

EFFICIENT MAXIMUM LIKELIHOOD ESTIMATION OF MULTIPLE MEMBERSHIP LINEAR MIXED MODELS, WITH AN APPLICATION TO EDUCATIONAL VALUE-ADDED ASSESSMENTS

ANDREW T. KARL AND YAN YANG

Arizona State University

SHARON L. LOHR

Westat

ABSTRACT. The generalized persistence (GP) model, developed in the context of estimating “value added” by individual teachers to their students’ current and future test scores, is one of the most flexible value-added models in the literature. Although developed in the educational setting, the GP model can potentially be applied to any structure where each sequential response of a lower-level unit may be associated with a different higher-level unit, and the effects of the higher-level units may persist over time. The flexibility of the GP model, however, and its multiple membership random effects structure lead to computational challenges that have limited the model’s availability. We develop an EM algorithm to compute maximum likelihood estimates efficiently for the GP model, making use of the sparse structure of the random effects and error covariance matrices. The algorithm is implemented in the package GPvbm in R statistical software. We give examples of the computations and illustrate the gains in computational efficiency achieved by our estimation procedure.

NOTICE

This is the author’s version of a work that was accepted for publication in *Computational Statistics & Data Analysis*. Changes resulting from the publishing process, such as peer review, editing, corrections, structural formatting, and other quality control mechanisms may not be reflected in this document. Changes may have been made to this work since it was submitted for publication. A definitive version was subsequently published in *Computational Statistics & Data Analysis*, [VOL59, March, (2013)] DOI:10.1016/j.csda.2012.10.004

1. INTRODUCTION

Multilevel mixed models are popular for describing data with complex dependence structure. The units on which primary measurements are taken (usually those at the lowest level) each belong to one or more units at higher levels. In

a nested (hierarchical) two-level model, each unit at the lowest level belongs to exactly one higher-level unit. In a multiple membership structure (Browne et al., 2001), a lower-level unit may be associated with multiple higher-level units. This structure is common with non-static populations, and we study multiple membership models in which a lower-level unit is sequentially associated with different higher-level units. Thus, a child in foster care may live with multiple families; a patient may see multiple doctors; a deer may visit multiple salt licks; a worker may have multiple employers; a person may attend multiple therapy groups; a student may have multiple teachers. Fielding and Goldstein (2006) describe multiple membership models and give examples of their use. The multiple membership structure induces a complex dependence structure in the data. Lower-level units are correlated whenever they share any higher-level unit, so the covariance matrix will not have a block diagonal structure as in the nested model.

The complex covariance structure of multiple membership mixed models makes computations challenging, particularly with large data sets. Computational methods that have been developed for nested hierarchical models and other special cases of linear mixed models often will not work. In this paper we develop an EM algorithm to compute maximum likelihood estimates for a class of longitudinal multiple membership models that are applicable in many settings. In the class of models considered, lower-level units are associated with multiple higher-level units in sequence, and a response is recorded on a lower-level unit after the association with each higher-level unit. If the population contains a large number of higher-level units, and the number of lower-level units associated with each higher-level unit is bounded, the covariance matrix will be sparse. The algorithm exploits sparseness of the covariance matrix to speed computations. This sparseness is achieved in many multiple membership settings since, for example, there are upper bounds on the number of patients a doctor can see or the number of students in a teacher’s class.

The application motivating this research comes from value-added models (VAMs) in educational evaluation. A VAM score for a teacher is intended to estimate the “value added” by that teacher to students’ knowledge—how much more (or less) students’ scores changed under that teacher than they would be expected to change under an “average” teacher—by apportioning students’ progress on standardized tests to the teachers or schools that have taught those students. Braun et al. (2010) describe some of the potential uses of VAMs and discuss issues associated with using them to evaluate teachers and schools.

While a variety of different models are used (see Lohr (2012) for a review of common VAMs), in this paper we primarily consider the generalized persistence (GP) model developed by Mariano et al. (2010), one of the most flexible models in the literature. In the GP model, each student is followed over T grades with a different teacher in each grade, and receives a score on a standardized test at the end of each grade. Each student therefore “belongs” to up to T different teachers, resulting in a multiple membership structure. The GP model, like other mixed models used in the value-added context (Sanders et al., 1997; Rowan et al., 2002; McCaffrey et al., 2003, 2004, 2005; Lockwood et al., 2007), uses a longitudinal database of student scores and models the scores with random teacher intercepts. Under this scenario, the empirical best linear unbiased predictors (EBLUPs) for random teacher intercepts are the teacher VAM scores. In this paper we use the term “teacher effect” to represent the VAM score of a teacher but note that, as

observed by Lockwood et al. (2007), these teacher effects measure “unexplained heterogeneity at the classroom level,” and not necessarily the causal effect of the teacher.

The GP model is distinguished from others in the VAM literature by how it attributes a student’s performance to current and prior teachers. If the effects of good teaching persist, one would expect that students of a good teacher in year 1 would do well on the test in year 1 and would continue to do well on the tests in future years. The Educational Value-Added Assessment System (EVAAS) model (Sanders et al., 1997), a complete persistence model, assumes that the effect of a teacher persists undiminished over all subsequent years of his or her students’ achievement. This complete persistence assumption, also proposed by Raudenbush and Bryk (2002), implies that each teacher has one VAM score: the effect of a teacher in year g on his or her students’ test scores is the same for their tests in each of years g, \dots, T . The complete persistence assumption simplifies the covariance structure, and the EVAAS model is implemented in SAS software (Wright et al., 2010). A model proposed by McCaffrey et al. (2004) allows the effect of a teacher on students’ scores to decay in future years, though the effects are otherwise perfectly correlated. Lockwood et al. (2007) refer to this structure as variable persistence (VP). In the VP model, each teacher has one estimated effect, but the impact on students’ future year scores is reduced by a multiplicative factor in each year. The multipliers, called persistence parameters, are estimated from the data.

The GP model allows a much more general structure for the effects of a current teacher on future test scores. In the GP model, a teacher in year g has a different effect on his or her students’ scores in each year from $t = g, \dots, T$, and the $(T - g + 1)$ effects of that teacher have an unstructured covariance matrix to allow the effects to be correlated. The EVAAS model is a special case of the GP model in which the current and future effects of a teacher are assumed to be identical. The general correlation structure in the GP model allows much more detailed exploration of the patterns of teacher effects, but greatly complicates the problem of computing estimates.

Hill and Goldstein (1998) estimate a class of multiple membership models using an iterative generalized least squares algorithm, and Browne et al. (2001) employ Monte Carlo Markov chain techniques. Likewise, Mariano et al. (2010) use Bayesian methods to estimate the parameters for the GP model using data from a large urban school district. To obtain a proper posterior distribution, however, a Bayesian approach to computations requires that an informative prior distribution be adopted for the covariance parameters. As investigated in their paper, different priors often result in different estimates of model parameters and teacher effects. A maximum likelihood (ML) approach avoids the need for priors, although ML estimation of even the simpler VP model has been “practically infeasible for all but small data sets” (Lockwood et al., 2007) up to this point. In this paper we use the sparseness of the covariance and design matrices to develop an efficient EM algorithm for calculating ML estimates of parameters in the GP and VP models. We implement the method in the user-friendly GPvam package (Karl et al., 2012) in R statistical software (R Core Team, 2013). This development makes the GP and VP models more accessible for use in practice, and provides an alternative to the Bayesian calculations implemented by Mariano et al. (2010).

While the GP model was developed for educational applications, the model and the computational methods in this paper apply in many other settings as well. For example, Ash et al. (2012) note the similarity between the problems of evaluating teacher performance on the basis of student outcomes, and evaluating hospital and physician performance on the basis of patient outcomes. The multiple membership structure also arises in social network data (Airoldi et al., 2008). In another example, Browne et al. (2001) and Goldstein et al. (2000) describe a multiple membership model used to study Belgian household migration with complete persistence, measuring the propensity of individuals to change household membership. The GP model is a good candidate for the Belgian household data since the similarity of former roommates may decrease over time. Browne et al. (2001) also describe an application in which a hospital patient is cared for by different nurses and the contribution of each nurse to the patient's progress is estimated.

The paper is organized as follows. Section 2 reviews the models studied and lays out the foundation for ML estimation. Section 3 presents the EM algorithm for the estimation of the model. The details of the implementation of the model in R appear in Section 4. The computational methods are applied to a data set from a large urban school district in Section 5 to demonstrate the capabilities of the estimation procedure and software.

2. MODEL SPECIFICATION

The GP model (Mariano et al., 2010) and other structures considered in this paper model responses of the lower-level units as follows:

$$(1) \quad y_{ig} = \mathbf{x}'_{ig}\boldsymbol{\beta} + \mathbf{s}'_{ig}\boldsymbol{\eta} + \epsilon_{ig}$$

where y_{ig} is a response for unit i at time g for $i = 1, \dots, n$, and $g \in A_i$; A_i , a subset of $\{1, \dots, T\}$, is the set of times for which unit i is observed. The vector of all responses is $\mathbf{y} = (\mathbf{y}'_1, \dots, \mathbf{y}'_n)'$, where $\mathbf{y}_i = (y_{ig}, g \in A_i)$ is the vector of responses for unit i . The matrix \mathbf{X} , with rows \mathbf{x}'_{ig} for $g \in A_i$ and $i = 1, \dots, n$, is the design matrix of covariates for the fixed effects parameter vector $\boldsymbol{\beta}$.

The random effects vector $\boldsymbol{\eta} \sim N(0, \mathbf{G})$ contains random intercepts for the higher-level units (and also for the lower-level units if desired). Each measurement on a lower-level unit is associated with multiple higher-level units as specified by the design matrix \mathbf{S} which has rows \mathbf{s}'_{ig} for $g \in A_i$ and $i = 1, \dots, n$. The multiple membership structure arises because rows of the \mathbf{S} matrix may contain multiple nonzero values. The vector of error terms for lower-level unit i , $\boldsymbol{\epsilon}_i = \{\epsilon_{ig}, g \in A_i\}$, is assumed to be normally distributed with mean $\mathbf{0}$ and covariance matrix \mathbf{R}_i . The lower-level units are assumed to be independent conditionally on the random intercepts contained in $\boldsymbol{\eta}$, so $\boldsymbol{\epsilon} = (\boldsymbol{\epsilon}'_1, \dots, \boldsymbol{\epsilon}'_n)' \sim N(\mathbf{0}, \mathbf{R})$ and \mathbf{R} is block diagonal with blocks $\mathbf{R}_1, \dots, \mathbf{R}_n$. The error terms $\boldsymbol{\epsilon}$ are also assumed to be independent of the effects in $\boldsymbol{\eta}$.

The observations taken together thus have the general form of a linear mixed model:

$$(2) \quad \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{S}\boldsymbol{\eta} + \boldsymbol{\epsilon},$$

with $\text{Cov}(\mathbf{y}) = \mathbf{V} = \mathbf{S}\mathbf{G}\mathbf{S}' + \mathbf{R}$. The log-likelihood based on the observed data \mathbf{y} from model (2) is

$$(3) \quad l(\boldsymbol{\Psi}; \mathbf{y}) \propto -\frac{1}{2} \log |\mathbf{V}| - \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})' \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

where Ψ is a vector of the unique model parameters from β , G , and R . We assume throughout that the sufficient conditions for consistency and asymptotic normality of the ML estimates given by Broatch and Lohr (2012) are met. In addition to usual regularity conditions, Broatch and Lohr (2012) assume that T is bounded and that the number of lower-level units associated with each higher-level unit is bounded, i.e., that the sum of each column of S is bounded. The model is also assumed to be identifiable. Let ψ_1, \dots, ψ_q be the parameters in Ψ that are components of G , and R and write $V = \sum_{j=1}^q \psi_j \Sigma_j$. Then the model will be identifiable when the matrices Σ_j are linearly independent for $j = 1, \dots, q$. In practical terms, the GP model will be identifiable as long as there is sufficient mixing in the population so that students in a class progress to a variety of different teachers as they continue through school. Briggs and Weeks (2011), fitting a VP VAM with schools as higher-level units instead of teachers, find that not all parameters are identifiable because most students in their data set move from one grade to the next as a cohort within the same school with insufficient mixing.

Note that the model formulation allows lower-level units to be missing observations for some times. In this paper we assume that observations are missing at random and that the parameters governing the outcome process are distinct from those characterizing the missingness process, yielding a valid likelihood-based analysis under the specified model (Little and Rubin, 2002). McCaffrey and Lockwood (2011) and Karl et al. (2011) propose joint models for test score data and missingness indicators to accommodate data with informative missingness, but we do not consider such models here.

We handle the missing teacher links resulting from missing student observations by assuming that the student was taught by an average teacher in that year. For example, when modeling scores from grades 1, 2, and 3, if a student enters the school at grade 2, we do not link that student's second and third grade scores to any of the first grade teachers. This approach was also used by Lockwood et al. (2007).

To make these ideas concrete, in the remainder of this section we present specific models considered in the educational setting, in which the lower-level units are students, the higher-level units are teachers, and the measurement y_{ig} is a test score of student i in year g . Rather than introducing new notation for each model, we recast the model terms to match the notation of Equations (1) and (2) so that the definitions of η , G , R and S depend on the chosen model. This streamlines the discussion of the estimation of the parameters.

2.1. Generalized Persistence Model. The GP VAM (Mariano et al., 2010) models student scores using information about the history of observations on each student and each student's teacher-history. It estimates the effect of teachers on students in the year that they teach them, their lasting effect on the next year's score, and so on. Following the notation of Mariano et al. (2010), let $\theta_{g[jt]}$ represent the effect for the j -th grade- g teacher on a student's grade t score, for $t \geq g$. A grade g teacher has $K_g = T - g + 1$ effects, one each for grades g, \dots, T . Thus $\theta_{g[j\cdot]}$ gives the vector of current and future year effects of the j -th grade g teacher. The vector η concatenates the $\theta_{g[j\cdot]}$ effects for all grades and teachers. The model is able to distinguish between the persistence effect of former teachers and the current effect of the present teacher because the students are not nested in teachers.

We structure $\boldsymbol{\eta}$ so that \mathbf{G} will be block diagonal: if

$$(4) \quad \boldsymbol{\eta} = (\boldsymbol{\theta}'_{1[1\cdot]}, \dots, \boldsymbol{\theta}'_{1[m_1\cdot]}, \boldsymbol{\theta}'_{2[1\cdot]}, \dots, \boldsymbol{\theta}'_{2[m_2\cdot]}, \dots, \boldsymbol{\theta}'_{T[1\cdot]}, \dots, \boldsymbol{\theta}'_{T[m_T\cdot]})'$$

then $\boldsymbol{\eta} \sim N(\mathbf{0}, \mathbf{G})$, where

$$(5) \quad \mathbf{G} = \text{blockdiag}(\boldsymbol{\Gamma}_1, \dots, \boldsymbol{\Gamma}_1, \dots, \boldsymbol{\Gamma}_T, \dots, \boldsymbol{\Gamma}_T).$$

With m_g teachers in year g , there are m_g copies each of $\boldsymbol{\Gamma}_g$ and each $\boldsymbol{\Gamma}_g$ is unstructured. The matrix $\boldsymbol{\Gamma}_g$ is square with K_g rows and gives the covariance of current and future year effects for teachers of grade g . The vector \mathbf{s}_{ig} contains 1's in entries corresponding to teachers who could affect response g of student i . Thus, for a measurement y_{i2} , where student i had teacher 5 at time 1 and teacher 12 at time 2, \mathbf{s}_{i2} contains a 1 corresponding to the position of $\boldsymbol{\theta}_{1[5,2]}$ to include the lagged-year effect of teacher 5, and it contains a 1 corresponding to the position of $\boldsymbol{\theta}_{2[12,2]}$ to include the current-year effect of teacher 12. If, on average, teachers have more effect on current-year student scores than on subsequent scores of their students, we expect the first diagonal element of $\boldsymbol{\Gamma}_g$ to be larger than the other diagonal elements, reflecting the larger variability of current-year teacher effects.

The intra-student correlation is modeled in unstructured blocks of the conditional covariance matrix \mathbf{R} . After ordering the data by student and then by year, the error terms $\boldsymbol{\epsilon} = (\boldsymbol{\epsilon}'_1, \dots, \boldsymbol{\epsilon}'_n)'$ are distributed as $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{R})$ where \mathbf{R} is a block diagonal matrix with blocks

$$(6) \quad \mathbf{R}_i = \begin{pmatrix} \sigma_{11} & \cdots & \sigma_{T1} \\ \vdots & \ddots & \vdots \\ \sigma_{T1} & \cdots & \sigma_{TT} \end{pmatrix}.$$

If student i is missing an observation, then \mathbf{R}_i omits the corresponding row and column corresponding to the year in which the observation is missing. \mathbf{R}_i depends on i only through the dimension. We refer to this model as GP.R, indicating that the intra-student correlation is modeled in the \mathbf{R} matrix.

An advantage of this model is that the responses in different years can use different scales—the scaling is picked up in the covariance matrices \mathbf{G} and \mathbf{R} . The model has great flexibility for the relation between current- and future-year teacher effects, and for the within-student correlation. Note that this formulation assumes that the sets of teachers in different grades are distinct (or that if someone teaches in both grades 3 and 4, their effects on the grade-3 and grade-4 students are independent). The model can be modified to allow additional dependence for persons who teach multiple grades, but for simplicity here we consider the case with distinct teachers.

2.2. GP Model with a Single Future Year Effect. Some processes, including the educational data analyzed by Mariano et al. (2010), produce strongly correlated future year effects. Mariano et al. (2010) note in their application that, within each grade, the future year effects are strongly correlated with each other, but only moderately correlated to the current year effect. Following their idea of averaging the future year effects of each teacher after fitting the full model, we fit a reduced model that allocates a single future year effect to each teacher. This combines aspects of the GP model and the complete persistence model. This reduction requires that the scale of measurement be the same for each year of the study. We refer to the reduced model as rGP.R. We use $\theta_{t[j1]}$ and $\theta_{t[j2]}$ to represent the current and future

year effects, respectively, of the j -th grade- t teacher in rGP.R. Then \mathbf{s}_{ig} contains a 1 in the entry corresponding to the teacher in year g , $\theta_{g[j1]}$, and also contains a 1 in each entry corresponding to the future effects of the student's teachers in years $1, \dots, g-1$. An alternative model would be to impose an autoregressive structure on the Γ_i as in Paddock et al. (2011).

2.3. GP Model with Random Student Effects. When scores from each year are measured on the same scale, an alternative model specification is available. Using a variable persistence structure for teacher effects, McCaffrey and Lockwood (2011) modeled the intra-student correlation by using a random intercept for each student. We implement this alternative structure here, except we use the generalized persistence structure for teacher effects. We refer to this model as GP.G:

$$(7) \quad y_{ig} = \mathbf{x}'_{ig}\boldsymbol{\beta} + \mathbf{s}'_{ig}\boldsymbol{\eta}_* + \delta_i + \epsilon_{ig}$$

The terms in Equation (7) are defined the same as they were in Equation (1), with the exception of ϵ_{ig} and the new term δ_i . For this subsection, we use $\boldsymbol{\eta}_*$ to denote the vector of teacher effects. Instead of modeling $\boldsymbol{\epsilon}_i$ with an unstructured covariance matrix, GP.G includes a separate error variance in each year $\epsilon_{ig} \sim N(0, \sigma_g^2)$. As a result, \mathbf{R} is diagonal with entries from the set $\{\sigma_1^2, \dots, \sigma_T^2\}$, corresponding to the year of the observation. We likewise offer new definitions for \mathbf{G} , \mathbf{S} and $\boldsymbol{\eta}$ for GP.G.

The δ_i are random student intercepts, distributed as $\delta_i \sim N_1(0, \Gamma_{stu})$, with $\text{cov}(\epsilon_{ig}, \delta_i) = 0$. We may express GP.G in the form of Equation (1) by including the δ_i in the random effects vector $\boldsymbol{\eta}$,

$$(8) \quad \boldsymbol{\eta} = (\delta_1, \dots, \delta_n, \boldsymbol{\eta}'_*)'.$$

The vector $\boldsymbol{\eta}$ is then distributed as $\boldsymbol{\eta} \sim N(\mathbf{0}, \mathbf{G})$ where

$$(9) \quad \mathbf{G} = \text{blockdiag}(\Gamma_{stu}\mathbf{I}_n, \Gamma_1, \dots, \Gamma_1, \dots, \Gamma_T, \dots, \Gamma_T),$$

with m_g copies each of Γ_g , where each Γ_g is unstructured. To accommodate the new $\boldsymbol{\eta}$, the design matrix \mathbf{S} is composed of the blocks $[\mathbf{S}_1|\mathbf{S}_2]$, where \mathbf{S}_1 is the design matrix for the student effects and \mathbf{S}_2 is the design matrix for the teacher effects.

The same model could be fit without student random intercepts by modeling \mathbf{R} in Equation (6) as a compound-symmetric, block-diagonal matrix. However, the student-intercept formulation is useful when exploring sensitivity to the presence of potentially nonignorable missing data (McCaffrey and Lockwood, 2011; Karl et al., 2011) or when the random student intercepts are of interest. The GP.G formulation is also more easily extended to allow a random growth model where each student has his or her own slope and intercept.

2.4. Complete and Variable Persistence Models. Instead of modeling a separate effect in years g, \dots, T for each grade- g teacher, the variable persistence (VP) VAM models a single effect for each teacher. Let $\theta_{t[j]}$ denote the effect of the j -th grade- t teacher. The persistent effect of the j -th grade- t teacher on grade- g scores is modeled as a multiple of that teacher's effect, $\alpha_{gt}\theta_{t[j]}$. Lockwood et al. (2007) refer to the α_{gt} for $g = 1, \dots, T$ and $t = 1, \dots, g$ as persistence parameters. The persistence parameters for the current year are fixed at one, $\alpha_{gt} = 1$ for $t = g$, while the others are estimated. The complete and zero persistence VAMs are two special cases of the VP model, with fixed persistence parameters $\alpha_{gt} = 1$ and $\alpha_{gt} = 0$ (for $t \neq g$), respectively. The \mathbf{R} matrix of VP is the same as the one defined for

GP.R. SAS does not provide the ability to estimate the VP model: Lockwood et al. (2007) note that there are no available scalable implementations of the VP model. However, we will show that the EM algorithm can provide a scalable routine for the VP model.

The random teacher effects for the VP model are concatenated

$$\boldsymbol{\eta} = (\theta_{1[1]}, \dots, \theta_{1[m_1]}, \theta_{2[1]}, \dots, \theta_{2[m_2]}, \dots, \theta_{T[1]}, \dots, \theta_{T[m_T]}).$$

and distributed as $N(0, \mathbf{G})$. Since there is only one effect modeled for each teacher, \mathbf{G} is diagonal with m_g copies of Γ_g for $g = 1, \dots, T$,

$$\mathbf{G} = \text{diag}(\Gamma_1, \dots, \Gamma_1, \Gamma_2, \dots, \Gamma_2, \dots, \Gamma_T, \dots, \Gamma_T),$$

In the VP and CP models, the α_{tg} 's modify the covariance structure, since the effect of teacher j in year g , $\theta_{g[j]}$, appears in the model for that teacher's students in all subsequent years. To match the structure of the VP model to that of the GP model, let $\boldsymbol{\eta}_* = \mathbf{A}\boldsymbol{\eta}$ where \mathbf{A} has $\sum_g m_g$ columns, one for each teacher. $\mathbf{A} = \text{blockdiag}(\mathbf{A}_1, \dots, \mathbf{A}_T)$, where $\mathbf{A}_g = \mathbf{I}_{m_g} \otimes (\alpha_{1g} = 1, \dots, \alpha_{Tg})'$, for $g = 1, \dots, T$, with \otimes representing the Kronecker product. Using the same definition of \mathbf{S} as in Section 2.1, the VP model may be expressed as

$$(10) \quad \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{S}\boldsymbol{\eta}_* + \boldsymbol{\epsilon},$$

where $\boldsymbol{\eta}_* \sim N(\mathbf{0}, \mathbf{G}_*)$, with $\mathbf{G}_* = \mathbf{A}\mathbf{G}\mathbf{A}'$. The error terms are distributed as $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{R})$, $\text{cov}(\boldsymbol{\eta}, \boldsymbol{\epsilon}) = 0$, where \mathbf{R} is defined in Section 2.1. Briggs and Weeks (2011) discuss alternative formulations of the VP model that can be used if there are concerns about identifiability.

3. COMPUTING MAXIMUM LIKELIHOOD ESTIMATES

The degree of computational difficulty associated with estimating the parameters of Model (2) depends largely on the structure of the random effects, manifested through the pattern of nonzero entries in \mathbf{S} . In applications where the random effects are nested within subjects, the resulting \mathbf{V} matrix is block diagonal, and the log-likelihood in (3) may be factored over the subjects. However, for non-nested models, \mathbf{V} has no patterned structure, and its dimension is equal to the number of observations in the data set. As a result, a direct maximization of the likelihood function is highly inefficient or infeasible for large data sets. Wolfinger et al. (1994) develop a dimensionality-reduction technique—used with a Newton-Raphson (NR) routine in SAS[®] software (SAS Institute Inc., 2013)—that requires the manipulation of a square matrix with dimension depending on the number of levels of fixed and random effects, rather than the number of observations. Either the method of Wolfinger et al. (1994) or some other form of dimensionality reduction is necessary for scalable estimation of Model (2) when the random effects are not nested.

Even after a dimensionality reduction for \mathbf{V} , the matrices \mathbf{R} and \mathbf{S} grow with the size of the data set. For example, SAS PROC GLIMMIX allows users to specify a custom \mathbf{S} matrix via the multimember option of its EFFECT statement. However, the procedure does not currently take into account the sparse structure of the design and covariance matrices, and does not scale well to large data sets. SAS PROC HP MIXED does take sparseness into account and can be used to estimate a variation of the complete persistence model. However, HP MIXED is tailored to

a specific model and has limited choice of covariance structures. Broatch and Lohr (2012) show how to use SAS software to estimate parameters in a multiresponse VAM by specifying a user-defined covariance matrix, but this method also does not work well with large data sets.

Model (2) requires a positive definite \mathbf{G} matrix. A third major difficulty associated with the estimation of linear mixed models arises when random effects are highly correlated, producing a nearly singular \mathbf{G} matrix. The Newton-Raphson routines are prone to failure in these settings, frequently producing non-positive definite estimates for \mathbf{G} (Demidenko, 2004). One possible solution to this issue is by parameterizing the model according to the Cholesky root of \mathbf{G} . SAS offers functionality for such a parametrization, but it is only compatible with banded-unstructured covariance matrices (SAS Institute Inc., 2013).

The EM algorithm presented below overcomes these challenges by using a matrix of reduced dimension from that of \mathbf{V} , utilizing the sparseness of \mathbf{S} , \mathbf{G} , and \mathbf{R} , and achieving stability when the random effects covariance matrix is nearly singular.

3.1. The EM Algorithm. The EM algorithm (Dempster et al., 1977; McLachlan and Krishnan, 2008) provides a broad framework for maximum likelihood estimation in the presence of missing data. It was one of the first methods used to estimate linear mixed models by treating latent random effects as missing data (Laird and Ware, 1982). Its use for estimation of mixed models has lagged behind the popularity of the often-faster NR algorithms. The EM algorithm has a linear rate of convergence (which depends on the number and structure of the random effects) near a local maximum (Dempster et al., 1977), whereas the NR algorithms provide a quadratic rate of convergence.

Nevertheless, an advantage of the EM algorithm is that no restrictions need to be placed on the \mathbf{G} matrix to ensure that it is positive definite, as shown in the Appendix. For some models, the usual advantages of NR over EM (Lindstrom and Bates, 1988) are negated by the presence of highly correlated random effects. Furthermore, the EM algorithm naturally depends on the manipulation of matrices of dimension equal to the number of random effects rather than the number of observations so that additional dimensionality reduction techniques are not necessary. When taking advantage of sparse matrix computations, the EM algorithm can provide a viable method for estimating non-nested mixed models, especially those with highly correlated random effects.

We will refer to $f(\mathbf{y}; \Psi)$ as the observed data density function and $f(\mathbf{y}, \boldsymbol{\eta}; \Psi) = f(\mathbf{y}|\boldsymbol{\eta}; \Psi)f(\boldsymbol{\eta}; \Psi)$ as the complete data density function, where

$$\begin{aligned} f(\mathbf{y}|\boldsymbol{\eta}; \Psi) &\propto |\mathbf{R}|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{S}\boldsymbol{\eta})' \mathbf{R}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{S}\boldsymbol{\eta}) \right\} \\ f(\boldsymbol{\eta}; \Psi) &\propto |\mathbf{G}|^{-1/2} \exp \left\{ -\frac{1}{2} \boldsymbol{\eta}' \mathbf{G}^{-1} \boldsymbol{\eta} \right\} \end{aligned}$$

Given initial values for the parameters and the random effects, the EM algorithm alternates between an expectation (E) step and a maximization (M) step. At iteration $(k + 1)$, the E step calculates the conditional expectation of the complete data log-likelihood, given the observed data, \mathbf{y} , and parameter estimates obtained in the k -th step, $\Psi^{(k)}$. That is, the E step computes

$$Q(\Psi; \Psi^{(k)}) = \int \{ \log f(\mathbf{y}|\boldsymbol{\eta}; \Psi) + \log f(\boldsymbol{\eta}; \Psi) \} f(\boldsymbol{\eta}|\mathbf{y}; \Psi^{(k)}) d\boldsymbol{\eta}.$$

The M step then maximizes $Q(\Psi; \Psi^{(k)})$ with respect to Ψ , resulting in the updated parameter vector $\Psi^{(k+1)}$ satisfying

$$(11) \quad \int \frac{\partial}{\partial \Psi} \{ \log f(\mathbf{y}|\boldsymbol{\eta}; \Psi) + \log(f(\boldsymbol{\eta}; \Psi)) \} f(\boldsymbol{\eta}|\mathbf{y}; \Psi^{(k)}) d\boldsymbol{\eta} \Big|_{\Psi=\Psi^{(k+1)}} = \mathbf{0},$$

provided that differentiation and integration are interchangeable, which is valid because the complete data likelihood $f(\mathbf{y}, \boldsymbol{\eta}; \Psi)$ is a member of the exponential family (Lehmann and Romano, 2010). Note that the expression on the left side of Equation (11) is equivalent to the observed data score vector $S(\Psi; \mathbf{y}) = (\partial/\partial \Psi) l(\Psi; \mathbf{y})$ (Louis, 1982).

In Sections 3.2 to 3.4 we derive the M step for each of the models developed in Section 2. The E step described in Section 3.5 is the same for all of the models discussed, using the appropriate definitions of $\boldsymbol{\eta}$, \mathbf{G} , \mathbf{S} , and \mathbf{R} .

3.2. M-Step for GP.R and rGP.R. The M-step updates appearing in this section apply to both the generalized persistence model GP.R and its reduced version rGP.R. The only differences that must be kept in mind are the definitions of $\boldsymbol{\eta}$ and \mathbf{G}_g . Using the definition of \mathbf{G} in Equation (5), which applies to both GP.R and rGP.R, we may write the density of $\boldsymbol{\eta}$ as

$$\begin{aligned} f(\boldsymbol{\eta}; \Psi) &\propto \det(\mathbf{G})^{-1/2} \exp\left(-\frac{\boldsymbol{\eta}' \mathbf{G}^{-1} \boldsymbol{\eta}}{2}\right) \\ &= \left[\prod_{g=1}^T \det(\mathbf{G}_g)^{-m_g/2} \right] \exp\left(-\sum_{g=1}^T \sum_{j=1}^{m_g} \frac{\boldsymbol{\theta}'_{g[j\cdot]} \mathbf{G}_g^{-1} \boldsymbol{\theta}_{g[j\cdot]}}{2}\right) \end{aligned}$$

We use Petersen and Pedersen (2008) and Harville (2008) for matrix differentiation, and note that each \mathbf{G}_g is symmetric. Referring to Equation (11), the score vector with respect to \mathbf{G}_g is

$$\begin{aligned} S(\mathbf{G}_g) &= \int \frac{\partial}{\partial \mathbf{G}_g} \log \left[\det(\mathbf{G})^{-1/2} \exp\left(-\frac{\boldsymbol{\eta}' \mathbf{G}^{-1} \boldsymbol{\eta}}{2}\right) \right] f(\boldsymbol{\eta}|\mathbf{y}; \Psi) d\boldsymbol{\eta} \\ &= -\frac{1}{2} \int \frac{\partial}{\partial \mathbf{G}_g} \left\{ m_g \log [\det(\mathbf{G}_g)] + \sum_{j=1}^{m_g} \boldsymbol{\theta}'_{g[j\cdot]} \mathbf{G}_g^{-1} \boldsymbol{\theta}_{g[j\cdot]} \right\} f(\boldsymbol{\eta}|\mathbf{y}; \Psi) d\boldsymbol{\eta} \\ &= \text{matrix with components } \begin{cases} d_{ij} & \text{if } i = j \\ 2d_{ij} & \text{if } i \neq j \end{cases} \end{aligned}$$

where d_{ij} is the ij -th component of the matrix

$$\mathbf{D} = -\frac{1}{2} \left\{ m_g \mathbf{G}_g^{-1} - \mathbf{G}_g^{-1} \left(\sum_{j=1}^{m_g} \mathbf{E} \left[\boldsymbol{\theta}_{g[j\cdot]} \boldsymbol{\theta}'_{g[j\cdot]} | \mathbf{y}; \Psi \right] \right) \mathbf{G}_g^{-1} \right\}$$

Let

$$(12) \quad \tilde{\boldsymbol{\eta}} = \mathbf{E}[\boldsymbol{\eta}|\mathbf{y}; \Psi]$$

$$(13) \quad \tilde{\mathbf{v}} = \text{var}[\boldsymbol{\eta}|\mathbf{y}; \Psi]$$

represent the conditional expectation and variance, respectively, of $\boldsymbol{\eta}$. These quantities are calculated in the E-step and remain fixed during the M-step. Likewise, let the sub-vector of $\tilde{\boldsymbol{\eta}}$ corresponding to $\mathbf{E}[\boldsymbol{\theta}_{g[j\cdot]}|\mathbf{y}; \Psi]$ be denoted $\tilde{\boldsymbol{\theta}}_{g[j\cdot]}$, and the block

of the matrix $\tilde{\mathbf{v}}$ corresponding to $E[\boldsymbol{\theta}_{g[j\cdot]}\boldsymbol{\theta}'_{g[j\cdot]}|\mathbf{y};\boldsymbol{\Psi}]$ be denoted $\tilde{\mathbf{v}}_{g[j\cdot]}$. Now, since $\tilde{\mathbf{v}} = E[\boldsymbol{\eta}\boldsymbol{\eta}'|\mathbf{y};\boldsymbol{\Psi}] - \tilde{\boldsymbol{\eta}}\tilde{\boldsymbol{\eta}}'$, setting $S(\boldsymbol{\Gamma}_g) = \mathbf{0}$ implies

$$m_g \boldsymbol{\Gamma}_g^{-1} = \boldsymbol{\Gamma}_g^{-1} \sum_{j=1}^{m_g} \left(\tilde{\mathbf{v}}_{g[j\cdot]} + \tilde{\boldsymbol{\theta}}_{g[j\cdot]} \tilde{\boldsymbol{\theta}}'_{g[j\cdot]} \right) \boldsymbol{\Gamma}_g^{-1}$$

Thus the M-step update for $\boldsymbol{\Gamma}_g$ is

$$(14) \quad \hat{\boldsymbol{\Gamma}}_g = \frac{1}{m_g} \sum_{j=1}^{m_g} \left(\tilde{\mathbf{v}}_{g[j\cdot]} + \tilde{\boldsymbol{\theta}}_{g[j\cdot]} \tilde{\boldsymbol{\theta}}'_{g[j\cdot]} \right)$$

Equation (14) calculates an average of the blocks of $\tilde{\mathbf{v}} + \tilde{\boldsymbol{\eta}}\tilde{\boldsymbol{\eta}}'$ that correspond to teachers who taught in year g .

The M-step update for $\boldsymbol{\beta}$ is the value that solves $S(\boldsymbol{\beta}) = 0$, where

$$\begin{aligned} S(\boldsymbol{\beta}) &= \int \frac{\partial}{\partial \boldsymbol{\beta}} \left[-\frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{S}\boldsymbol{\eta})' \mathbf{R}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{S}\boldsymbol{\eta}) \right] f(\boldsymbol{\eta}|\mathbf{y}; \boldsymbol{\Psi}) d\boldsymbol{\eta} \\ &= \mathbf{X}' \mathbf{R}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{S}\tilde{\boldsymbol{\eta}}), \end{aligned}$$

namely,

$$(15) \quad \hat{\boldsymbol{\beta}} = (\mathbf{X}' \mathbf{R}^{-1} \mathbf{X})^{-1} \mathbf{X}' \mathbf{R}^{-1} (\mathbf{y} - \mathbf{S}\tilde{\boldsymbol{\eta}})$$

The calculation of the M-step update for \mathbf{R} from Equation (6) is complicated by the fact that the structure of \mathbf{R} changes in the presence of unbalanced data. The M-step update for the component σ_{kl} of \mathbf{R} is the value that solves $S(\sigma_{kl}) = 0$, where

$$\begin{aligned} S(\sigma_{kl}) &= \int \frac{\partial}{\partial \sigma_{kl}} \left[\log(|\mathbf{R}|^{-1/2}) \right. \\ &\quad \left. - \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{S}\boldsymbol{\eta})' \mathbf{R}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{S}\boldsymbol{\eta}) \right] f(\boldsymbol{\eta}|\mathbf{y}; \boldsymbol{\Psi}) d\boldsymbol{\eta}. \end{aligned}$$

If the observations are sorted by students and then by year, \mathbf{R} is block-diagonal with block sizes depending on the number of observations on each student. For T years, there are $2^T - 1$ possible combinations of years in which a student may be observed, although not all of these patterns may appear in a given data set. To parameterize these combinations, we treat the ordered, binary observed-test-score (OTS) indicators for each student as a number in base-2. So in a study over three years, each student will have an OTS pattern from the first column of Table 1.

For example, a student with observations in each year has pattern 7, with the corresponding block of \mathbf{R} given by

$$\begin{pmatrix} \sigma_{11} & \sigma_{21} & \sigma_{31} \\ \sigma_{21} & \sigma_{22} & \sigma_{32} \\ \sigma_{31} & \sigma_{32} & \sigma_{33} \end{pmatrix}.$$

The matrices corresponding to the other patterns are subsets of this matrix, using the rows and columns suggested by the OTS indicator. A student who is missing an observation in year 2 has pattern 5 and corresponding error covariance matrix

$$\begin{pmatrix} \sigma_{11} & \sigma_{31} \\ \sigma_{31} & \sigma_{33} \end{pmatrix}.$$

TABLE 1. Parameterizing the OTS patterns for example with 3 years

OTS	
indicators	Pattern
001	1
010	2
011	3
100	4
101	5
110	6
111	7

Let p denote the OTS pattern, n_p be the number of students with that pattern, and $\mathbf{R}_{(p)}$ represent the covariance matrix corresponding to the p -th pattern. In addition, let P_{kl} denote the set of patterns p whose covariance matrix $\mathbf{R}_{(p)}$ contains σ_{kl} . Furthermore, let $b(p)$ denote the b -th student with pattern p . We may write

$$|\mathbf{R}| = \prod_p |\mathbf{R}_{(p)}|^{n_p}.$$

Thus the score function may be expressed as

$$S(\sigma_{kl}) = -\frac{1}{2} \int \frac{\partial}{\partial \sigma_{kl}} \left\{ \sum_p n_p \log |\mathbf{R}_{(p)}| + \sum_p \sum_b \left[\left(\mathbf{y}_{b(p)} - \mathbf{X}_{b(p)}\boldsymbol{\beta} - \mathbf{S}_{b(p)}\boldsymbol{\eta} \right)' \mathbf{R}_{(p)}^{-1} \left(\mathbf{y}_{b(p)} - \mathbf{X}_{b(p)}\boldsymbol{\beta} - \mathbf{S}_{b(p)}\boldsymbol{\eta} \right) \right] \right\} \\ \times f(\boldsymbol{\eta}|\mathbf{y}; \boldsymbol{\Psi}) d\boldsymbol{\eta}$$

where $\mathbf{y}_{b(p)}$ is the vector of observations from student $b(p)$, with corresponding design matrices for fixed and random effects $\mathbf{X}_{b(p)}$ and $\mathbf{S}_{b(p)}$. The derivative will be 0 for all terms that do not contain the parameter σ_{kl} . This includes observations on students who do not have observations in both years k and l . Then, taking the derivative and letting $1_{\{C\}}$ be the indicator function that takes the value 1 if condition C is true and 0 otherwise,

$$S(\sigma_{kl}) = - \left(1_{\{i \neq j\}} + \frac{1}{2} \times 1_{\{i=j\}} \right) \sum_{p \in P_{kl}} \left\{ n_p \left(\mathbf{R}_{(p)}^{-1} \right)_{\{kl\}} - \right. \\ \left. \int \sum_b \left[\mathbf{R}_{(p)}^{-1} \left(\mathbf{y}_{b(p)} - \mathbf{X}_{b(p)}\boldsymbol{\beta} - \mathbf{S}_{b(p)}\boldsymbol{\eta} \right) \right. \right. \\ \left. \left. \times \left(\mathbf{y}_{b(p)} - \mathbf{X}_{b(p)}\boldsymbol{\beta} - \mathbf{S}_{b(p)}\boldsymbol{\eta} \right)' \mathbf{R}_{(p)}^{-1} \right]_{\{kl\}} f(\boldsymbol{\eta}|\mathbf{y}; \boldsymbol{\Psi}) d\boldsymbol{\eta} \right\}.$$

The notation $\{kl\}$ indicates the matrix component corresponding to the position of the parameter σ_{kl} in $R_{(p)}$. Again using the relationship $\tilde{\mathbf{v}} = E[\boldsymbol{\eta}\boldsymbol{\eta}'|\mathbf{y}; \boldsymbol{\Psi}] - \tilde{\boldsymbol{\eta}}\tilde{\boldsymbol{\eta}}'$,

$$\begin{aligned}
 S(\sigma_{kl}) = & - \left(1_{\{k \neq l\}} + \frac{1}{2} \times 1_{\{k=l\}} \right) \sum_{p \in P_{kl}} \left\{ n_p \mathbf{R}_{(p)}^{-1} \right. \\
 & - \mathbf{R}_{(p)}^{-1} \sum_b \left[\left(\mathbf{y}_{b(p)} - \mathbf{X}_{b(p)} \boldsymbol{\beta} \right) \left(\mathbf{y}_{b(p)} - \mathbf{X}_{b(p)} \boldsymbol{\beta} \right)' \right. \\
 & - \left(\mathbf{y}_{b(p)} - \mathbf{X}_{b(p)} \boldsymbol{\beta} \right) \left(\mathbf{S}_{b(p)} \tilde{\boldsymbol{\eta}} \right)' - \mathbf{S}_{b(p)} \tilde{\boldsymbol{\eta}} \left(\mathbf{y}_{b(p)} - \mathbf{X}_{b(p)} \boldsymbol{\beta} \right)' \\
 & \left. \left. + \mathbf{S}_{b(p)} \left(\tilde{\mathbf{v}} + \tilde{\boldsymbol{\eta}}\tilde{\boldsymbol{\eta}}' \right) \mathbf{S}_{b(p)}' \right] \mathbf{R}_{(p)}^{-1} \right\}_{\{kl\}}.
 \end{aligned} \tag{16}$$

If there were no missing observations then there would only be one OTS pattern and the calculation of the M-step update for \mathbf{R} would have a solution that followed the same pattern as the M-step update for \mathbf{G} . However, the presence of unbalanced student profiles disrupts the structure of \mathbf{R} , and score functions must be calculated for each of the unique model parameters in \mathbf{R} . The closed form solution for $S(\sigma_{kl}) = 0$ depends on the number of years and on the OTS patterns that are present in the data set. One option is to use a Newton-Raphson routine to solve the score equations. We suggest such a method in Section 4.

3.3. M-step for GP.G. The M-step update for $\boldsymbol{\beta}$ in GP.G is the same as the update for GP.R appearing in Equation (15), given the appropriate definition of \mathbf{R} . Likewise, the M-step updates for the $\boldsymbol{\Gamma}_g$ appearing in Equation (14) are unchanged.

The new work required for GP.G in Equation (7) is the calculation of the M-step update for the student variance component Γ_{stu} and the yearly error variances σ_g^2 , for $g = 1, \dots, T$. The M-step update for Γ_{stu} is derived in the same way as the update for $\boldsymbol{\Gamma}_g$, and is equal to the mean of the first n diagonal elements of $\tilde{\mathbf{v}} + \tilde{\boldsymbol{\eta}}\tilde{\boldsymbol{\eta}}'$. For the purpose of calculating $\hat{\sigma}_g^2$, let B_g be the set of students that are observed in year g .

$$S(\sigma_g^2) = \int \frac{\partial}{\partial \sigma_g^2} \left[\log \left(\prod_{j=1}^T \prod_{i \in B_j} \sigma_j^{-1} \exp \left[-\frac{(y_{ij} - \mathbf{x}_{ij}' \boldsymbol{\beta} - \mathbf{s}_{ij}' \boldsymbol{\eta})^2}{2\sigma_j^2} \right] \right) \right] f(\boldsymbol{\eta}|\mathbf{y}; \boldsymbol{\Psi}) d\boldsymbol{\eta}$$

Setting the score function equal to 0 and then using the fact that

$$E[\boldsymbol{\eta}' \mathbf{s}_{ig} \mathbf{s}_{ig}' \boldsymbol{\eta} | \mathbf{y}; \boldsymbol{\Psi}] = \text{tr}(\mathbf{s}_{ig} \mathbf{s}_{ig}' \tilde{\mathbf{v}}) + \tilde{\boldsymbol{\eta}}' \mathbf{s}_{ig} \mathbf{s}_{ig}' \tilde{\boldsymbol{\eta}}$$

yields

$$(17) \quad \hat{\sigma}_g^2 = \frac{1}{n_g} \sum_{i \in B_g} \{ (y_{ig} - \mathbf{x}_{ig}' \boldsymbol{\beta}) (y_{ig} - \mathbf{x}_{ig}' \boldsymbol{\beta} - 2\mathbf{s}_{ig}' \tilde{\boldsymbol{\eta}}) + \mathbf{s}_{ig}' (\tilde{\mathbf{v}} + \tilde{\boldsymbol{\eta}}\tilde{\boldsymbol{\eta}}') \mathbf{s}_{ig} \}$$

3.4. M-step for VP and CP Models. Although the covariance structure of the VP model in Equation (10) is different, the parameters may be estimated in much the same way as for the GP model. The EM algorithm requires a positive definite covariance matrix for the random effects. Since \mathbf{G}_* in Equation (10) is singular, we work instead with the diagonal \mathbf{G} matrix defined in Section 2.4 and the associated vector $\boldsymbol{\eta}$ of current year teacher effects. This is done operationally

by forming $\mathbf{S}_* = \mathbf{S}\mathbf{A}$, so that the “design” matrix \mathbf{S}_* includes the parameters α_{gt} , and then iteratively updating \mathbf{S}_* as the parameter estimates are updated during the estimation procedure. This is merely an algebraic distinction, since $\mathbf{S}\mathbf{A}\boldsymbol{\eta} = \mathbf{S}\boldsymbol{\eta}_* = \mathbf{S}_*\boldsymbol{\eta}$, where $\boldsymbol{\eta}_*$ is the vector defined in Section 2.4.

The M-step updates for $\boldsymbol{\beta}$ and \mathbf{R} in the VP and CP models appear in Equations (15) and (16), given the appropriate definitions of $\boldsymbol{\eta}$ and \mathbf{S} . The estimates for Γ_g appearing in Equation (14) apply as well, except in this case the Γ_g are all scalars.

The VP model estimates the persistence parameters α_{uv} , whereas the CP model fixes them at 1. Let $\partial\mathbf{S}_*/\partial\alpha_{gt} = \boldsymbol{\Delta}^{gt}$. For example, when each student is linked to only one teacher in each year, $\boldsymbol{\Delta}^{gt}$ will be a sparse matrix with 1’s in rows corresponding to year- g observations, under columns corresponding to year- t teachers. The score function for α_{gt} in the VP model is

$$S(\alpha_{gt}) = (\mathbf{y}' - \boldsymbol{\beta}'\mathbf{X}')\mathbf{R}^{-1}\boldsymbol{\Delta}^{gt}\tilde{\boldsymbol{\eta}} - \text{tr}[\mathbf{S}'_*\mathbf{R}^{-1}\boldsymbol{\Delta}^{gt}(\tilde{\mathbf{v}} + \tilde{\boldsymbol{\eta}}\tilde{\boldsymbol{\eta}}')]$$

The score function is linear in $\boldsymbol{\alpha} = \{\alpha_{uv}\}_{g,t}$, meaning that a single Newton step provides an exact solution for $S(\boldsymbol{\alpha}) = \mathbf{0}$.

3.5. E-step. The E-step updates for all of the discussed models are identical, using the appropriate definitions of \mathbf{S} , \mathbf{G} , $\boldsymbol{\eta}$, and \mathbf{R} . Calculation of the components of observed data score vector requires the first two moments, $\tilde{\boldsymbol{\eta}}$ and $\tilde{\mathbf{v}}$, of $f(\boldsymbol{\eta}|\mathbf{y}; \boldsymbol{\Psi})$. Using the method of Henderson (1950, 1975), the moments are obtained from the gradient and Hessian of $f(\mathbf{y}, \boldsymbol{\eta})$ with respect to $\boldsymbol{\eta}$. The resulting estimates are

$$(18) \quad \tilde{\mathbf{v}} = (\mathbf{S}'\mathbf{R}^{-1}\mathbf{S} + \mathbf{G}^{-1})^{-1}$$

$$(19) \quad \tilde{\boldsymbol{\eta}} = \tilde{\mathbf{v}}\mathbf{S}'\mathbf{R}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

The expression for the EBLUP in Equation (19) is equivalent, via a matrix identity, to the perhaps more familiar expression

$$(20) \quad \tilde{\boldsymbol{\eta}} = \mathbf{G}\mathbf{S}'\mathbf{V}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

However, from a computational standpoint, (19) is much more efficient than (20) since it does not require calculation of the full marginal covariance matrix \mathbf{V} . The calculation of (19) is relatively fast despite the large dimension of \mathbf{R} because both \mathbf{S}' and \mathbf{R}^{-1} are sparse.

3.6. EM Standard Errors. One criticism of the EM algorithm is that it does not produce the Hessian of the MLE $\hat{\boldsymbol{\Psi}}$ as a byproduct. The work we have already done, however, makes it possible for us to compute the observed data information matrix directly without working through a correction to the complete-data information matrix, as done by Louis (1982). Equation (11) expresses the observed data score vector $S(\boldsymbol{\Psi})$ as the conditional expectation of the complete data likelihood. We derived the components of the observed data score vector in order to calculate the M-step equations. Together with the values $\tilde{\boldsymbol{\eta}}$ and $\tilde{\mathbf{v}}$ from the E-step, our expression for the score vector allows us to calculate the observed information matrix,

$$(21) \quad -\partial S(\boldsymbol{\Psi})/\partial\boldsymbol{\Psi}|_{\boldsymbol{\Psi}=\hat{\boldsymbol{\Psi}}}.$$

with a central difference approximation at the MLE $\hat{\boldsymbol{\Psi}}$. This method is suggested by Jamshidian and Jennrich (2000), who propose using either a forward or central difference approximation, or a Richardson extrapolation (Lindfield and Penny, 1988).

It is also useful to calculate standard errors for the predicted random effects. The matrix $\tilde{\mathbf{v}}$ provides the covariance matrix for $\boldsymbol{\eta}$; however, since $\boldsymbol{\eta}$ is random, $\tilde{\mathbf{v}}$ underestimates the prediction variance of $\tilde{\boldsymbol{\eta}} - \boldsymbol{\eta}$ (Littell et al., 2006). As demonstrated by McLean et al. (1991), the prediction variance matrix of the random effects appears in block \mathbf{C}_{22} of

$$\mathbf{C} = \begin{pmatrix} \mathbf{C}_{11} & \mathbf{C}_{12} \\ \mathbf{C}_{21} & \mathbf{C}_{22} \end{pmatrix} = \begin{pmatrix} \mathbf{X}'\mathbf{R}^{-1}\mathbf{X} & \mathbf{X}'\mathbf{R}^{-1}\mathbf{Z} \\ \mathbf{Z}'\mathbf{R}^{-1}\mathbf{X} & \mathbf{S}'\mathbf{R}^{-1}\mathbf{S} + \mathbf{G}^{-1} \end{pmatrix}^{-1}$$

This procedure also yields the standard errors for $\hat{\boldsymbol{\beta}}$. The standard errors obtained by this method for $\hat{\boldsymbol{\beta}}$ are the same as those obtained by the central difference approximation: the central difference approximation is needed only for the standard errors of the covariance parameters.

3.7. Convergence and Initial Values of the EM Algorithm. The EM algorithm converges to a stationary value of the observed data likelihood as long as $Q(\boldsymbol{\Psi}; \boldsymbol{\Psi}')$ is continuous in both $\boldsymbol{\Psi}$ and $\boldsymbol{\Psi}'$, and the parameter space is compact (Wu, 1983). Although the parameter space for $\boldsymbol{\Psi}$ is not compact for Model (2), this regularity condition can be satisfied by a truncation of the parameter space (McCulloch, 1994; Demidenko, 2004).

One possible convergence criterion is to stop the algorithm when the relative change in the log-likelihood at iteration k , $l(\boldsymbol{\Psi}^{(k)})$, is less than a fixed tolerance,

$$\frac{l(\boldsymbol{\Psi}^{(k)}) - l(\boldsymbol{\Psi}^{(k-1)})}{l(\boldsymbol{\Psi}^{(k)})} < w.$$

In general, we use $w = 10^{-7}$. Verification that the EM algorithm has converged to a local maximum of the likelihood function is possible by checking that the Hessian of the observed data likelihood is negative definite. As with any iterative maximization routine, there is no way to guarantee that the EM algorithm will converge to the global maximum of the likelihood, given a single set of initial values $\boldsymbol{\Psi}_0$. It is advisable to compare the results of the algorithm after starting from different sets of initial values. For the VAMs in Section 2, we did not find any sensitivity to the choice of $\boldsymbol{\Psi}_0$.

4. IMPLEMENTATION OF THE EM ALGORITHM

We have implemented estimation of the GP, VP and CP models in the R (R Core Team, 2013) package GPvam (Karl et al., 2012) for educational value-added assessments as an example of applying the proposed EM algorithm. Our program takes advantage of the sparseness of the design and certain covariance matrices, and handles large data sets relatively well. Because the program was custom-designed for these VAMs, it requires minimal input. The user must supply a data frame with columns for test scores, year of observation, student ID, and teacher ID. Optionally, other columns may be included for additional covariates in the \mathbf{X} matrix; these are declared to the program through an R `formula` statement. Sparse matrices are constructed and handled via the R package Matrix (Bates and Maechler, 2012).

The GP model requires specification of a complex random effects structure. Doran and Lockwood (2006) provide a tutorial to the implementation of VAMs in R using the functions `lme` and `lmer`. However, Lockwood et al. (2003) explain that, for a less complicated multi-membership model, data sets with more than

TABLE 2. Run times in minutes

Model	GPvam	SAS
CP	4.2	55
VP	5.5	N.A.
GP.G	12	Failed
rGP.R	80	Failed
GP.R	114	Failed

200 teachers require several tricks to program with `lme`, and often fail to converge. GPvam automatically builds the sparse design matrix for the random effects, and performs well in the application in Section 5 which contains 4781 teacher effects for GP.R.

Although GPvam has been tailored to the estimation of VAMs, the R code may be generalized to other applications involving linear mixed models. New code would need to be written to build the application-specific \mathbf{S} and \mathbf{G} matrices but iterative parts of the program would not need to be adjusted. The EM algorithm may also be extended for efficient estimation of non-nested, nonlinear mixed models. The use of a nonlinear link function will require an integral approximation in the E step. It would also be possible to impose structure on the \mathbf{R} matrix, such as autoregressive, compound symmetric, or Toeplitz. The procedures for obtaining the score functions of the parameters of an unstructured \mathbf{R} may serve as a template for these other situations.

As mentioned in Section 3.2, the M-step update for \mathbf{R} in GP.R, rGP.R, and VP requires extra computational work due to the lack of a readily-available closed form solution. We use a Newton-Raphson algorithm to calculate the M-step update. The standard NR algorithm for solving $S(\mathbf{R}) = \mathbf{0}$ often diverges when the initial \mathbf{R}_0 is too far from the maximum. To improve the stability of the routine, we modify the appropriate Hessian by adding a scaled diagonal matrix during the first few M-step updates for \mathbf{R} . This results in a hybrid of a Newton and a gradient descent method that produces more reliable convergence when the initial value is far away from the critical point (Nocedal and Wright, 1999).

To compare the performance of GPvam and SAS (with the EFFECT statement of PROC GLIMMIX) in implementing the models presented in Section 2, we consider a data set with 6236 observations on 2834 students over 3 years, with 102, 104, and 98 teachers in each year, respectively. Table 2 gives the results. GP.R and rGP.R each failed in SAS after encountering a negative-definite covariance matrix, while GP.G ran out of memory in SAS after a few minutes. The application in Section 5 involves a much larger data set than the one used in this example.

5. APPLICATION

We apply the models to the data set analyzed by Mariano et al. (2010), which is available in the supplementary material of McCaffrey and Lockwood (2011). According to McCaffrey and Lockwood (2011), the data come from vertically linked mathematics standardized test scores from grades 1–5 for a cohort of students from a large urban US school district.

FIGURE 1. Estimated \mathbf{G} and \mathbf{R} matrices from GP.R. The covariance matrix is on the left, and the correlation matrix is on the right.

$$\begin{aligned}
 &\mathbf{R}: \\
 &\begin{pmatrix} 0.741 & 0.478 & 0.463 & 0.456 & 0.392 \\ 0.478 & 0.705 & 0.523 & 0.516 & 0.449 \\ 0.463 & 0.523 & 0.736 & 0.563 & 0.484 \\ 0.456 & 0.516 & 0.563 & 0.688 & 0.509 \\ 0.392 & 0.449 & 0.484 & 0.509 & 0.565 \end{pmatrix} \begin{pmatrix} 1.000 & 0.661 & 0.626 & 0.639 & 0.606 \\ 0.661 & 1.000 & 0.726 & 0.740 & 0.711 \\ 0.626 & 0.726 & 1.000 & 0.791 & 0.750 \\ 0.639 & 0.740 & 0.791 & 1.000 & 0.817 \\ 0.606 & 0.711 & 0.750 & 0.817 & 1.000 \end{pmatrix} \\
 &\mathbf{\Gamma}_1: \\
 &\begin{pmatrix} 0.443 & 0.121 & 0.120 & 0.107 & 0.095 \\ 0.121 & 0.100 & 0.088 & 0.084 & 0.077 \\ 0.120 & 0.088 & 0.087 & 0.083 & 0.076 \\ 0.107 & 0.084 & 0.083 & 0.080 & 0.074 \\ 0.095 & 0.077 & 0.076 & 0.074 & 0.069 \end{pmatrix} \begin{pmatrix} 1.000 & 0.575 & 0.610 & 0.568 & 0.541 \\ 0.575 & 1.000 & 0.941 & 0.941 & 0.920 \\ 0.610 & 0.941 & 1.000 & 0.994 & 0.986 \\ 0.568 & 0.941 & 0.994 & 1.000 & 0.995 \\ 0.541 & 0.920 & 0.986 & 0.995 & 1.000 \end{pmatrix} \\
 &\mathbf{\Gamma}_2: \\
 &\begin{pmatrix} 0.281 & 0.059 & 0.039 & 0.042 \\ 0.059 & 0.025 & 0.023 & 0.020 \\ 0.039 & 0.023 & 0.024 & 0.020 \\ 0.042 & 0.020 & 0.020 & 0.017 \end{pmatrix} \begin{pmatrix} 1.000 & 0.703 & 0.478 & 0.593 \\ 0.703 & 1.000 & 0.951 & 0.967 \\ 0.478 & 0.951 & 1.000 & 0.979 \\ 0.593 & 0.967 & 0.979 & 1.000 \end{pmatrix} \\
 &\mathbf{\Gamma}_3: \\
 &\begin{pmatrix} 0.248 & 0.032 & 0.024 \\ 0.032 & 0.015 & 0.015 \\ 0.024 & 0.015 & 0.015 \end{pmatrix} \begin{pmatrix} 1.000 & 0.516 & 0.394 \\ 0.516 & 1.000 & 0.979 \\ 0.394 & 0.979 & 1.000 \end{pmatrix} \\
 &\mathbf{\Gamma}_4: \\
 &\begin{pmatrix} 0.130 & 0.038 \\ 0.038 & 0.030 \end{pmatrix} \begin{pmatrix} 1.000 & 0.612 \\ 0.612 & 1.000 \end{pmatrix} \\
 &\mathbf{\Gamma}_5 : 0.146
 \end{aligned}$$

The data have been pre-processed by McCaffrey and Lockwood (2011), and we further process the data by removing observations with no student link, as well as observations missing both the test score and the teacher link. The resulting data set consists of 26019 observations on 9295 students over 5 years. For grades 1 through 5, there are 338, 318, 306, 321, and 259 teachers, respectively. This results in a total of 4781 teacher effects for GP.R. The data set does not contain any additional covariates, so the fixed effects modeled include a mean for each year.

We fit each of the models GP.R, rGP.R, GP.G, VP, and CP to this data set using the program GPvam. Table 3 lists the estimated yearly means from GP.R: the results from the other models are similar. Figure 1 gives the maximum likelihood estimates of the covariance parameters from GP.R. Models rGP.R and GP.G are valid for this data set because the scores from each year are on the same scale of measurement. The estimates of current-year teacher effects are nearly identical for the three variations of the GP model, with correlations of 0.998 or higher among the estimated effects. The agreement between GP.R and rGP.R is not surprising

FIGURE 2. Standard errors for current year teacher ratings for year-3 teachers.

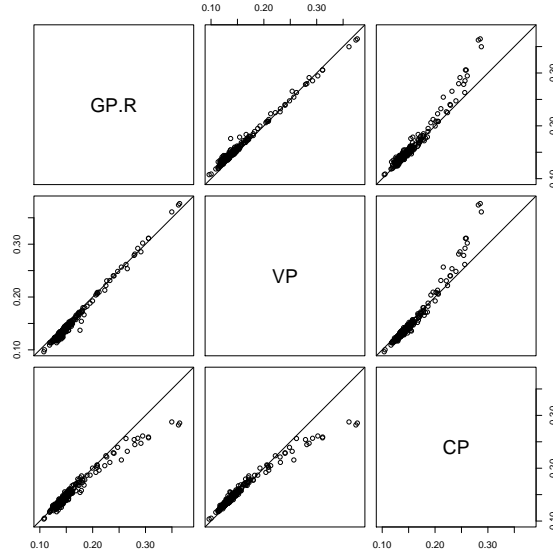


TABLE 3. Estimates for yearly means from GP.R

	Estimate	Std. Error
Year 1	3.395	0.030
Year 2	3.996	0.029
Year 3	4.726	0.023
Year 4	5.309	0.022
Year 5	5.984	0.025

TABLE 4. Persistence Parameters from VP

	Estimate	S.E.
α_{21}	0.18	0.02
α_{31}	0.19	0.02
α_{41}	0.17	0.02
α_{51}	0.15	0.02
α_{32}	0.22	0.02
α_{42}	0.14	0.02
α_{52}	0.16	0.02
α_{43}	0.13	0.02
α_{53}	0.09	0.02
α_{54}	0.29	0.03

given the extremely high correlations between future year effects seen in the $\mathbf{\Gamma}_g$ matrices of Figure 1. With a simplified covariance structure, rGP.R converges after 48 iterations in 7% of the time it takes the 431 iterations needed for GP.R to converge. For this data set, the reduced model rGP.R appears to provide a good alternative. Slow convergence in the neighborhood of a maximum that lies near the boundary of the parameter space is a well-known property of the EM algorithm (Demidenko, 2004). However, the time used by the EM algorithm is worthwhile, because the faster NR algorithms are prone to failure in these settings.

Fitting the CP and VP models with the EM algorithm results in approximately the same correlations among the GP.R, VP, and CP estimates as found by Lockwood et al. (2007) and Mariano et al. (2010). The persistence parameters for the VP model in Table 4 are similar to those provided by McCaffrey and Lockwood (2011), who modeled students with random effects. The persistence parameters are all significantly different from 1, indicating that the assumption of complete persistence is not compatible with this data set.

Using the EM algorithm, we obtain correlation patterns for GP.R that are similar to those in Figures 2 and 3 of Mariano et al. (2010). However, we note that Mariano et al. (2010) obtained these results after careful choice of an informative prior that allowed for strong correlations between future year effects. In simulation studies, they found that a minimally informative Wishart prior for covariance parameters could result in posterior credible intervals for the correlations that did not include the true values. The EM algorithm gives maximum likelihood estimates that do not need any specifications of prior distributions.

Figure 2 compares the standard errors associated with the predicted teacher effects for GP.R, VP, and CP. As stated in Section 3.6, the standard errors are calculated as a by-product of the EM algorithm. The values for the larger standard errors, which likely correspond to teachers with relatively fewer observations, are inflated when moving from the CP to the VP or GP models. This is interesting because the prediction intervals are used by some researchers to classify teachers as below-average, average, or above average (Draper, 1995; Lockwood et al., 2007). Despite the inflation of some of the standard errors seen when moving from the CP to the VP model, most of the prediction errors are smaller in the VP model. An advantage of maximum likelihood estimation is that the standard errors for the teacher effects, derived in Section 3.6, are free from the influence of potentially informative prior distributions.

The models fit to this data set have a number of assumptions that were stated in Section 2. All models assume that the teacher effects in $\boldsymbol{\eta}$ and the residual student effects in $\boldsymbol{\epsilon}$ are independent. If data are collected in a designed experiment, with students randomly assigned to teachers, this assumption is reasonable. Most data used in VAMs are observational, though, so this assumption is violated if, say, some teachers are regularly assigned the best students. In that case, the effects ascribed to teachers may actually be more properly attributed to the students who take those teachers. The models are also assumed to be “correct” in that they are assumed to include all factors relevant to the response. The data set analyzed here did not contain information on student-level covariates such as socioeconomic status, for example, and it is possible that including such covariates in the model would change the estimated teacher effects. These models further assume that missing test scores are missing at random. This amounts to assuming that the probability a test

score is missing does not depend on the student’s latent ability, teacher history, or what the student’s test score would have been if observed. Finally, although the multiresponse models considered here relate later test scores to earlier test scores through the within-student correlations, the models imply that, conditionally on $y_{i,g-1}$, the relationship between y_{ig} and $y_{i,g-1}$ is linear.

The usefulness of VAM scores for measuring teacher effectiveness depends on the quality of the tests as measures of student achievement (Koretz, 2008), and many aspects of teacher contributions may be unrelated by standardized tests (Braun et al., 2010). A relevant discussion of the utility and limitations of multilevel models is given by Draper (1995), who urges a careful examination of the nature of the sampling in the study. Thus, caution is needed when interpreting a teacher’s VAM score.

With these limitations in mind, VAMs can provide valuable information for improving the educational system (Harris and McCaffrey, 2010). In this example, the GP model indicates that the current-year teacher has the highest effect on student scores, and that the current-year effects are, on average, stronger in earlier grades than in later grades. The estimates of relative sources of variability also provide valuable information: in grade 1, current-year teachers (or other classroom effects that are associated with teachers) account for approximately 36% of the variability in student scores but that percentage drops to 16–20% in grades 4 and 5. Ballou et al. (2004) and Lockwood et al. (2007) note that teacher effects from the first year are most susceptible to bias resulting from nonrandom student assignment to classrooms.

6. CONCLUSIONS

In this article, we have developed a method for computing maximum likelihood estimates for a class of multiple membership models in which lower-level units progress through sequential higher-level units (Mariano et al., 2010). The EM algorithm offers an efficient method of computation, taking advantage of matrix sparsity and requiring inversion of a matrix whose dimension depends on the number of random effects, rather than on the total number of observations as in other implementations. The algorithm produces stable behavior even when the covariance matrix for the random effects is nearly singular. We have implemented the proposed methods in the R package GPvam. The availability of maximum-likelihood estimates should be useful for those preferring Bayesian estimation as well, providing a sensitivity analysis to their choice of priors. We hope that this user-friendly implementation of the model will facilitate further empirical study of the model’s properties.

In the educational context, the GP and other models fit provide a great deal of flexibility for studying relative effects of teachers on their students’ current and future achievement. In some cases the full flexibility of the GP model may be needed to summarize the data structure; in others, being able to examine the estimates from the GP model may show that a simpler structure adequately describes the data. Our algorithm readily calculates standard errors for the predicted teacher effects (VAM scores). In many applications, the standard errors of the teacher effects are quite large (Braun et al., 2010, p. 45), so that including the standard errors along with the point estimates can help distinguish “real” effects from random variation.

Although the computational methods and software were developed in the educational setting, they can be used in many other applications as well, substituting the lower-level units for “students” and the higher-level units for “teachers”. Similar models have been considered for studying the relative contributions of health care professionals or clinics to patient outcomes (Zaslavsky et al., 2004), and the models can be applied in any setting where different higher-level units sequentially affect the outcomes of lower-level units.

The algorithm presented here and the code in package GPvam may be extended to other multiple membership models. One extension would be to allow an expanded covariance structure in which the same higher-level unit may be associated with multiple responses. This would allow the model to better fit situations in which a student had the same teacher for more than one year, or in which a patient returned to a previous doctor. This extension would require careful bookkeeping to track which doctors are repeats, but the basic M-step and E-step of the algorithm would remain the same. Our models and applications used only intercepts for the random effects, but the implementation may be extended to include random slopes as well, and other more general multiple membership models described in Browne et al. (2001).

APPENDIX

We show that the EM algorithm produces positive definite (PD) \mathbf{G} matrix after each iteration. This assumes that the \mathbf{R} matrix is PD after each iteration, which is true for the GP.G model specification in Section 2.3, and can be demonstrated for GP.R from Section 2.1 in the absence of incomplete data where the M-step update for \mathbf{R} has an easily-obtainable solution. However, even in the presence of missing data, \mathbf{R} will usually not be near the boundary of the parameter space, since the error variances on the diagonal of \mathbf{R} will be positive as long as the model does not fit the data perfectly, and the intra-student effects are not likely to be perfectly correlated.

A much more common problem in estimating mixed models occurs when the estimated \mathbf{G} matrix is not PD (Verbeke and Molenberghs, 2000, Section 5.6.1). This may happen when 0 variance components are estimated, or when random effects are perfectly correlated, the later being a significant concern for the future teacher effects of the GP VAM. Notice that \mathbf{G} is a block-diagonal portion of $\tilde{\mathbf{v}} + \tilde{\boldsymbol{\eta}}\tilde{\boldsymbol{\eta}}'$. The matrix $\tilde{\mathbf{v}}$ is defined by Equation (18). Thus, $\tilde{\mathbf{v}}$ is PD as long as the initial \mathbf{G} is PD, because $\mathbf{S}'\mathbf{R}^{-1}\mathbf{S}$ is positive semi-definite. Furthermore, $\tilde{\boldsymbol{\eta}}\tilde{\boldsymbol{\eta}}'$ is positive semi-definite, and the sum of a PD and a positive semi-definite matrix is PD.

ACKNOWLEDGMENTS

This research was partially supported by the National Science Foundation under grant DRL-0909630. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author and do not reflect the views of the National Science Foundation or Arizona State University.

REFERENCES

- Airoldi, E. M., Blei, D. M., Fienberg, S. E., and Xing, E. P. (2008), “Mixed Membership Stochastic Blockmodels,” *Journal of Machine Learning Research*, 9, 1981–2014.

- Ash, A., Fienberg, S., Louis, T., Normand, S.-L., Stukel, T., and Utts, J. (2012), “Statistical Issues in Assessing Hospital Performance,” *COPPS-CMS White Paper*.
- Ballou, D., Sanders, W., and Wright, P. (2004), “Controlling for Student Background in Value-Added Assessment of Teachers,” *Journal of Educational and Behavioral Statistics*, 29, 37–65.
- Bates, D. and Maechler, M. (2012), *Matrix: Sparse and Dense Matrix Classes and Methods*, <http://cran.r-project.org/web/packages/Matrix/index.html>, R package version 1.0-4.
- Braun, H. I., Chudowsky, N., and Koenig, J. (2010), *Getting Value Out of Value-Added*, Washington, DC: National Academies Press.
- Briggs, D. C. and Weeks, J. P. (2011), “The Persistence of School-Level Value-Added,” *Journal of Educational and Behavioral Statistics*, 36, 616–637.
- Broatch, J. and Lohr, S. (2012), “Multidimensional Assessment of Value Added by Teachers to Real-World Outcomes,” *Journal of Educational and Behavioral Statistics*, 37, 256–277.
- Browne, W. J., Goldstein, H., and Rasbash, J. (2001), “Multiple Membership Multiple Classification (MMMC) Models,” *Statistical Modelling*, 1, 103–124.
- Demidenko, E. (2004), *Mixed Models Theory and Applications*, Hoboken: Wiley-Interscience.
- 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.
- Doran, H. C. and Lockwood, J. R. (2006), “Fitting Value-Added Models in R,” *Journal of Educational and Behavioral Statistics*, 31, 205–230.
- Draper, D. (1995), “Inference and Hierarchical Modeling in the Social Sciences,” *Journal of Educational and Behavioral Statistics*, 20, 115–117.
- Fielding, A. and Goldstein, H. (2006), *Cross-classified and Multiple Membership Structures in Multi-level Models: An Introduction and Review. Research Report No. 791*, Birmingham, UK: University of Birmingham, Department of Education and Skills.
- Goldstein, H., Rasbash, J., Browne, W., Woodhouse, G., and Poulain, M. (2000), “Multilevel Models in the Study of Dynamic Household Structures,” *European Journal of Population*, 16, 373–387.
- Harris, D. N. and McCaffrey, D. F. (2010), “Value-Added: Assessing Teachers’ Contributions to Student Achievement,” in *Teacher Assessment and the Quest for Teacher Quality*, ed. Kennedy, M. M., San Francisco: Jossey-Bass, pp. 251–282.
- Harville, D. A. (2008), *Matrix Algebra from a Statistician’s Perspective*, New York: Springer.
- Henderson, C. R. (1950), “The Estimation of Genetic Parameters,” *The Annals of Mathematical Statistics*, 21, 309–310.
- (1975), “Best Linear Unbiased Estimation and Prediction under a Selection Model,” *Biometrics*, 31, 423.
- Hill, P. and Goldstein, H. (1998), “Multilevel Modeling of Educational Data with Cross-Classification and Missing Identification for Units,” *Journal of Educational and Behavioral Statistics*, 23, 117–128.

- Jamshidian, M. and Jennrich, R. I. (2000), "Standard Errors for EM Estimation," *Journal of the Royal Statistical Society, Series B*, 62, 257–270.
- Karl, A., Yang, Y., and Lohr, S. (2011), "Exploring Missing Data in Value-Added Models in Education," in *JSM Proceedings*, Social Statistics Section, Alexandria, VA: American Statistical Association, 2449–2460.
- Karl, A. T., Yang, Y., and Lohr, S. (2012), *GPvam: Maximum Likelihood Estimation of Multiple Membership Mixed Models Used in Value-Added Modeling*, <http://cran.r-project.org/web/packages/GPvam/index.html>, R package version 2.0-0.
- Koretz, D. M. (2008), *Measuring Up: What Educational Testing Really Tells Us*, Cambridge, MA: Harvard University Press.
- Laird, N. M. and Ware, J. H. (1982), "Random-Effects Models for Longitudinal Data," *Biometrics*, 38, 963–974.
- Lehmann, E. L. and Romano, J. (2010), *Testing Statistical Hypotheses*, New York: Springer, 3rd ed.
- Lindfield, G. and Penny, J. E. T. (1988), *Microcomputers in Numerical Analysis*, New York: Halsted Press.
- Lindstrom, M. J. and Bates, D. M. (1988), "Newton-Raphson and EM Algorithms for Linear Mixed-Effects Models for Repeated-Measures Data," *Journal of the American Statistical Association*, 83, 1014–1022.
- Littell, R. C., Milliken, G. A., Stroup, W. W., Wolfinger, R. D., and Schabenberger, O. (2006), *SAS for Mixed Models*, Cary: SAS Institute, Inc., 2nd ed.
- Little, R. and Rubin, D. (2002), *Statistical Analysis with Missing Data*, New York: John Wiley, 2nd ed.
- Lockwood, J., McCaffrey, D., Mariano, L., and Setodji, C. (2007), "Bayesian Methods for Scalable Multivariate Value-Added Assessment," *Journal of Educational and Behavioral Statistics*, 32, 125–150.
- Lockwood, J. R., Doran, H., and McCaffrey, D. F. (2003), "Using R for Estimating Longitudinal Student Achievement Models," *The Newsletter of the R Project*, 3, 17–23.
- Lohr, S. (2012), "The Value Deming's Ideas Can Add to Educational Evaluation," *Statistics, Politics, and Policy*, 3(2), article 1.
- Louis, T. A. (1982), "Finding the Observed Information Matrix when Using the EM Algorithm," *Journal of the Royal Statistical Society Series B*, 44, 226–233.
- Mariano, L. T., McCaffrey, D. F., and Lockwood, J. (2010), "A Model for Teacher Effects from Longitudinal Data Without Assuming Vertical Scaling," *Journal of Educational and Behavioral Statistics*, 35, 253–279.
- McCaffrey, D., Lockwood, J., Mariano, L. T., and Setodji, C. (2005), "Challenges for Value-Added Assessment of Teacher Effects," in *Value Added Models in Education: Theory and Applications*, ed. Lissitz, R., Maple Grove, MN: JAM Press, pp. 111–144.
- McCaffrey, D., Lockwood, J. R., Koretz, D., Louis, T., and Hamilton, L. (2004), "Models for Value-Added Modeling of Teacher Effects," *Journal of Educational and Behavioral Statistics*, 29, 67–101.
- McCaffrey, D. F., Lockwood, J., Koretz, D. M., and Hamilton, L. S. (2003), *Evaluating Value-Added Models for Teacher Accountability*, Pittsburgh: The RAND Corporation.

- McCaffrey, D. F. and Lockwood, J. R. (2011), “Missing Data in Value-Added Modeling of Teacher Effects,” *Annals of Applied Statistics*, 5, 773–797.
- McCulloch, C. E. (1994), “Maximum Likelihood Variance Components Estimation for Binary Data,” *Journal of the American Statistical Association*, 89, 330–335.
- McLachlan, G. J. and Krishnan, T. (2008), *The EM Algorithm and Extensions*, Hoboken: John Wiley & Sons, 2nd ed.
- McLean, R. A., Sanders, W. L., and Stroup, W. W. (1991), “A Unified Approach to Mixed Linear Models,” *The American Statistician*, 45, 54–64.
- Nocedal, J. and Wright, S. J. (1999), *Numerical Optimization*, New York: Springer.
- Paddock, S. M., Hunter, S. B., Watkins, K. E., and McCaffrey, D. F. (2011), “Analysis of Rolling Group Therapy Data Using Conditionally Autoregressive Priors,” *Annals of Applied Statistics*, 5, 605–627.
- Petersen, K. B. and Pedersen, M. S. (2008), “The Matrix Cookbook,” Version 20081110.
- R Core Team (2013), *R: A Language and Environment for Statistical Computing*, R Foundation for Statistical Computing, Vienna, Austria.
- Raudenbush, S. and Bryk, A. (2002), *Hierarchical Linear Models: Applications and Data Analysis Methods*, Thousand Oaks, CA: Sage, 2nd ed.
- Rowan, B., Correnti, R., and Miller, R. J. (2002), “What Large-Scale, Survey Research Tells Us About Teacher Effects on Student Achievement: Insights From the Prospects Study of Elementary schools,” *Teachers College Record*, 104, 1525–1567.
- Sanders, W., Saxton, A., and Horn, B. (1997), “The Tennessee Value-Added Assessment System: A Quantitative Outcomes-Based Approach to Educational Assessment.” in *Grading Teachers, Grading Schools. Is Student Achievement a Valid Evaluation Measure?*, ed. Millman, J., Thousand Oaks, CA: Corwin Press, Inc, pp. 137–162.
- SAS Institute Inc. (2013), *SAS 9.3 Help and Documentation*, Cary: SAS Institute, Inc.
- Verbeke, G. and Molenberghs, G. (2000), *Linear Mixed Models for Longitudinal Data*, New York: Springer.
- Wolfinger, R., Tobias, R., and Sall, J. (1994), “Computing Gaussian Likelihoods and Their Derivatives for General Linear Mixed Models,” *SIAM Journal of Scientific Computing*, 15:6, 1294–1310.
- Wright, S. P., White, J. T., and Sanders, W. L. (2010), *SAS EVAAS Statistical Models*, Cary, NC: SAS Institute, www.sas.com/resources/asset/SAS-EVAAS-Statistical-Models.pdf.
- Wu, C.F.J. (1983), “On the Convergence Properties of the EM Algorithm,” *Annals of Statistics*, 11, 95–103.
- Zaslavsky, A. M., Zaborski, L. B., and Cleary, P. D. (2004), “Plan, Geographical, and Temporal Variation of Consumer Assessments of Ambulatory Health Care,” *Health Services Research*, 39, 1467–1486.