **2021**, *13*, 1286. <https://doi.org/10.3390/sym13071286>

Academic Editors: Christophe Chesneau and José Carlos R. Alcantud

Received: 3 June 2021

Accepted: 15 July 2021

Published: 17 July 2021

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Copyright: © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

The latent factor multivariate linear mixed model (LFLMM) is a combination between the Factor Analysis (FA) and the Linear Mixed Model (LMM), as proposed by [1]. The model aims to analyze longitudinal data sets with large numbers of multivariate responses, i.e., high-dimensional longitudinal data. The authors proposed estimation of the LFLMM by means of the EM algorithm, which is a closed-form solution. They showed by way of simulation that EM estimation of the LFLMM provides accurate parameter estimates and is more efficient in terms of adding variables other than time variables in the model than alternatives like the structural equation model. As shown by [2,3], the combination of fixed and random effects and the interaction of covariates with time can be straightforwardly handled by the LFLMM estimated by EM.

The LFLMM assumes that the responses are continuous and that the number of latent variables is known [1]. Moreover, convergence of the EM algorithm is sometimes slow. The main disadvantage, however, is that the EM algorithm does not produce standard errors of the estimator of the regression coefficients because it does not calculate the derivatives of the likelihood function, which are often complicated and tedious to derive [4,5]. Thus, it is difficult to study the effects of different covariates or fixed variables for different latent factors simultaneously.

The general Supplemented EM algorithm was proposed by [6] to obtain the standard errors by calculating the complete information matrix as the base of the variance-covariance

matrix of the estimator. The Supplemented EM algorithm has been applied to various kinds of models, notably item response models [6–9]. However, its suitability and features in the case of application to the LFLMM have not been investigated yet. In this study, we extend the work of [1] by employing the Supplemented EM algorithm as a by-product of the EM estimator for the case of fixed variables. We used simulation studies to investigate the computational aspects of the Supplemented EM algorithm and used a real data example to illustrate the practical usefulness of this work.

The remainder of this study is organized as follows. In Section 2, we specify the LFLMM and summarize the EM algorithm to estimate it. Section 3 presents the Supplemented EM algorithm. Sections 4 and 5 present the results of the simulation and real data example. Conclusions follow in Section 6.

2. The LFLMM and the EM Algorithm

Following [1], an LFLMM can be composed of two parts. The first is the factor analysis model, which represents the relationships between the observed and latent variables. This is similar to the structural equation model, which explains the relationship of the latent variables and the measurement indicators carried out through factor analysis [10]. This part can be written as:

$$Y_{it} = \Lambda \eta_{it} + \epsilon_{it} \tag{1}$$

Specifically, for the i -th of N individuals, we observe $j = 1, \dots, J$ responses which characterize d latent factors $(\eta_{it} = (\eta_{it}^1, \dots, \eta_{it}^d), d < J)$ at time $t, t = 1, \dots, T_i$, where T_i is the number of time periods for subject i . Λ is the matrix of factor loadings and $\epsilon_{it} = (\epsilon_{it1}, \dots, \epsilon_{itJ})$ the vector of measurement errors for subject i at time t . It is assumed that $\epsilon_{itj} \sim N(0, \tau_j^2)$ and $\epsilon_{itj} \perp \epsilon_{ith}, j \neq h$. In matrix notation, Equation (1) reads:

$$Y_i = (I_{T_i} \otimes \Lambda) \eta_i + \epsilon_i \tag{2}$$

where

$$\begin{aligned} Y_i &= (y'_{i1}, \dots, y'_{iT_i})'_{[J \times T_i, 1]} \\ \eta_i &= (\eta'_{i1}, \dots, \eta'_{iT_i})'_{[d \times T_i, 1]} \\ \Lambda_{[J \times d]} &= \begin{pmatrix} \lambda'_1 \\ \vdots \\ \lambda'_J \end{pmatrix} \end{aligned}$$

The second part of the LFLMM is a multivariate linear mixed model containing the fixed and random effects for each latent variable (η_{it}). For individual $i, i = 1, 2, \dots, N$, at time $t, t = 1, 2, \dots, T_i$, and latent variable $l, l = 1, 2, \dots, d$, we thus have:

$$\eta_{it}^l = x_{it}^l \beta^l + z_{it}^l a_i^l + \epsilon_{it}^l \tag{3}$$

where x_{it}^l and z_{it}^l are the elements of design matrices of the p fixed variables and q random effects, respectively. β^l is an unknown coefficient, $a_i^l = (a_{i1}^l, \dots, a_{iq}^l)$ and $\epsilon_{it}^l = (\epsilon_{it1}^l, \dots, \epsilon_{itT_i}^l), l = 1, 2, \dots, d$, are the random effects and errors for subject i and factor l , respectively. The random effects are assumed to be normally distributed with mean 0 and variance-covariance matrix $V(a) = \Sigma_a$. It is assumed that Σ_a captures the changes among the latent variables [1]. For example, a positive covariance between the random effects for the latent variables 1 and 2 means that if for a given individual i the latent variable 1 increases over time, the latent variable 2 also increases for that individual. Note that in this setting, the covariates are included in the multivariate linear mixed model (MLMM) of Equation (3) but not in the factor analysis model of Equation (1).

In matrix notation, Equation (3) reads:

$$\boldsymbol{\eta}_i = \mathbf{X}_i\boldsymbol{\beta} + \mathbf{Z}_i\mathbf{a}_i + \boldsymbol{\varepsilon}_i \tag{4}$$

where

$$\begin{aligned} \mathbf{X}_i &= \begin{pmatrix} \mathbf{x}_{i1} \\ \vdots \\ \mathbf{x}_{iT_i} \end{pmatrix}_{[d \times T_i, p \times d]}, & \mathbf{x}_{it} &= \begin{pmatrix} x_{it}^1 & 0 & \dots & 0 \\ 0 & x_{it}^2 & \dots & 0 \\ \dots & \dots & \ddots & \vdots \\ 0 & 0 & \dots & x_{it}^d \end{pmatrix}_{[d, p \times d]} \\ \mathbf{Z}_i &= \begin{pmatrix} \mathbf{z}_{i1} \\ \vdots \\ \mathbf{z}_{iT_i} \end{pmatrix}_{[d \times T_i, q \times d]}, & \mathbf{z}_{it} &= \begin{pmatrix} z_{it}^1 & 0 & \dots & 0 \\ 0 & z_{it}^2 & \dots & 0 \\ \dots & \dots & \ddots & \vdots \\ 0 & 0 & \dots & z_{it}^d \end{pmatrix}_{[d, q \times d]} \\ \boldsymbol{\beta} &= (\boldsymbol{\beta}^{1'}, \dots, \boldsymbol{\beta}^{d'})'_{[p \times d, 1]} \\ \mathbf{a}_i &= (\mathbf{a}_i^{1'}, \dots, \mathbf{a}_i^{d'})'_{[q \times d, 1]} \sim N(\mathbf{0}, \boldsymbol{\Sigma}_a) \\ \boldsymbol{\varepsilon}_i &= (\boldsymbol{\varepsilon}'_{i1}, \dots, \boldsymbol{\varepsilon}'_{iT_i})'_{[d \times T_i, 1]}, \boldsymbol{\varepsilon}_{it} \sim N(\mathbf{0}, \boldsymbol{\Sigma}_\varepsilon) \end{aligned}$$

The marginal distribution of \mathbf{Y}_i is assumed multivariate normal with mean:

$$E(\mathbf{Y}_i) = (\mathbf{I}_{T_i} \otimes \boldsymbol{\Lambda})\mathbf{X}_i\boldsymbol{\beta}$$

and variance-covariance matrix

$$\begin{aligned} V(\mathbf{Y}_i) &= (\mathbf{I}_{T_i} \otimes \boldsymbol{\Lambda})V(\boldsymbol{\eta}_i)(\mathbf{I}_{T_i} \otimes \boldsymbol{\Lambda})' + \mathbf{I}_{T_i} \\ &\quad \otimes \text{diag}(\boldsymbol{\tau}_1^2, \dots, \boldsymbol{\tau}_d^2) \end{aligned}$$

The first term in $V(\mathbf{Y}_i)$ denotes the variances and covariance of the latent factors and the last term the variances of the error term, $\boldsymbol{\varepsilon}_{it}$. The mean and variance-covariance matrix of $\boldsymbol{\eta}_i$ are $E(\boldsymbol{\eta}_i) = \mathbf{X}_i\boldsymbol{\beta}$ and $V(\boldsymbol{\eta}_i) = \mathbf{Z}_i\boldsymbol{\Sigma}_a\mathbf{Z}_i' + \mathbf{I}_{T_i} \otimes \boldsymbol{\Sigma}_\varepsilon$, respectively.

To estimate the LFLMM by EM, we summarize it below, as proposed by [1]. Before going into detail, we observe that $\{\boldsymbol{\eta}_i, \mathbf{a}_i\}$ is treated as missing data. Hence, the complete dataset is $\{\mathbf{Y}_i, \mathbf{X}_i, \mathbf{Z}_i, \boldsymbol{\eta}_i, \mathbf{a}_i\}$ whereas the observed data is $\{\mathbf{Y}_i, \mathbf{X}_i, \mathbf{Z}_i\}$. It follows that the complete data likelihood is:

$$L = \prod_{i=1}^N P(\mathbf{Y}_i | \boldsymbol{\eta}_i, \boldsymbol{\Lambda}, \boldsymbol{\tau}^2) P(\boldsymbol{\eta}_i | \mathbf{X}_i, \mathbf{Z}_i, \mathbf{a}_i, \boldsymbol{\beta}, \boldsymbol{\Sigma}_\varepsilon) P(\mathbf{a}_i | \boldsymbol{\Sigma}_a) \tag{5}$$

The corresponding complete data loglikelihood is:

$$\log L = \sum_{i=1}^N \left[\log P(\mathbf{Y}_i | \boldsymbol{\eta}_i, \boldsymbol{\Lambda}, \boldsymbol{\tau}^2) + \log P(\boldsymbol{\eta}_i | \mathbf{X}_i, \mathbf{Z}_i, \mathbf{a}_i, \boldsymbol{\beta}, \boldsymbol{\Sigma}_\varepsilon) + \log P(\mathbf{a}_i | \boldsymbol{\Sigma}_a) \right] \tag{6}$$

where

$$\begin{aligned} \sum_{i=1}^N \log P(\mathbf{Y}_i | \boldsymbol{\eta}_i, \boldsymbol{\Lambda}, \boldsymbol{\tau}^2) &= \sum_{i=1}^N \sum_{j=1}^J \log P(\mathbf{Y}_{ij} | \boldsymbol{\eta}_i, \boldsymbol{\Lambda}_j, \boldsymbol{\tau}_j^2) \\ &= \sum_{i=1}^N \sum_{j=1}^J \left[-\frac{n_i}{2} \log \boldsymbol{\tau}_j^2 - \frac{1}{2\boldsymbol{\tau}_j^2} (\mathbf{Y}_{ij} - \boldsymbol{\eta}'_i \boldsymbol{\Lambda}_j)' (\mathbf{Y}_{ij} - \boldsymbol{\eta}'_i \boldsymbol{\Lambda}_j) \right] \end{aligned} \tag{7}$$

$$\begin{aligned} \sum_{i=1}^N \log P(\boldsymbol{\eta}_i | \mathbf{X}_i, \mathbf{Z}_i, \boldsymbol{\beta}, \mathbf{a}_i, \boldsymbol{\Sigma}_\varepsilon) &= \sum_{i=1}^N \sum_{t=1}^{T_i} \log P(\boldsymbol{\eta}_{it} | \mathbf{X}_i, \mathbf{Z}_i, \mathbf{a}_i, \boldsymbol{\beta}, \boldsymbol{\Sigma}_\varepsilon) \\ &= \sum_{i=1}^N \sum_{t=1}^{T_i} \left[-\frac{1}{2} \log |\boldsymbol{\Sigma}_\varepsilon| - \frac{1}{2} (\boldsymbol{\eta}_{it} - \mathbf{X}_{it}\boldsymbol{\beta} - \mathbf{Z}_i\mathbf{a}_i)' \boldsymbol{\Sigma}_\varepsilon^{-1} (\boldsymbol{\eta}_{it} - \mathbf{X}_{it}\boldsymbol{\beta} - \mathbf{Z}_i\mathbf{a}_i) \right] \end{aligned} \tag{8}$$

$$\sum_{i=1}^N \log P(\mathbf{a}_i | \Sigma_a) = -\frac{1}{2} \sum_{i=1}^N \left[\log |\Sigma_a| - \frac{1}{2} \mathbf{a}_i' \Sigma_a^{-1} \mathbf{a}_i \right] \tag{9}$$

Let the θ denote the parameter vector $(\Lambda, \tau^2, \beta, \text{ and } \Sigma_\varepsilon)$, $\theta^{(w)}$ be the ML estimate of θ at the w th iteration for $w = 0, 1, \dots$, and $Q(\theta | \theta^{(w)})$ the expectation of the joint loglikelihood for the complete data $\{Y_i, X_i, Z_i, \eta_i, \mathbf{a}_i\}$ conditional on the observed data $\{Y_i, X_i, Z_i\}$:

$$Q(\theta | \theta^{(w)}) = E \left\{ \log L(\theta | Y_i, X_i, Z_i, \eta_i, \mathbf{a}_i) \middle| Y_i, X_i, Z_i, \theta^{(w)} \right\} \tag{10}$$

Then the $(w + 1)$ th iteration of the EM algorithm consists of (i) the E-step, which is the expectation of the joint loglikelihood computed according to (10) and (ii) the M-step, which maximizes $Q(\theta | \theta^{(w)})$ to yield $\theta^{(w+1)}$. Further details on EM estimation of the LFLMM can be found in [1].

3. The Supplemented EM

Below we discuss the Supplemented EM algorithm, denoted as SEM. Before going into detail, we observe that the main purpose of this study is to estimate the standard errors of the fixed effects, β .

Consider the mapping M defined by iteration w of the EM algorithm:

$$\beta^{(w+1)} = M(\beta^{(w)}), \text{ for } w = 0, 1, \dots$$

when the parameter vector converges to β^* , we obtain $\beta^* = M(\beta^*)$. For $M(\beta)$ continuous we have by Taylor expansion in the neighbourhood of β^*

$$\beta^{(w+1)} = M(\beta^{(w)}) \approx M(\beta^*) + DM(\beta^{(w)} - \beta^*) = \beta^* + DM(\beta^{(w)} - \beta^*) \tag{11}$$

where

$$DM = \left(\frac{\partial M_h(\beta)}{\partial \beta_g} \right) \bigg|_{\beta=\beta^*} \tag{12}$$

$g = 1, 2, \dots, k$ and $h = 1, 2, \dots, k$ is the $k \times k$ Jacobian matrix of $M(\beta) = (M_1(\beta), \dots, M_k(\beta))$ evaluated at the ML estimate of β with $k = p \times d$. DM is known as the rate matrix. To obtain the loglikelihood of β , we consider the complete data density of the LFLMM:

$$\begin{aligned} & f(\{Y_i, X_i, Z_i, \eta_i, \mathbf{a}_i\} | \theta) \\ & = f(\{Y_i, X_i, Z_i\} | \theta) f(\{\eta_i, \mathbf{a}_i\} | \{Y_i, X_i, Z_i\}, \theta) \end{aligned}$$

where $f(\{Y_i, X_i, Z_i\} | \theta)$ is the density of the observed data and $f(\{\eta_i, \mathbf{a}_i\} | \{Y_i, X_i, Z_i\}, \theta)$ the density of missing data, given the observed data. Thus, the loglikelihood of β given the complete data is:

$$\log L(\beta | \{Y_i, X_i, Z_i, \eta_i, \mathbf{a}_i\}) = \log L(\beta | \{Y_i, X_i, Z_i\}) + \log f(\{\eta_i, \mathbf{a}_i\} | \{Y_i, X_i, Z_i\}, \beta) \tag{13}$$

where $\log L(\beta | \{Y_i, X_i, Z_i\})$ is the observed-data loglikelihood and $\log L(\beta | \{Y_i, X_i, Z_i, \eta_i, \mathbf{a}_i\})$ is the complete data loglikelihood.

The asymptotic variance-covariance matrix of β , $V(\beta)$, is the inverse of the observed information matrix (I_0). In the case of the LFLMM, the observed data is $\{Y_i, X_i, Z_i\}$ so that $V(\beta)$ is:

$$V(\beta) = I_0^{-1}(\beta | \{Y_i, X_i, Z_i\}) \tag{14}$$

where $I_o(\beta|\{Y_i, X_i, Z_i\})$ is the information matrix of the observed data loglikelihood (which is assumed to exist). That is [6,11]:

$$I_o(\beta|\{Y_i, X_i, Z_i\}) = -E \left[\frac{\partial^2 \log L(\beta|\{Y_i, X_i, Z_i\})}{\partial \beta \cdot \partial \beta} \right] \tag{15}$$

Equation (15) is difficult to evaluate directly using the EM algorithm [6,11]. As a way out, [7] suggested to evaluate the complete data information matrix:

$$I_o(\beta|\{Y_i, X_i, Z_i, \eta_i, a_i\}) = -E \left[\frac{\partial^2 \log L(\beta|\{Y_i, X_i, Z_i, \eta_i, a_i\})}{\partial \beta \cdot \partial \beta} \right] \tag{16}$$

The conditional complete data information, given the observed data evaluated at $\beta = \beta^*$, is:

$$I_{oc} = E[I_o(\beta|\{Y_i, X_i, Z_i, \eta_i, a_i\})|\{Y_i, X_i, Z_i\}, \beta^*] \tag{17}$$

After taking second derivatives, averaging over $f(\{\eta_i, a_i\}|\{Y_i, X_i, Z_i\}, \beta)$, and evaluating at $\beta = \beta^*$, Equation (13) implies:

$$I_o(\beta^*|\{Y_i, X_i, Z_i\}) = I_{oc} - I_{om} \tag{18}$$

where the missing information matrix (I_{om}) is

$$I_{om} = E \left[-\frac{\partial^2 \log f(\{\eta_i, a_i\}|\{Y_i, X_i, Z_i\}, \beta^*)}{\partial \beta \cdot \partial \beta} \right] \tag{19}$$

ref [12] interpreted Equation (18) as

$$\textit{observed information} = \textit{complete information} - \textit{missing information}$$

and called it the “missing information principle”. Equation (18) can be written as:

$$I_o(\beta^*|\{Y_i, X_i, Z_i\}) = (I - I_{om}I_{oc}^{-1})I_{oc}, \tag{20}$$

where I is the $k \times k$ identity matrix and $I_{om}I_{oc}^{-1}$ is the matrix of the fraction of missing information [7,11]. According to [13], the rate of convergence of the EM algorithm is determined by the fraction of missing information in the neighborhood of β^* :

$$DM = I_{om}I_{oc}^{-1} \tag{21}$$

Substituting $DM = I_{om}I_{oc}^{-1}$ into Equation (20) and inverting, the asymptotic variance-covariance matrix of β^* , $V(\beta^*)$ is:

$$V(\beta^*) = I_{oc}^{-1}(I - DM)^{-1} \tag{22}$$

From the equality $(I - P)^{-1} = (I - P + P)(I - P)^{-1} = I + P(I - P)^{-1}$ it follows that:

$$V(\beta^*) = I_{oc}^{-1} \left\{ I + DM(I - DM)^{-1} \right\} = I_{oc}^{-1} + I_{oc}^{-1}DM(I - DM)^{-1} \tag{23}$$

or

$$V(\beta^*) = I_{oc}^{-1} + \Delta V(\beta^*) \tag{24}$$

where $\Delta V(\beta^*)$ is the increase of the diagonal elements of $V(\beta^*)$ related to missing information.

Calculation of the DM matrix can be done using the code and output of the original EM algorithm as follows [6,7]. The DM matrix represents the differential of the parameter mappings during the EM algorithm. Hence, each element of the DM matrix represents

a component-wise increase of the rate of convergence per iteration of the EM algorithm. Let r_{gh} be the (g, h) th element of the *DM* matrix. From Equation (13), we have:

$$r_{gh} = \frac{\partial M_h(\beta^*)}{\partial \beta_g} = \lim_{\beta_g \rightarrow \beta_g^*} \frac{M_h(\beta_1^*, \dots, \beta_{g-1}^*, \beta_g, \beta_{g+1}^*, \dots, \beta_k^*) - M_h(\beta^*)}{\beta_g - \beta_g^*} \tag{25}$$

$$= \lim_{w \rightarrow \infty} \frac{M_h(\beta^{(w)}(g)) - M_h(\beta^*)}{\beta_g^{(w)} - \beta_g^*} \equiv \lim_{w \rightarrow \infty} r_{gh}^{(w)}$$

$g = 1, 2, \dots, k$ and $h = 1, 2, \dots, k$
 where $\beta^{(w)}(g)$ is called the semi-active parameter set

$$\beta^{(w)}(g) = (\beta_1^*, \dots, \beta_{g-1}^*, \beta_g^{(w)}, \beta_{g+1}^*, \dots, \beta_k^*), \quad w = 1, 2, \dots \tag{26}$$

which converges to β_g^* . Note that only the g th component in $\beta^{(w)}(g)$ takes a value different from its maximum likelihood estimate.

To calculate r_{gh} , the Supplemented EM algorithm requires $\theta^* = \{\Lambda^*, \tau^{2*}, \beta^*$ and $\Sigma_\epsilon^*\}$ and $\theta^{(w)} = \{\Lambda^{(w)}, \tau^{2(w)}, \beta^{(w)}$ and $\Sigma_\epsilon^{(w)}\}$ for $w = 1, 2, \dots$ as input. θ^* can be obtained by the EM algorithm using a set of arbitrarily chosen initial parameters θ^{init} including $\theta^{(w)}$ for $w = 1$, i.e., $\theta^{(1)}$. The starting point $\theta^{(1)}$ may, but need not, be close θ^* . The algorithm below closely follows [14,15].

1. Select input: $\theta^{(w)}$ and θ^*
2. Set $\theta^{(w)}$ for $w = 1, 2, \dots$. Then take the E step and M step of the LFLMM EM algorithm to produce $\theta^{(w+1)}$.
3. For rows = $1, 2, \dots, k$:
 - (i) Set $\tilde{\beta}^{(w)}(g)$ be equal to β^* , except for the g th element:
 $(\tilde{\beta}^{(w)}(g) = (\beta_1^*, \dots, \beta_{g-1}^*, \beta_g^{(w)}, \beta_{g+1}^*, \dots, \beta_k^*))$
 - (ii) Run the LFLMM EM algorithm with $\tilde{\beta}^{(w)}(g)$ as the current estimate of β to obtain $\tilde{\beta}^{(w+1)}(g)$.
 - (iii) Calculate the g th row of $r_{gh}^{(w)}$ as

$$r_{gh}^{(w)} = \frac{\tilde{\beta}_h^{(w+1)}(g) - \beta_h^*}{\beta_g^{(w)} - \beta_g^*}, \quad \text{for } h = 1, 2, \dots, k$$

The output after a single run of the Supplemented EM algorithm (Step 1 and 2) are β^{w+1} and $r_{gh}^{(w)}$ $g = 1, 2, \dots, k$ and $h = 1, 2, \dots, k$. Based on the final estimates of DM, $V(\beta^*)$ is calculated using (24). The diagonal elements of V are the variance of β^* .

4. Simulation

To evaluate the statistical properties and computational aspects of the SEM, we set up a simulation study. The number of subjects (N) is set at 500, 1000, and 1500 with six time periods. The number of simulations (S) is set at 50 and 250. The other set-up of the simulations is adopted from [1]. Particularly, we use the same initial values of the parameters of the LFLMM model (12 items, 2 latent factors, and a simple structure to model the relationship between the items and the latent factors). It is done to check if the bias resulting from the Supplemented EM algorithm on LFLMM is in line with the results presented in [1]. Table 1 presents the absolute difference for the true parameters, and the averages of the SEM estimates are calculated as a measure of performance.

Table 1 shows that the absolute difference of $\sigma_{a,11}$ has a range from 0 to 0.0444 ($N = 500$ and $S = 250$). The results are in line with [1] the parameters of the measurement model (factor loadings and error variances) are estimated more precisely than those of the latent mixed regression model. Overall, these results indicate that with the increasing number of

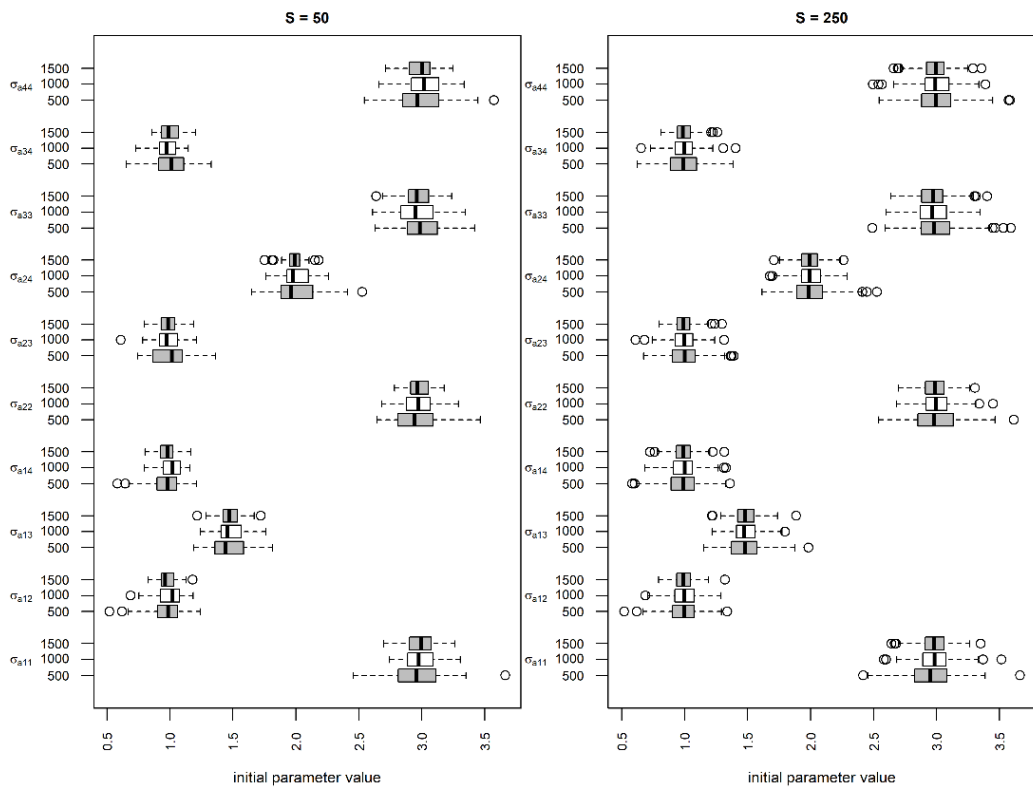
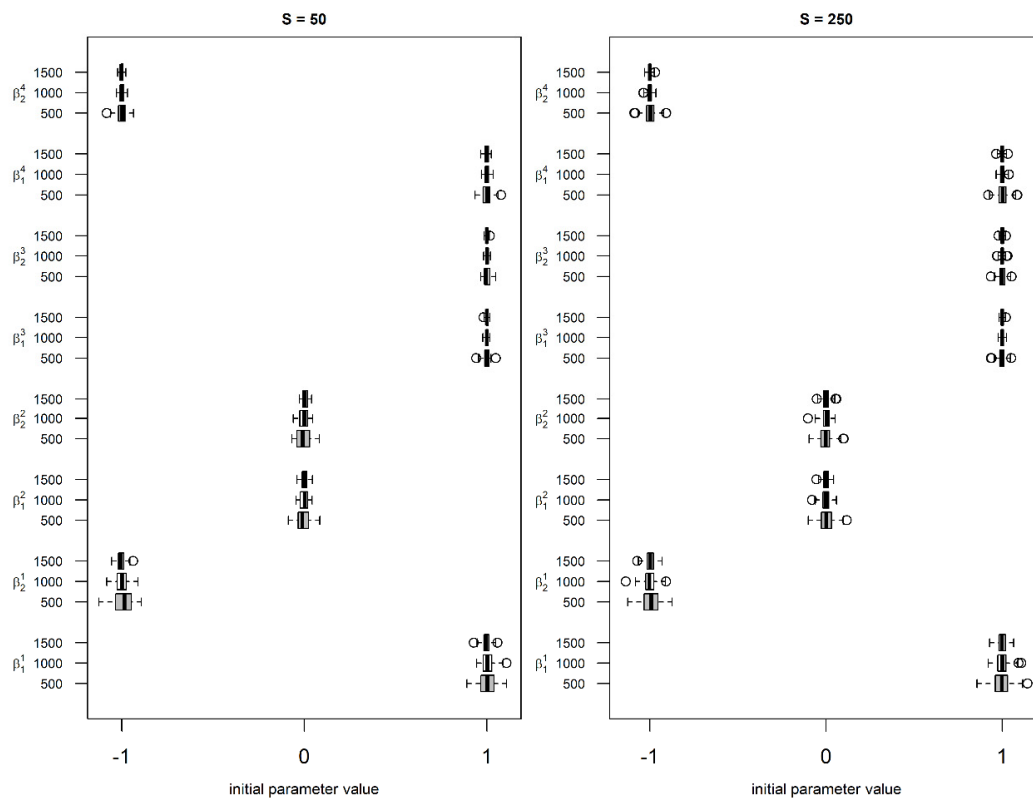


Figure 1. Cont.

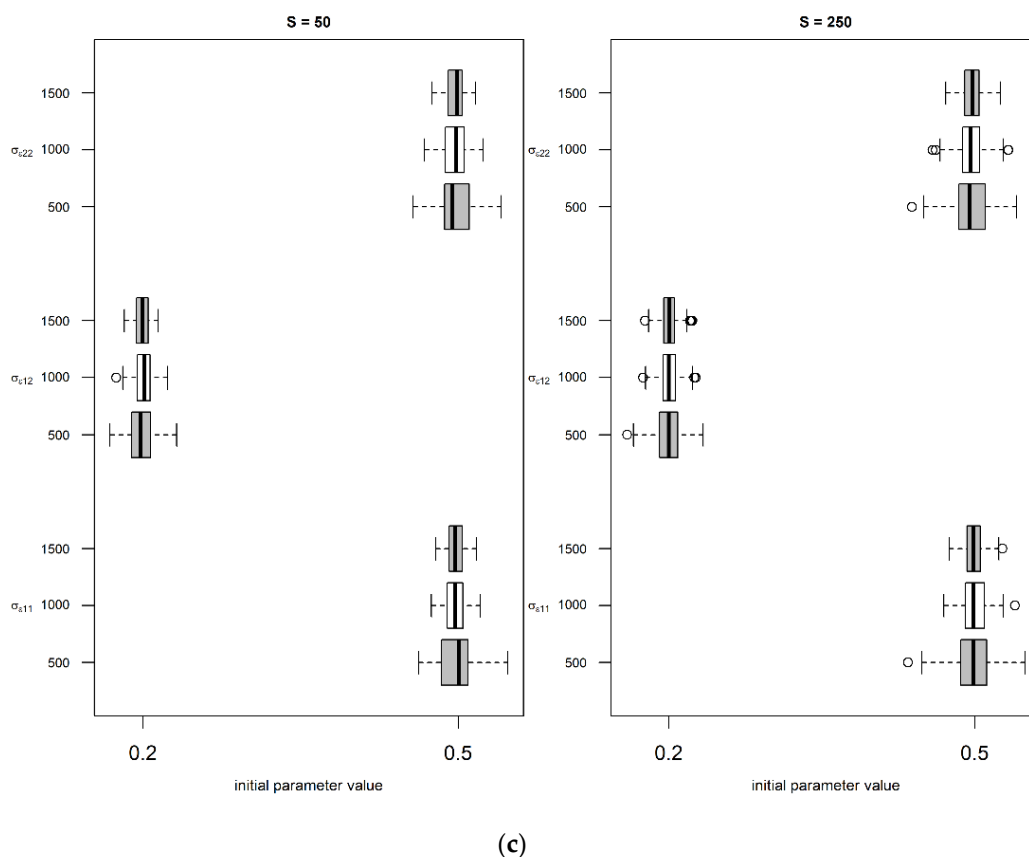


Figure 1. (a) Boxplot the parameters estimate of the latent mixed regression models (β). (b) Boxplot the parameters estimate of the latent mixed regression models (σ_a). (c) Boxplot the parameters estimate of the latent mixed regression models (σ_e).

Table 2. The parameter estimates for $\sqrt{V(\hat{\beta})}$.

Number of Subjects	Parameter	The 2nd Moment		SEM	
		50	250	50	250
500	β_1^1	0.0512	0.0528	0.0249	0.0251
	β_1^2	0.0587	0.0530	0.0187	0.0195
	β_2^1	0.0378	0.0392	0.0110	0.0126
	β_2^2	0.0391	0.0379	0.0164	0.0158
	β_3^1	0.0184	0.0190	0.0245	0.0268
	β_3^2	0.0184	0.0190	0.0212	0.0200
	β_4^1	0.0283	0.0283	0.0105	0.0105
	β_4^2	0.0288	0.0308	0.0145	0.0145
1000	β_1^1	0.0338	0.0327	0.0179	0.0184
	β_1^2	0.0342	0.0339	0.0138	0.0134
	β_2^1	0.0232	0.0232	0.0071	0.0077
	β_2^2	0.0253	0.0232	0.0105	0.0100
	β_3^1	0.0095	0.0100	0.0176	0.0184
	β_3^2	0.0089	0.0095	0.0130	0.0130
	β_4^1	0.0141	0.0130	0.0077	0.0077
	β_4^2	0.0130	0.0130	0.0110	0.0105

Table 2. Cont.

Number of Subjects	Parameter	The 2nd Moment		SEM	
		50	250	50	250
1500	β_1^1	0.0276	0.0253	0.0152	0.0152
	β_1^2	0.0270	0.0265	0.0105	0.0110
	β_2^1	0.0187	0.0182	0.0063	0.0063
	β_2^2	0.0164	0.0192	0.0084	0.0084
	β_3^1	0.0071	0.0071	0.0145	0.0152
	β_3^2	0.0063	0.0071	0.0105	0.0105
	β_4^1	0.0110	0.0105	0.0063	0.0063
	β_4^2	0.0110	0.0105	0.0084	0.0084

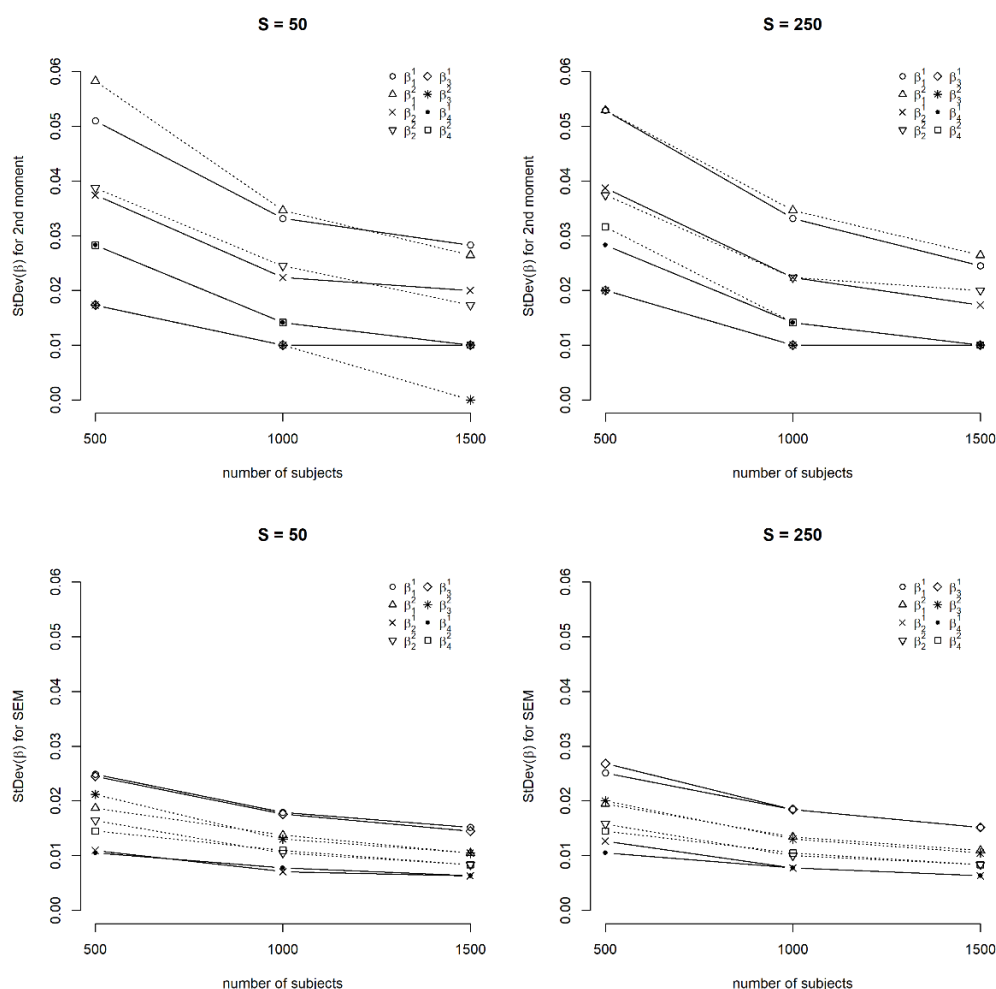


Figure 2. Line plot for the standard deviation of beta.

5. Real Data Example

The real data-set that we used to illustrate the development of the Supplemented EM algorithm is the political attitudes and behavior data of Flemish. The data was designed to include a representative sample of the target population under the Belgian electorate. The Flemish data set (Flemish and Dutch speaking respondents from Brussels Capital Region) consists of 1274 respondents, who have been interviewed three times (1991, 1995, and 1999) [16–18]. There are four latent factors measured on political attitudes of Flemish

used, i.e., Individualism, Nationalism, Ethnocentrism, and Authoritarianism. This data has been analyzed using various methods by several authors, including [19–23]. There are three interesting questions in this real data case, i.e., how Individualism, Nationalism, Ethnocentrism, and Authoritarianism of the Flemish develop over time; whether there is an association between these four developments, and whether the gender of the respondent affects the change patterns of latent developments.

I, N, E, and A in Table 3 correspond to Individualism, Nationalism, Ethnocentrism, and Authoritarianism, respectively. a_{11} and a_{12} are the random intercept and random slope for Individualism. a_{21} and a_{22} are the random intercept and random slope for Nationalism. a_{31} and a_{32} are the random intercept and random slope for Ethnocentrism. a_{41} and a_{42} are the random intercept and random slope for Authoritarianism. The positive correlation of random intercept between a_{11} and a_{21} , a_{11} and a_{31} , a_{11} and a_{41} suggests that the development of Individualism and other political attitudes is highly related, which highest correlated with Ethnocentrism. The results indicate that those who have a better sense of Individualism tend to have a better sense of Nationalism, Ethnocentrism, and Authoritarianism. The results find a positive correlation of random intercept between a_{21} and a_{31} , a_{21} and a_{41} . It suggests that those who have a better sense of Nationalism tend to have a better sense of Ethnocentrism and Authoritarianism, as well as those who have a better sense of Ethnocentrism tend to have a better sense of Authoritarianism. There is also a positive correlation of random slope between a_{12} and a_{22} . It means that if one subject's Individualism decreases over time, then it is reasonable to expect that his or her Nationalism will decrease over time and vice versa. This also holds between Individualism and Ethnocentrism and between Individualism and Authoritarianism. The positive correlation of random slope between a_{22} and a_{32} , meaning that if one subject's Nationalism decreases over time, then it is reasonable to expect that his or her Ethnocentrism will decrease over time. The correlation matrix of random effects confirms that all latent factors have a positive correlation over time.

Table 3. Correlation matrix of random effects.

Random Effects		I		N		E		A	
		a_{11}	a_{12}	a_{21}	a_{22}	a_{31}	a_{32}	a_{41}	a_{42}
I	a_{11}	1	0.892	0.832	0.903	0.966	0.930	0.956	0.880
	a_{12}	0.892	1	0.854	0.878	0.937	0.931	0.913	0.895
N	a_{21}	0.832	0.854	1	0.311	0.864	0.883	0.856	0.835
	a_{22}	0.903	0.878	0.311	1	0.918	0.893	0.899	0.861
E	a_{31}	0.966	0.937	0.864	0.918	1	0.946	0.973	0.919
	a_{32}	0.930	0.931	0.883	0.893	0.946	1	0.946	0.918
A	a_{41}	0.956	0.913	0.856	0.899	0.973	0.946	1	0.873
	a_{42}	0.880	0.895	0.835	0.861	0.919	0.918	0.873	1

The significance of parameter estimate of β is analyzed via the z -values. By using the Supplemented EM algorithm, the standard errors of β for all parameters can be calculated. The standard errors of β are listed in Table 4. Using a 95 percent confidence interval of β , almost all confidence intervals do not include the null value, except the slope of Male on Authoritarianism. Hence there are statistically significant differences in the parameter estimate of β . In other words, all latent factors of Flemish people decrease over time, with Ethnocentrism having the highest rate of decline over time (-0.252) and Nationalism the lowest (-0.177). On average, the Individualism and Nationalism of the male respondent are higher than that of the female. However, Ethnocentrism of the male respondent is lower than that of the female.

Table 4. Parameter estimates of β .

Parameter	Estimate	$SE(\hat{\beta})$	Lower	Upper
Individualism	−0.195	0.004	−0.203	−0.187
Nationalism	−0.177	0.028	−0.231	−0.123
Ethnocentrism	−0.252	0.011	−0.273	−0.231
Authoritarianism	−0.186	0.002	−0.190	−0.182
Slope of Male on Individualism	0.110	0.010	0.091	0.129
Slope of Male on Nationalism	0.219	0.017	0.185	0.253
Slope of Male on Ethnocentrism	−0.038	0.003	−0.045	−0.031
Slope of Male on Authoritarianism	0.022	0.017	−0.011	0.055

6. Conclusions

This paper proposed the Supplemented EM algorithm for LFLMM in estimating the asymptotic variance-covariance matrix as a by-product of the EM estimator for the case of fixed variables in the model. Results from simulation studies suggest that the Supplemented EM algorithm can estimate the model very close to the initial parameters.

As a result of the development of EM algorithm of LFLMM, the Supplemented EM algorithm is very slow to converge, as stated by [1], especially when the number of simulations is 250 times with 1500 subjects. For this reason, further research is needed to find techniques that can be used to accelerate the speed of the algorithm. Several approaches to speed the EM algorithm have been proposed and can be found in [24–26] (the ECM algorithm), [27] (the ECME algorithm), and [28] (the Parameter-Expanded EM algorithm).

Author Contributions: Conceptualization, Y.A.; methodology, Y.A.; software, Y.A.; validation, K.A.N. and H.F.; formal analysis, Y.A.; data curation, Y.A. and T.T.; writing—original draft preparation, Y.A.; writing—review and editing, K.A.N., H.F., A.S. and T.T.; visualization, T.T.; supervision, K.A.N., H.F. and A.S. All authors have read and agreed to the published version of the manuscript.

Funding: This research was funded by RUG and IPB University.

Institutional Review Board Statement: Not applicable.

Informed Consent Statement: Not applicable.

Data Availability Statement: The data that support the findings of this study are available from ISPO which were used under license, and so are not publicly available. Data are however available from the authors upon reasonable request and with permission of ISPO.

Conflicts of Interest: The authors declare no conflict of interest.

References

- An, X.; Yang, Q.; Bentler, P.M. A latent factor linear mixed model for high-dimensional longitudinal data analysis. *Stat. Med.* **2013**, *32*, 4229–4239. [[CrossRef](#)] [[PubMed](#)]
- Kondaurova, M.V.; Bergeson, R.R.; Xu, H.; Kitamura, C. Affective Properties of Mothers' Speech to Infants with Hearing Impairment and Cochlear Implants. *J. Speech Lang. Hear. Res.* **2015**, *58*, 590–600. [[CrossRef](#)]
- Wang, J.; Luo, S. Multidimensional latent trait linear mixed model: An application in clinical studies with multivariate longitudinal outcomes. *Stat. Med.* **2017**, *36*, 3244–3256. [[CrossRef](#)]
- Ng, S.K.; Krishnan, T.; McLachlan, G. The EM algorithm. In *Handbook of Computational Statistics*; Springer: Berlin, Germany, 2004; pp. 137–168, ISBN 9783642215513.
- McLachlan, G.J.; Krishnan, T. *The EM Algorithm and Extensions Second Edition*, 2nd ed.; Wiley: New York, NY, USA, 2007; ISBN 9780471201700.
- Meng, A.X.; Rubin, D.B. Using EM to Obtain Asymptotic Variance-Covariance Matrices: The SEM Algorithm. *J. Am. Stat. Assoc.* **1991**, *86*, 899–909. [[CrossRef](#)]
- Cai, L. SEM of another flavour: Two new applications of the supplemented EM algorithm. *Br. J. Math. Stat. Psychol.* **2008**, *61*, 309–329. [[CrossRef](#)] [[PubMed](#)]
- Cai, L.; Lee, T.; Lee, T. Covariance Structure Model Fit Testing Under Missing Data: An Application of the Supplemented EM Algorithm. *Multivar. Behav. Res.* **2009**, *44*, 281–304. [[CrossRef](#)] [[PubMed](#)]

9. Tian, W.; Cai, L.; Thissen, D.; Xin, T. Numerical Differentiation Methods for Computing Error Covariance Matrices in Item Response Theory Modeling: An Evaluation and a New Proposal. *Educ. Psychol. Meas* **2012**, *73*, 412–439. [[CrossRef](#)]
10. Caraka, R.E.; Noh, M.; Chen, R.C.; Lee, Y.; Gio, P.U.; Pardamean, B. Connecting climate and communicable disease to penta helix using hierarchical likelihood structural equation modelling. *Symmetry* **2021**, *13*, 657. [[CrossRef](#)]
11. Pritikin, J.N. A comparison of parameter covariance estimation methods for item response models in an expectation-maximization framework. *Cogent Psychol.* **2017**, *4*, 1–11. [[CrossRef](#)]
12. Orchard, T.; Woodbury, M. A Missing Information Principle: Theory and Applications. In *Theory of Statistics*; University of California Press: Berkeley, CA, USA, 1972; Volume 1, pp. 697–715. Available online: https://projecteuclid.org/download/pdf_1/euclid.bsm/1200514117 (accessed on 1 June 2021).
13. Dempster, A.P.; Laird, N.M.; Rubin, D.B. Maximum Likelihood from Incomplete Data via the EM Algorithm A. *J. R. Stat. Soc. Ser. B* **1977**, *39*, 1–38.
14. Little, R.J.A.; Rubin, D.B. *Statistical Analysis with Missing*; Wiley: New York, NY, USA, 2002; ISBN 3175723993.
15. Abel, G.J. International Migration Flow Table Estimation. Ph.D. Thesis, University of Southampton, Southampton, UK, 2009.
16. Interuniversitair Steunpunt Politieke-Opinieonderzoek. *General Election Study: Codebook and Questionnaire*; ISPO: Leuven, Belgium, 1991; ISBN 9067841161.
17. Interuniversitair Steunpunt Politieke-Opinieonderzoek. *General Election Study: Codebook and Questionnaire*; ISPO: Leuven, Belgium, 1995; ISBN 9067841366.
18. Interuniversitair Steunpunt Politieke-Opinieonderzoek. *General Election Study: Codebook and Questionnaire*; ISPO: Leuven, Belgium, 1999.
19. Billiet, J. Church Involvement, Individualism, and Ethnic Prejudice among Flemish Roman Catholics: New Evidence of a Moderating Effect. *J. Sci. Study Relig.* **1995**, *34*, 224–233. [[CrossRef](#)]
20. Billiet, J.; Coffe, H.; Maddens, B. Een Vlaams-nationale identiteit en de houding tegenover allochtonen in een longitudinaal perspectief. In *Proceedings of the Paper Presented at the Marktdag Sociologie*; Universitaire Pers Leuven: Leuven, Belgium, 2005.
21. Toharudin, T.; Oud, J.H.L.; Billiet, J.B. Assessing the relationships between Nationalism, Ethnocentrism, and Individualism in Flanders using Bergstrom’s approximate discrete model. *Stat. Neerl.* **2008**, *62*, 83–103. [[CrossRef](#)]
22. Toharudin, T.; Oud, J.H.L.; Billiet, J.; Folmer, H. Measuring Authoritarianism with Different Sets of Items in a Longitudinal Study. In *Methods, Theories, Andempirical Applications in the Social Sciences*; Salzborn, S., Davidov, E., Reinecke, J., Eds.; Springer: Heidelberg, Germany, 2012; pp. 193–200, ISBN 9783531188980.
23. Angraini, Y.; Toharudin, T.; Folmer, H.; Oud, J.H.L. The Relationships between Individualism, Nationalism, Ethnocentrism, and Authoritarianism in Flanders: A Continuous Time-Structural Equation Modeling Approach. *Multivar. Behav. Res.* **2014**, *49*, 41–53. [[CrossRef](#)] [[PubMed](#)]
24. Meng, X.; Rubin, D.B. Maximum likelihood estimation via the ECM algorithm: A general framework. *Biometrika* **1993**, *80*, 267–278. [[CrossRef](#)]
25. Van Dyk, D.A.; Meng, X.; Rubin, D.B. Maximum Likelihood Estimation via the ECM Algorithm: Computing The Asymptotic Variance. *Stat. Sin.* **1995**, *5*, 55–75.
26. Li, H.; Tian, W. Slashed lomax distribution and regression model. *Symmetry* **2020**, *12*, 1877. [[CrossRef](#)]
27. Liu, B.Y.C.; Rubin, D.B. The ECME algorithm: A simple extension of EM and ECM with faster monotone convergence. *Biometrika* **1994**, *81*, 633–648. [[CrossRef](#)]
28. Liu, C.; Rubin, D.B.; Wu, Y.N. Parameter Expansion to Accelerate EM: The PX-EM Algorithm. *Biometrika* **1998**, *85*, 755–770. [[CrossRef](#)]