

ESTIMATION AND INFERENCE FOR VERY LARGE LINEAR MIXED EFFECTS MODELS

Katelyn Gao and Art B. Owen

Intel Inc. and Stanford University

Abstract: Linear mixed models with large imbalanced crossed random effects structures pose severe computational problems for maximum likelihood estimation and for Bayesian analysis. The costs can grow as fast as $N^{3/2}$ when there are N observations. Such problems arise in any setting where the underlying factors satisfy a many-to-many rather than a nested relationship. The former are common in electronic commerce applications, where N can be quite large. Methods that do not account for the correlation structure can greatly underestimate the uncertainty. Thus, we propose a method of moments approach that takes account of the correlation structure and that can be computed at a cost of $O(N)$. The method of moments can be parallelized easily, because it is based on sums and it does not require parametric distributional assumptions, tuning parameters, or convergence diagnostics. For the regression coefficients, we give conditions for consistency and asymptotic normality, as well as a consistent variance estimate. We also provide the conditions necessary for a consistent estimation of the variance components, as well as consistent estimates of a mildly conservative upper bound on the variance of the variance component estimates. All of these computations require a total processing time of $O(N)$. We illustrate the algorithm using data from Stitch Fix, where the crossed random effects correspond to clients and items. Here, a naive analysis can overestimate the effective sample size by hundreds and, thus yield unreliable conclusions about the parameters.

Key words and phrases: Crossed random effects, linear mixed models, scalable inference.

1. Introduction

The field of statistics is confronting two important challenges at present. The first is the increasing prevalence of ever larger data sets, sometimes described as ‘big data’; see, for instance, Provost and Fawcett (2013) and Varian (2014). The second is the reproducibility crisis, in which published findings cannot be replicated. This problem was presented clearly by Ioannidis (2005) among others, and has led to the American Statistical Association releasing a statement on p -values (Wasserstein and Lazar (2016)).

We might naively hope that the first problem would somehow remove the second problem, owing to the decrease in uncertainty. However, this may not be true in reality, where difficulties remain when data are not independent. We consider one such situation, in which a crossed random effects structure exists in the data. This structure introduces a dense tangle of correlations that can sharply reduce the effective sample size of the data at hand. If, as we suspect, most data scientists treat these large data sets as independent and identically distributed (i.i.d.) samples, then they will greatly underestimate the uncertainty in their fitted models. The usual methods for solving this problem, whether by maximum likelihood or restricted maximum likelihood (REML), or Bayes, have a cost that grows superlinearly in the sample size and thus cannot be run on the largest data sets. We present and study a method of moments approach with a cost that scales linearly in the problem size, among other advantages.

The sort of data that motivates us arise in e-commerce applications, and include factors such as cookies, customer IDs, query strings, IP addresses, product IDs (e.g., SKUs), URLs, and so on. The most direct way to handle them is to treat them as categorical variables that simply happen to have a large number of levels, including many that have not yet appeared in the data. We think that a random effects model is more appropriate (McCulloch, Searle and Neuhaus (2008)). For instance, internet cookies are cleared regularly, and hence any specific cookie is likely to disappear relatively quickly. It is therefore more appropriate to consider the specific cookies in a data set as a sample from some distribution, that is, as a random effect. Similarly there is turnover in popular products and queries, which motivates treating them as random effects too.

While the largest crossed random effect data sets we know of occur in e-commerce and social media (for example, the Netflix data set (Bennett and Lanning (2007))), we expect the problem to arise in other settings where data set sizes are growing. The crossed random effects structure is fundamental. Any setting with a many-to-many mapping of factor levels involves crossed effects that one might want to model as random. In agriculture and genomics, there are gene-by-environment or gene-by-patient crosses. In education, neither schools nor neighborhoods are perfectly nested within the other (Raudenbush (1993)), and in multiyear data sets there is a many-to-many relationship between teachers and students.

When our chosen model involves only one of these random effect entities then a hierarchical model, based on Bayes or empirical Bayes, can be quite effective (Yu and Meng (2011); Gelman et al. (2012)). Things change considerably when

we want to use two or more crossed random effects. In this study, we consider the following model,

Model 1. Two-factor linear mixed effects:

$$\begin{aligned}
 Y_{ij} &= x_{ij}^\top \beta + a_i + b_j + e_{ij}, \quad x_{ij} \in \mathbb{R}^p, \quad i, j \in \mathbb{N} \quad \text{where,} \\
 a_i &\overset{i.i.d.}{\sim} (0, \sigma_A^2), \quad b_j \overset{i.i.d.}{\sim} (0, \sigma_B^2), \quad e_{ij} \overset{i.i.d.}{\sim} (0, \sigma_E^2) \quad (\text{independently}) \text{ and,} \quad (1.1) \\
 E(a_i^4) &< \infty, \quad E(b_j^4) < \infty, \quad E(e_{ij}^4) < \infty.
 \end{aligned}$$

For instance, customer i might assign a score Y_{ij} to product j . Then x_{ij} contains features about the customer or product or some joint properties of both, β is of interest to the company choosing a product to recommend, b_j measures some general appeal of the product not captured by the features in x_{ij} , a_i captures the variation in which customers are harder or easier to please, and e_{ij} is an error term. This is a mixed effects model because it contains both random effects a_i , b_j and fixed effects x_{ij} .

Model 1 describes any ij pair, but the given data set will only contain some finite number N of them. If the available data are laid out as rows i and columns j with R distinct rows and C distinct columns, then the cost of fitting a generalized least squares regression model for β scales as $O((R + C)^3)$ because it solves a $p \times p$ system of equations with $p \geq R + C$. See Searle, Casella and McCulloch (1992), Raudenbush (1993) and Bates (2014). Now because $RC \geq N$ we have $\max(R, C) \geq \sqrt{N}$ and $(R + C)^3 > N^{3/2}$.

Gao and Owen (2017) consider an intercept-only version of Model 1 where $x_{ij}^\top \beta$ is simply a constant $\mu \in \mathbb{R}$ for all i and j . They find that Markov chain Monte Carlo (MCMC) method does not solve the inference problem under the assumption that the random effects are normally distributed. All of the MCMC methods considered either failed to mix, or converged to the wrong answer, even at modest sample sizes. For the specific case of a Gibbs sampler and Gaussian a_i , b_j , and e_{ij} , using methods from Roberts and Sahu (1997) they prove that it will take $O(N^{1/2})$ iterations costing $O(N)$ each to converge, for a total cost of $O(N^{3/2})$. Fox (2013) presents a very general equivalence between the convergence rate of an iterative equation solver and the convergence rate of an associated MCMC scheme. Therefore these identical rates may be a sign of a deeper connection. Consensus Bayes (Scott et al. (2016)) splits the data into shards, one per processor. However the data given to each shard has to be independent and here data sets corresponding to a subset of rows will have correlations owing to their commonly sampled columns (and vice versa). As an alternative to MCMC,

variational Bayesian methods can be used to approximate the posterior distribution of the parameters by minimizing the divergence of a chosen parametric density from the posterior (Blei, Kucukelbir and McAuliffe (2017)). However, it is not straightforward to choose the parametric density and little is known about the theoretical properties of such methods.

Likelihood-based approaches are commonly used to analyze Model 1 under a Gaussian model for the random effects. For example, see Jiang (2007). Maximum likelihood (ML) maximizes the log-likelihood of the data with respect to the parameters (Demidenko (2013); Hartley and Rao (1967)), while REML mitigates the bias of the ML estimates by estimating the variance components using ML on some residuals that do not depend on β . See (McCulloch, Searle and Neuhaus (2008, Chap. 6.9)). Various optimization algorithms have been applied to compute ML and REML estimates, including gradient-free algorithms such as BOBYQA (Powell (2009)) and Expectation-Maximization (Dempster, Laird and Rubin (1977)) and gradient-based algorithms such as Gauss-Newton (Bates (2014)). The asymptotic variances of these estimates are readily obtained from the Fisher information matrix. The main disadvantage of likelihood-based approaches is that even computing the value of the likelihood at given values of the parameters requires $O(N^{3/2})$ time (Bates (2014)).

Thus, we find that existing Bayes and likelihood methods are not effective for this problem. Here we present an approach based on the method of moments. We seek estimates $\hat{\beta}$, $\hat{\sigma}_A^2$, $\hat{\sigma}_B^2$ and $\hat{\sigma}_E^2$ along with variance estimates for these quantities. We have three criteria:

- 1) the total computational cost *must* be $O(N)$ time and $O(R + C)$ space,
- 2) the variance estimates *should* be reliable or conservative, and
- 3) we *prefer* $\hat{\beta}$ to be statistically efficient.

We regard the first criterion as a constraint that must be met. For the second criterion, a mild over-estimate of $\text{Var}(\hat{\beta})$ is acceptable in order to keep the costs in $O(N)$. The third criterion is to be met as well as we can, subject to constraints given by the first two. Computational efficiency is more important than statistical efficiency in this context. For very large N , requiring $O(N^{3/2})$ computation is highly unrealistic.

With an apt choice of the estimating equations, the method of moments meets our $O(N)$ time and $O(R + C)$ space criteria, and we show that it can also yield reliable variance estimates. Further advantages of the method of moments are that it does not require parametric distributional assumptions (e.g., Gaussianity), there are no tuning parameters to choose, and most importantly for

large N , it is very well suited to parallel computation. The method of moments is not without drawbacks. Sometimes it yields parameter estimates that are out of bounds, such as negative variance estimates, which are often set to zero Searle, Casella and McCulloch (1992). The method of moments has been highly successful in the case of nested random effects. In Wu, Stute and Zhu (2012), the authors propose estimators with linear computational complexity for not only for the regression coefficients and the variance components, but also for the higher moments of the random effects. See also Wu and Zhu (2010) and Perry (2017).

The remainder of the paper proceeds as follows. Section 2 introduces most of the notation for Model 1, especially the pattern of missingness in the data, and presents some of the asymptotic assumptions. Section 3 presents our algorithm and shows that it takes $O(N)$ time and $O(R+C)$ space. We compute a generalized least squares (GLS) estimate for a model with either row or column variance components, but not both. We choose based on an efficiency criterion. Then we estimate $\text{Var}(\hat{\beta})$ accounting for all three error terms including the one left out of the GLS estimate. Section 4 illustrates our algorithm using ratings data from Stitch Fix. There Y_{ij} is a rating, from a 10-point scale, by customer i on item j , with features x_{ij} . Compared to ordinary least squares (OLS) estimates, the random effects model leads to standard errors on coefficients β_j that can be more than 10 times higher. This may be interpreted as an effective sample size which is less than 1% of the nominal sample size. Section 5 gives conditions under which $\hat{\beta}$ and the variance components are consistent. There is also a central limit theorem (CLT) for $\hat{\beta}$. Section 6 compares our method of moments estimator to a state-of-the-art GLMM code (Bates (2016)) written in Julia (Bezanson et al. (2017)). That algorithm takes $O(N^{3/2})$ cost per iteration, with a number of iterations that, in our simulations, depends on N . For problems where the GLMM code gives an answer we find it more statistically efficient for β and σ_E^2 but not for σ_A^2 or σ_B^2 . Lastly, Section 7 discusses some future work and related literature.

Our method of moments approach is similar to Henderson's classical methods Henderson (1953) for Gaussian data, as presented in Searle, Casella and McCulloch (1992). For an intercept-only model, Gao and Owen (2017) uses U -statistics to find a counterpart to the Henderson I estimator that can be computed in $O(N)$ time and $O(R+C)$ space. We also obtain a variance estimator for the variance components, without assuming a Gaussian distribution. The variance estimator can be computed in $O(N)$ time. It targets a mildly conservative upper bound on the variance as the variance itself, like the one for Henderson's estimates, takes more than $O(N)$ computation. In this study we incorporate fixed effects

along with the random effects, just as Henderson II does in generalizing Henderson I, by transforming the original model to one with only random effects. Like Henderson II, our algorithm alternates between estimating the regression coefficients and the variance components. However, our estimators are different, having linear computational complexity instead of superlinear. Henderson III uses estimating equations that are based on the residual sum of squares when treating subsets of the random effects as fixed effects and fitting OLS models. In addition, Henderson III allows for interactions between fixed and random effects. We believe such interactions are very reasonable in our motivating applications, but incorporating them is beyond the scope of this study.

Our analysis is for a fixed dimension p . This is reasonable for our motivating data from Stitch Fix, where $p \ll N$. It remains to develop methods for cases where $p \rightarrow \infty$ with N .

Another issue that we do not address is the selection bias in the available observations. Sometimes ratings are biased towards the high end because customers seek products that they expect to like and companies endeavor to recommend such products. In other data sets, such as restaurant reviews, customers may be more likely to make a rating when they are either very unhappy or very happy. For such data, the ratings will be biased towards both extremes and away from the middle. Accounting for selection bias requires assumptions or information from outside the given data. Propensity weighting (Imbens and Rubin (2015, Chap. 13)) may fit within our framework. This too is left for future research.

2. Notation and Asymptotic Conditions

Here, we give a fuller presentation of our notation. Equation (1.1) describes the distribution of observed and future data. We call the first index of Y_{ij} the ‘row’ and the second the ‘column’. We use integers i, i', r, r' to index rows and j, j', s, s' for columns, but the actual indices may be URLs, customer IDs, or query strings. The index sets are countably infinite to always leave room for unseen levels in the future.

The variable Z_{ij} takes the value one if (x_{ij}, Y_{ij}) is observed and zero otherwise. We assume that there is at most one observation in position (i, j) . For customer rating data, we suppose that if i has rated j multiple times, then only the most recent rating is retained. We believe that in most other settings, only a negligible fraction of ij pairs will have been duplicated.

The sample size is $N = \sum_{ij} Z_{ij} < \infty$. The number of observations in row i

is $N_{i\bullet} = \sum_j Z_{ij}$ and the number in column j is $N_{\bullet j} = \sum_i Z_{ij}$. The number of distinct rows is $R = \sum_i 1_{N_{i\bullet} > 0}$ and there are $C = \sum_j 1_{N_{\bullet j} > 0}$ distinct columns. In the following, summing over rows i means summing over just the R rows i with $N_{i\bullet} > 0$, and sums over columns are defined similarly. This convention corresponds to what happens when one makes a pass through the whole data set.

Let Z be the matrix containing Z_{ij} . Then $(ZZ^T)_{ii'} = \sum_j Z_{ij}Z_{i'j}$ is the number of columns for which we have data in both rows i and i' . Similarly, $(Z^TZ)_{jj'}$ is the number of rows in which both columns j and j' are observed. Note that $(ZZ^T)_{ii'} \leq N_{i\bullet}$ and $(Z^TZ)_{jj'} \leq N_{\bullet j}$. We will use the following identities:

$$\sum_{ir} (ZZ^T)_{ir} = \sum_j N_{\bullet j}^2, \quad \text{and} \quad \sum_{js} (Z^TZ)_{js} = \sum_i N_{i\bullet}^2.$$

This notation allows for an arbitrary pattern of observations. We mention three special cases. A balanced crossed design has $Z_{ij} = 1_{i \leq R} 1_{j \leq C}$. If $\max_i N_{i\bullet} = 1$ but $\max_j N_{\bullet j} > 1$ then the data have a hierarchical structure with rows nested in columns. If $\max_i N_{i\bullet} = \max_j N_{\bullet j} = 1$, then the observed Y_{ij} have i.i.d. errors. Some of these patterns cause problems for parameter estimation. For example, if the errors are i.i.d., then the variance components are not identifiable. Our assumptions rule these out in order to focus on large genuinely crossed data sets.

The following vectors are useful in our subsequent analyses. Let $v_{1,i}$ be the length- N vector with ones in entries $\sum_{r=1}^{i-1} N_{r\bullet} + 1$ to $\sum_{r=1}^i N_{r\bullet}$ and zeros elsewhere. Similarly, let $v_{2,j}$ be the length- N vector with ones in entries $\sum_{s=1}^{j-1} N_{\bullet s} + 1$ to $\sum_{s=1}^j N_{\bullet s}$ and zeros elsewhere.

Next, we describe our asymptotic assumptions. First

$$\epsilon_R = \max_i \frac{N_{i\bullet}}{N} \rightarrow 0, \quad \text{and} \quad \epsilon_C = \max_j \frac{N_{\bullet j}}{N} \rightarrow 0, \tag{2.1}$$

such that no single row or column dominates. The average row size can be measured by N/R or by $\sum_i N_{i\bullet}^2/N$; the latter is $E(N_{i\bullet})$ when choosing one of the N data points (i, j, x_{ij}, Y_{ij}) at random (uniformly). Similar formulae hold for the average column size. These average row and column sizes are $o(N)$, because

$$\frac{1}{N^2} \sum_i N_{i\bullet}^2 \leq \epsilon_R \rightarrow 0, \quad \text{and} \quad \frac{1}{N^2} \sum_j N_{\bullet j}^2 \leq \epsilon_C \rightarrow 0.$$

We often expect the average row and column sizes to diverge, while growing more

slowly than N :

$$\min\left(\frac{N}{R}, \frac{N}{C}\right) \rightarrow \infty, \quad \text{and}$$

$$\min\left(\frac{1}{N} \sum_i N_{i\bullet}^2, \frac{1}{N} \sum_j N_{\bullet j}^2\right) \rightarrow \infty.$$

We do not however impose these conditions.

Even for large average row and columns sizes, there can still be numerous new or rare entities with $N_{i\bullet} = 1$ or $N_{\bullet j} = 1$. Our analysis can include such small rows and columns without requiring that they be deleted. When there are covariates x_{ij} we need to rule out degenerate settings where the sample variance of x_{ij} does not grow with N or where it is dominated by a handful of observations. We add some such conditions when we prove some CLTs in Section 5.2.

The finite fourth moments $E(a_i^4)$, $E(b_j^4)$ and $E(e_{ij}^4)$ are conveniently described through finite kurtoses κ_A , κ_B and κ_E , respectively. Some of the variance expressions in Gao and Owen (2017) are dominated by terms proportional to $\kappa + 2$ for one of these kurtoses. Following Gao and Owen (2017) we assume that $\min(\kappa_A, \kappa_B, \kappa_E) > -2$. This lower bound rules out some symmetric binary distributions for a_i , b_j and e_{ij} . However, these cases seem unrealistic for our motivating applications.

The randomness in Y_{ij} comes from a_i , b_j and e_{ij} . In some places we combine them into $\eta_{ij} \equiv a_i + b_j + e_{ij}$.

We use the method of moment estimators $\hat{\sigma}_A^2$, $\hat{\sigma}_B^2$ and $\hat{\sigma}_E^2$ for σ_A^2 , σ_B^2 and σ_E^2 , respectively, from Gao and Owen (2017), who provide exact finite sample formulae for the variances of their estimators. They also give asymptotic variance expressions, while letting ϵ_R , ϵ_C , R/N and C/N approach zero. The Stitch Fix data that we consider in Section 4 does not have a very small value for R/N . Here we develop nonasymptotic magnitude bounds for the bias and variance that do not require R/N and C/N to be close to zero. They need only be bounded away from one.

Theorem 1. *Suppose that $\max(R/N, C/N) \leq \theta$ for some $\theta < 1$ and let $\epsilon = \max(\epsilon_R, \epsilon_C)$. Then the moment-based estimators from Gao and Owen (2017) satisfy*

$$E(\hat{\sigma}_A^2) = (\sigma_A^2 + \Upsilon)(1 + O(\epsilon)),$$

$$E(\hat{\sigma}_B^2) = (\sigma_B^2 + \Upsilon)(1 + O(\epsilon)), \quad \text{and}$$

$$E(\sigma_E^2) = (\sigma_E^2 + \Upsilon)(1 + O(\epsilon)),$$

where

$$\Upsilon \equiv \sigma_A^2 \frac{\sum_i N_{i\bullet}^2}{N^2} + \sigma_B^2 \frac{\sum_j N_{\bullet j}^2}{N^2} + \frac{\sigma_E^2}{N} = O(\epsilon).$$

Furthermore

$$\max(\text{Var}(\hat{\sigma}_A^2), \text{Var}(\hat{\sigma}_B^2), \text{Var}(\hat{\sigma}_E^2)) = O\left(\frac{\sum_i N_{i\bullet}^2}{N^2} + \frac{\sum_j N_{\bullet j}^2}{N^2}\right) = O(\epsilon).$$

Proof. See Section S1 in the Supplementary Material.

Theorem 1 has the same variance rate for all variance components. In our computed examples $\text{Var}(\hat{\sigma}_E^2) \ll \min(\text{Var}(\hat{\sigma}_A^2), \text{Var}(\hat{\sigma}_B^2))$ because $N \gg \max(R, C)$, a condition not imposed in Theorem 1. The bias and variance are both $O(\epsilon)$. Therefore, a (conservative) effective sample size is $O(1/\epsilon)$. The quantity Υ appearing in Theorem 1 is $\text{Var}(\bar{Y}_{\bullet\bullet})$ where $\bar{Y}_{\bullet\bullet} = (1/N) \sum_{ij} Z_{ij} Y_{ij}$. The variances of the estimated variance components contain similar quantities to Υ although kurtoses and other quantities appear in their implied constants.

3. An Alternating Algorithm

Our estimation procedure for Model 1 is given in Algorithm 1. We alternate twice between finding $\hat{\beta}$ and finding the variance component estimates $\hat{\sigma}_A^2$, $\hat{\sigma}_B^2$, and $\hat{\sigma}_E^2$. One can continue iterating if desired, but our theory shows that two iterations suffice. In our experience, additional iterations did not change the estimates much. Further details of these steps, including the way we choose generalized least squares (GLS) estimator for step 3, are given in the next two subsections.

The data are a collection of (i, j, x_{ij}, Y_{ij}) tuples. A pass over the data proceeds via iteration over all tuples in the data set. Such a pass may generate $O(R + C)$ intermediate values, which are retained for future computations.

3.1. Algorithm 1

Step 1

The first step of Algorithm 1 is to compute the OLS estimate of β . Let $X \in \mathbb{R}^{N \times p}$ have rows x_{ij} in some order and let $Y \in \mathbb{R}^N$ be elements Y_{ij} in the

Algorithm 1 Alternating Algorithm

Estimate β via ordinary least squares (OLS): $\hat{\beta} = \hat{\beta}_{\text{OLS}}$.

Let $\hat{\sigma}_A^2$, $\hat{\sigma}_B^2$, and $\hat{\sigma}_E^2$ be the method of moments estimates from (Gao and Owen (2017)) defined on the data $(i, j, \hat{\eta}_{ij})$, where $\hat{\eta}_{ij} = Y_{ij} - x_{ij}^\top \hat{\beta}_{\text{OLS}}$.

Compute a more efficient $\hat{\beta}$ using $\hat{\sigma}_A^2$, $\hat{\sigma}_B^2$, and $\hat{\sigma}_E^2$. If $\hat{\sigma}_A^2 \max_i N_{i\bullet} \geq \hat{\sigma}_B^2 \max_j N_{\bullet j}$, estimate β via GLS accounting for row correlations: $\hat{\beta} = \hat{\beta}_{\text{RLS}}$. Otherwise, estimate it via GLS accounting for column correlations: $\hat{\beta} = \hat{\beta}_{\text{CLS}}$.

Repeat step 2 using $\hat{\eta}_{ij} = Y_{ij} - x_{ij}^\top \hat{\beta}$ with $\hat{\beta}$ from step 3.

Compute an estimate $\widehat{\text{Var}}(\hat{\beta})$ for $\hat{\beta}$ from step 3 using $\hat{\sigma}_A^2$, $\hat{\sigma}_B^2$ and $\hat{\sigma}_E^2$ from step 4.

same order. Then,

$$\hat{\beta}_{\text{OLS}} = (X^\top X)^{-1} X^\top Y = \left(\sum_{ij} Z_{ij} x_{ij} x_{ij}^\top \right)^{-1} \sum_{ij} Z_{ij} x_{ij} Y_{ij}. \quad (3.1)$$

In one pass over the data, we can compute $X^\top X$ and $X^\top Y$ and solve for $\hat{\beta}$. Solving the normal equations this way is easy to parallelize but typically incurs a larger roundoff error than the usual alternative based on computing the SVD of X . The numerical conditioning of the SVD computation essentially doubles the number of floating point bits available in comparison to solving the normal equations. One can compensate by solving normal equations in extended precision. It costs $O(p^3)$ to compute $\hat{\beta}_{\text{OLS}}$ and so the cost of step 1 is $O(Np^2 + p^3)$. The space cost is $O(p^2)$.

Step 2

Step 2 uses the algorithm from Gao and Owen (2017) to compute the variance component estimates $\hat{\sigma}_A^2$, $\hat{\sigma}_B^2$ and $\hat{\sigma}_E^2$ in $O(N)$ time and $O(R + C)$ space; see Section 3.2. This takes $O(Np)$ time to recompute $\hat{\eta}_{ij}$.

Step 3

GLS estimators: First we define and compare GLS estimators of β accounting for row correlations, column correlations, or both. These estimators are most easily presented through a reordering of the data. Our algorithm does not have to sort the data which would be a major inconvenience in our motivating applications. We work with one row ordering of the data, in which ij precedes $i'j'$ whenever $i < i'$ and with one column ordering of the data. Let P be the $N \times N$ permutation matrix corresponding to the transformation of the column ordering to the row ordering. Let $A_R \in \mathbb{N}^{N \times N}$ be the block diagonal matrix with i 'th block $1_{N_i} \cdot 1_{N_i}^\top$.

and let $B_C \in \mathbb{N}^{N \times N}$ the block diagonal matrix with j 'th block $1_{N_{\bullet j}} 1_{N_{\bullet j}}^\top$.

If Y is given in the row ordering, then

$$\text{Cov}(Y) = V_R \equiv \sigma_E^2 I_N + \sigma_A^2 A_R + \sigma_B^2 B_R, \quad \text{for } B_R = P B_C P^\top. \quad (3.2)$$

For Y in the column ordering,

$$\text{Cov}(Y) = V_C \equiv \sigma_E^2 I_N + \sigma_A^2 A_C + \sigma_B^2 B_C, \quad \text{for } A_C = P^\top A_R P. \quad (3.3)$$

GLS algorithms based on (3.2) or (3.3) have computational complexity $O(N^{3/2})$. This is better than the $O(N^3)$ that we might have faced had V_R or V_C been arbitrary dense matrices, rather than being comprised of the identity and some low rank block diagonal matrices. However, it is still too slow for large scale applications.

In a hierarchical model where only row correlations were present we could take $\sigma_B^2 = 0$ and define

$$\hat{\beta}_{\text{RLS}} = (X^\top \hat{V}_A^{-1} X)^{-1} X^\top \hat{V}_A^{-1} Y, \quad \text{for } \hat{V}_A = \hat{\sigma}_E^2 I_N + \hat{\sigma}_A^2 A_R, \quad (3.4)$$

using sample estimates $\hat{\sigma}_A^2$ and $\hat{\sigma}_E^2$ of σ_A^2 and σ_E^2 , respectively. This GLS estimator of β accounts for the intra-row correlations in the data. Similarly, the GLS estimator of β accounting for the intra-column correlations is

$$\hat{\beta}_{\text{CLS}} = (X^\top \hat{V}_B^{-1} X)^{-1} X^\top \hat{V}_B^{-1} Y, \quad \text{for } \hat{V}_B = \hat{\sigma}_E^2 I_N + \hat{\sigma}_B^2 B_C. \quad (3.5)$$

We show next that $\hat{\beta}_{\text{RLS}}$ and $\hat{\beta}_{\text{CLS}}$ can be computed in $O(N)$ time.

GLS Computations in $O(N)$ cost: From the Woodbury formula (Hager (1989)) and defining $Z_a \in \{0, 1\}^{N \times R}$ as the matrix with i th column $v_{1,i}$ (from Section 2), we have

$$\begin{aligned} & X^\top \hat{V}_A^{-1} X \\ &= X^\top (\hat{\sigma}_E^2 I_N + \hat{\sigma}_A^2 Z_a Z_a^\top)^{-1} X \\ &= \frac{X^\top X}{\hat{\sigma}_E^2} - \frac{\hat{\sigma}_A^2}{\hat{\sigma}_E^2} X^\top Z_a \text{diag} \left(\frac{1}{\hat{\sigma}_E^2 + \hat{\sigma}_A^2 N_{i\bullet}} \right) Z_a^\top X \\ &= \frac{1}{\hat{\sigma}_E^2} \sum_{ij} Z_{ij} x_{ij} x_{ij}^\top - \frac{\hat{\sigma}_A^2}{\hat{\sigma}_E^2} \sum_i \frac{1}{\hat{\sigma}_E^2 + \hat{\sigma}_A^2 N_{i\bullet}} \left(\sum_j Z_{ij} x_{ij} \right) \left(\sum_j Z_{ij} x_{ij} \right)^\top. \end{aligned}$$

Similarly, $X^\top \hat{V}_A^{-1} Y$ is equal to

$$\frac{1}{\hat{\sigma}_E^2} \sum_{ij} Z_{ij} x_{ij} Y_{ij} - \frac{\hat{\sigma}_A^2}{\hat{\sigma}_E^2} \sum_i \frac{1}{\hat{\sigma}_E^2 + \hat{\sigma}_A^2 N_{i\bullet}} \left(\sum_j Z_{ij} x_{ij} \right) \left(\sum_j Z_{ij} Y_{ij} \right).$$

One pass over the data allows us to compute $\sum_{ij} Z_{ij} x_{ij} x_{ij}^\top$ and $\sum_{ij} Z_{ij} x_{ij} Y_{ij}$, as well as $N_{i\bullet}$, and the row sums $\sum_j Z_{ij} x_{ij}$ and $\sum_j Z_{ij} Y_{ij}$ for $i = 1, \dots, R$. The cost is $O(Np^2)$ time and $O(Rp)$ space. None of these quantities require that we sort the data. Next, we compute $X^\top \hat{V}_A^{-1} X$ and $X^\top \hat{V}_A^{-1} Y$ in time $O(Rp^2)$. Then, $\hat{\beta}_{\text{RLS}}$ is computed in $O(p^3)$. Hence, $\hat{\beta}_{\text{RLS}}$ can be found within $O(Rp)$ space and $O(Np^2 + p^3) = O(Np^2)$ time. Clearly $\hat{\beta}_{\text{CLS}}$ costs $O(Cp)$ space and $O(Np^2)$ time.

Efficiencies: We can compute either $\hat{\beta}_{\text{RLS}}$ or $\hat{\beta}_{\text{CLS}}$ in our computational budget. We choose RLS if the variance component associated with the rows is dominant and CLS otherwise. The choice could be made dependent on X but in many applications one considers numerous different X matrices and we prefer to have a single choice for all regressions. Accordingly, we find a lower bound on the efficiency of RLS when X is a single nonzero vector $\mathbf{x} \in \mathbb{R}^{N \times 1}$. We choose RLS if that lower bound is higher than the corresponding bound for CLS, in this $p = 1$ setting.

The full GLS estimator is $\hat{\beta}_{\text{GLS}} = (X^\top V_R^{-1} X)^{-1} X^\top V_R^{-1} Y$ when the data are ordered by rows and $(X^\top V_C^{-1} X)^{-1} X^\top V_C^{-1} Y$ when the data are ordered by columns. For data ordered by rows, the efficiency of $\hat{\beta}_{\text{RLS}}$ is

$$\text{eff}_{\text{RLS}} = \frac{\text{Var}(\hat{\beta}_{\text{GLS}})}{\text{Var}(\hat{\beta}_{\text{RLS}})} = \frac{(\mathbf{x}^\top V_A^{-1} \mathbf{x})^2}{(\mathbf{x}^\top V_A^{-1} V_R V_A^{-1} \mathbf{x})(\mathbf{x}^\top V_R^{-1} \mathbf{x})}. \tag{3.6}$$

For data ordered by columns, the corresponding efficiency of $\hat{\beta}_{\text{CLS}}$ is

$$\text{eff}_{\text{CLS}} = \frac{\text{Var}(\hat{\beta}_{\text{GLS}})}{\text{Var}(\hat{\beta}_{\text{CLS}})} = \frac{(\mathbf{x}^\top V_B^{-1} \mathbf{x})^2}{(\mathbf{x}^\top V_B^{-1} V_C V_B^{-1} \mathbf{x})(\mathbf{x}^\top V_C^{-1} \mathbf{x})}. \tag{3.7}$$

The next two theorems establish lower bounds on these efficiencies.

Theorem 2. *Let A be a positive-definite Hermitian matrix and \mathbf{u} be a unit vector. If the eigenvalues of A are bounded below by $m > 0$ and above by $M < \infty$, then*

$$(\mathbf{u}^\top A \mathbf{u})(\mathbf{u}^\top A^{-1} \mathbf{u}) \leq \frac{(m + M)^2}{4mM}.$$

Equality may hold, for example when $\mathbf{u}^\top A \mathbf{u} = (M + m)/2$ and the only roots of

A are m and M .

Proof. This is Kantorovich’s inequality (Marshall and Olkin (1990)).

From the two applications of Theorem 2 on (3.6) and (3.7) we prove the following.

Theorem 3. For $p = 1$ and $\sigma_E^2 > 0$, let eff_{RLS} and eff_{CLS} be defined as in (3.6) and (3.7), respectively. Then,

$$\begin{aligned} \text{eff}_{\text{RLS}} &\geq \frac{4\sigma_E^2(\sigma_E^2 + \sigma_B^2 \max_j N_{\bullet j})}{(2\sigma_E^2 + \sigma_B^2 \max_j N_{\bullet j})^2} \quad \text{and} \\ \text{eff}_{\text{CLS}} &\geq \frac{4\sigma_E^2(\sigma_E^2 + \sigma_A^2 \max_i N_{i\bullet})}{(2\sigma_E^2 + \sigma_A^2 \max_i N_{i\bullet})^2}. \end{aligned}$$

Both inequalities are tight.

Proof. See Section S2.1 in the Supplementary Material.

After some algebra, we see that the worst case efficiency of $\hat{\beta}_{\text{RLS}}$ is higher than that of $\hat{\beta}_{\text{CLS}}$ when $\sigma_A^2 \max_i N_{i\bullet} > \sigma_B^2 \max_j N_{\bullet j}$. We set $\hat{\beta}$ to be $\hat{\beta}_{\text{RLS}}$ when $\hat{\sigma}_A^2 \max_i N_{i\bullet} \geq \hat{\sigma}_B^2 \max_j N_{\bullet j}$, and $\hat{\beta}_{\text{CLS}}$ otherwise.

Optimizing a lower bound does not necessarily optimize the quantity of interest, and so we expect that our choice here is not the only reasonable one. The efficiency of $\hat{\beta}_{\text{RLS}}$ depends only on the ratio $\hat{\sigma}_A^2/\hat{\sigma}_E^2$. We investigated GLS estimators of β based on $\hat{V}_A = \hat{\sigma}_A^2 A_R + (\hat{\sigma}_E^2 + \lambda \hat{\sigma}_B^2) I_N$ for λ chosen by the Kantorovich inequality. However, this did not appear to improve the accuracy over our default choice in some simulations. In practice, one can also compute both $\hat{\beta}_{\text{RLS}}$ and $\hat{\beta}_{\text{CLS}}$ and compare $\widehat{\text{Var}}(\hat{\beta}_{\text{RLS}})$ and $\widehat{\text{Var}}(\hat{\beta}_{\text{CLS}})$.

Steps 4 and 5

Step 4 is just like step 2 and it costs $O(Np)$ time. Step 5 is described in Section 5.3 where we derive $\text{Var}(\hat{\beta}_{\text{RLS}})$ and $\text{Var}(\hat{\beta}_{\text{CLS}})$.

3.2. Method of moments (Steps 2 and 4)

In this subsection, we discuss steps 2 and 4 of Algorithm 1 in more detail. The errors $Y_{ij} - x_{ij}^T \beta$ follow a two-factor crossed random effects model (Gao and Owen (2017)). If $\hat{\beta}$ is a good estimate of β , then the residuals $\hat{\eta}_{ij} = Y_{ij} - x_{ij}^T \hat{\beta}$ approximately follow a two-factor crossed random effects model with $\mu = 0$ and variance components σ_A^2 , σ_B^2 , and σ_E^2 .

We estimate σ_A^2 , σ_B^2 , and σ_E^2 , using the algorithm from Gao and Owen (2017) with data $(i, j, \hat{\eta}_{ij})$. That algorithm gives unbiased estimates of the variance

components in a two-factor crossed random effects model.

The algorithm of Gao and Owen (2017) applies the method of moments to three statistics; a weighted sum of within-row sample variances, a weighted sum of within-column sample variances, and a multiple of the full sample variance. For Algorithm 1, these are:

$$\begin{aligned} U_a(\hat{\beta}) &= \sum_i S_{i\bullet}, & S_{i\bullet} &= \sum_j Z_{ij}(\hat{\eta}_{ij} - \hat{\eta}_{i\bullet})^2, \\ U_b(\hat{\beta}) &= \sum_j S_{\bullet j}, & S_{\bullet j} &= \sum_i Z_{ij}(\hat{\eta}_{ij} - \hat{\eta}_{\bullet j})^2, & \text{and} \\ U_e(\hat{\beta}) &= \sum_{ij} Z_{ij}(\hat{\eta}_{ij} - \hat{\eta}_{\bullet\bullet})^2, \end{aligned} \quad (3.8)$$

where subscripts replaced by \bullet are averaged over. The variance component estimates are obtained by solving the system

$$M \begin{pmatrix} \hat{\sigma}_A^2 \\ \hat{\sigma}_B^2 \\ \hat{\sigma}_E^2 \end{pmatrix} = \begin{pmatrix} U_a(\hat{\beta}) \\ U_b(\hat{\beta}) \\ U_e(\hat{\beta}) \end{pmatrix}, \quad M = \begin{pmatrix} 0 & N - R & N - R \\ N - C & 0 & N - C \\ N^2 - \sum_i N_{i\bullet}^2 & N^2 - \sum_j N_{\bullet j}^2 & N^2 - N \end{pmatrix}. \quad (3.9)$$

The matrix M is nonsingular under very weak conditions. It suffices to have $R \geq 2$, $C \geq 2$, $\epsilon_R \leq 1/2$ and $\epsilon_C \leq 1/2$ (Gao and Owen (2017, Sec. 4.1)).

Gao and Owen (2017) compute the U -statistics in one pass over the data taking $O(N)$ time and $O(R+C)$ space. Solving (3.9) takes constant time. Thus, steps 2 and 4 each have computational complexity $O(N)$ and space complexity $O(R+C)$.

4. Stitch Fix Rating Data

Stitch Fix sells clothing, primarily women's clothing. They mail their clients a sample of clothing items. A client keeps and purchases some items and returns the others. It is important to predict which items a client will like. In the context of our model, client i might receive item j and then rate that item with a score Y_{ij} .

Stitch Fix provided us with some of their client ratings data. These data are fully anonymized and contain no personally identifying information. The data provided by Stitch Fix is a sample of their data, and consequently does not reflect their actual numbers of clients, items or their ratios, for example. Nonetheless this is an interesting data set with which to illustrate a linear mixed

effects model.

We received data on clients' ratings of items they received, as well as the following information about the clients and items. For client i and item j , the response is a composite rating Y_{ij} on a scale from 1 to 10. There was a categorical variable giving the item's material, with 23 categories. We also received a binary variable indicating whether the item style is considered to be 'edgy', and another indicating whether the client likes edgy styles. Similarly, there was another pair of binary variables indicating whether items were labeled 'boho' (Bohemian) and whether the client likes boho items. Finally, there was a match score, which is an estimate of the probability that the client keeps the item, predicted before it is actually sent. The match score is a prediction from a baseline model and is not representative of all algorithms used at Stitch Fix.

The observation pattern in the data is as follows. We received $N = 5,000,000$ ratings on $C = 6,318$ items by $R = 762,752$ clients. Thus $C/N \doteq 0.00126$ and $R/N \doteq 0.153$. The latter ratio indicates that only a relatively small number of ratings from each client are included in the data (their full shipment history is not included in the sampled data). The data are not dominated by a single row or column because $\epsilon_R \doteq 9 \times 10^{-6}$ and $\epsilon_C \doteq 0.0143$. Similarly

$$\begin{aligned} \frac{N}{\sum_i N_{i\bullet}^2} &\doteq 0.103, & \frac{\sum_i N_{i\bullet}^2}{N^2} &\doteq 1.95 \times 10^{-6}, \\ \frac{N}{\sum_j N_{\bullet j}^2} &\doteq 1.22 \times 10^{-4}, \quad \text{and} & \frac{\sum_j N_{\bullet j}^2}{N^2} &\doteq 0.00164. \end{aligned}$$

Our two-factor linear mixed effects model for this data is given as follows.

Model 2. For client i and item j ,

$$\begin{aligned} \text{rating}_{ij} = & \beta_0 + \beta_1 \text{match}_{ij} + \beta_2 \mathbb{I}\{\text{client edgy}\}_i + \beta_3 \mathbb{I}\{\text{item edgy}\}_j \\ & + \beta_4 \mathbb{I}\{\text{client edgy}\}_i * \mathbb{I}\{\text{item edgy}\}_j + \beta_5 \mathbb{I}\{\text{client boho}\}_i \\ & + \beta_6 \mathbb{I}\{\text{item boho}\}_j + \beta_7 \mathbb{I}\{\text{client boho}\}_i * \mathbb{I}\{\text{item boho}\}_j \\ & + \beta_8 \text{material}_{ij} + a_i + b_j + e_{ij}. \end{aligned}$$

Here material_{ij} is a categorical variable that is implemented via indicator variables for each type of material. We chose 'Polyester', the most common material, as the baseline.

In a regression analysis, Model 2 would be just one of many models one might consider. There would be numerous ways to encode the variables, and

Table 1. Stitch Fix Regression Results (omitting material type).

	$\hat{\beta}_{\text{OLS}}$	$\widehat{\text{se}}_{\text{OLS}}(\hat{\beta}_{\text{OLS}})$	$\widehat{\text{se}}(\hat{\beta}_{\text{OLS}})$	$\hat{\beta}$	$\widehat{\text{se}}(\hat{\beta})$
Intercept	4.635*	0.005397	0.05808	5.110*	0.01250
Match	5.048*	0.01174	0.1464	3.529*	0.02153
$\mathbb{I}\{\text{client edgy}\}$	0.001020	0.002443	0.004593	0.001860	0.003831
$\mathbb{I}\{\text{item edgy}\}$	-0.3358*	0.004253	0.03730	-0.3328*	0.01542
$\mathbb{I}\{\text{both edgy}\}$	0.3925*	0.006229	0.01352	0.3864*	0.006432
$\mathbb{I}\{\text{client boho}\}$	0.1386*	0.002264	0.004354	0.1334*	0.003622
$\mathbb{I}\{\text{item boho}\}$	-0.5499*	0.005981	0.03049	-0.6261*	0.01661
$\mathbb{I}\{\text{both boho}\}$	0.3822*	0.007566	0.01057	0.3837*	0.007697

the coefficients in any one model would depend on which other variables were included. The odds of settling on exactly this model are low. Our focus is on the estimated standard errors due to variance components and so we will work with a naive face-value interpretation of the coefficients β_j in Model 2. If the emphasis is on prediction, then one can use $x_{ij}^T \hat{\beta}$ perhaps adding shrunken row and/or column means of the residuals. See Gao and Owen (2017) for a discussion of how estimates of σ_A^2 , σ_B^2 , and σ_E^2 can be used to shrink row and/or column means in the intercept-only setting. Even in prediction, underestimating the uncertainty in $\hat{\beta}$ could be costly.

Suppose that one ignored client and item random effects and simply ran OLS. Table 1 shows the results for all coefficients except the material type indicator variables. Section S4 of the Supplementary Material has the complete table. The next column has $\widehat{\text{se}}_{\text{OLS}}(\hat{\beta}_{\text{OLS}})$, the standard error that OLS produces for the OLS coefficient estimate. Then $\widehat{\text{se}}_{\text{Mom}}(\hat{\beta}_{\text{OLS}})$ reports a moment-based standard error for $\hat{\beta}_{\text{OLS}}$ using the estimated variance components. The next two columns are the method of moments estimator $\hat{\beta}_{\text{Mom}}$ and its own standard error $\widehat{\text{se}}_{\text{Mom}}(\hat{\beta}_{\text{Mom}})$, respectively, based on the variance components.

Figure 1 shows a graphical presentation of these results. The leftmost panel shows that ignoring the random effects greatly underestimates the uncertainty in the regression coefficients. Furthermore, this underestimation can be tenfold or even hundredfold when interpreted via effective sample sizes. Even if N is ‘big data’, $N/100$ might not be. The right panel shows that properly accounting for uncertainty makes many material indicator variables change from significant to nonsignificant based on a threshold of $|t| \geq 2$. In other words, the difference in effective sample size could have left the user of this model with stronger conclu-

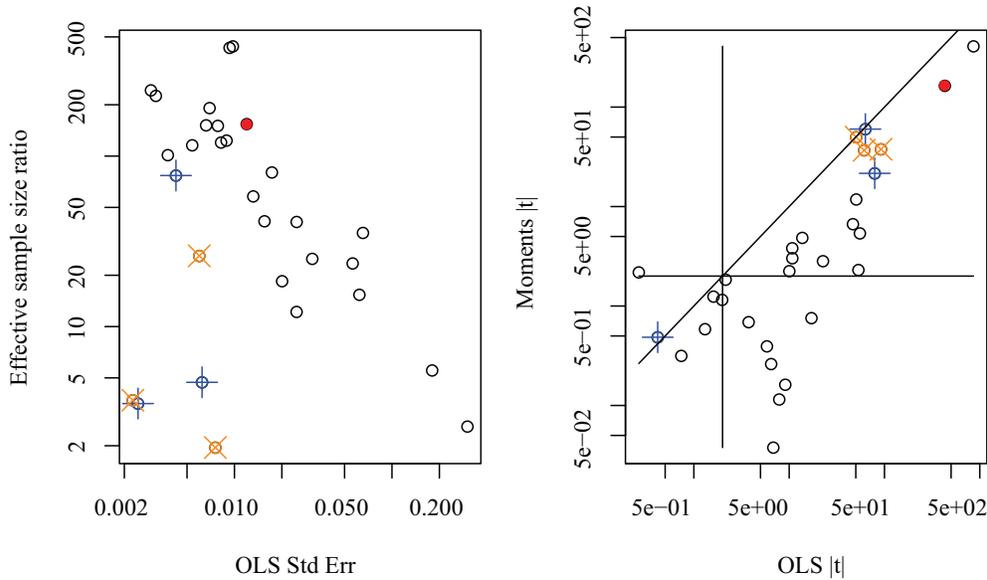


Figure 1. The left panel shows how OLS underestimates the variance. The horizontal axis is $\widehat{se}_{OLS}(\hat{\beta}_{OLS})$. The vertical axis is $[\widehat{se}_{Mom}(\hat{\beta}_{OLS})/\widehat{se}_{OLS}(\hat{\beta}_{OLS})]^2$, which we interpret as the extent to which OLS overestimates the effective sample size. The right panel plots absolute t statistics, $|\hat{\beta}_{Mom}|/\widehat{se}(\beta_{Mom})$ versus $|\hat{\beta}_{OLS}|/\widehat{se}(\beta_{OLS})$. There are reference lines at $|t| = 2$ and at 45 degrees. The Match variable is plotted as a solid circle. The edgy variables include a + and the boho variables include a \times . Material types have open circles.

sions and different decisions to those the or she would otherwise have had. It is likely that industry uses more elaborate models than our simple regression, but a lower than anticipated effective sample size will remain an issue.

The estimated variance components are $\hat{\sigma}_A^2 = 1.133$, $\hat{\sigma}_B^2 = 0.1463$, and $\hat{\sigma}_E^2 = 4.474$. Their standard errors are approximately 0.0046, 0.00089, and 0.0050 respectively, which means that these components are well determined. The error variance component is largest, and the client effect dominates the item effect by almost a factor of eight.

The ‘Match’ variable is significantly and positively associated with rating, indicating that the baseline prediction provided by Stitch Fix is a useful predictor in this data set. However the random effects model reduces its coefficient from about 5 to about 3.5, a change that represents quite a large number of estimated standard errors. We have seen that some clients tend to give higher ratings on average than others. That is, client indicator variables take away some of the explanatory power of the match variable.

Shipping an edgy item to a client who does not like edgy styles is associated with a rating decrease of about 0.33 points, but shipping such an item to a client who does like edgy styles is associated with a small increase in the rating.

The boho indicator variable also has a negative overall estimated coefficient $\hat{\beta}_6 < 0$. The modeled impact of a boho item sent to a boho client is $\hat{\beta}_5 + \hat{\beta}_6 + \hat{\beta}_7 < 0$, unlike the positive result we saw for sending an edgy item to an edgy client. This suggests that it is more difficult to make matches for boho items. Perhaps there is an interaction where ‘boho to boho’ has a positive impact for a sufficiently high value of the match variable. For large data sets, such an interaction can be conveniently handled by filtering the data to cases with $\text{Match}_{ij} \geq t$, and then refitting. We did so but did not find a threshold that yielded $\hat{\beta}_5 + \hat{\beta}_6 + \hat{\beta}_7 > 0$.

Of the materials, ‘Cotton’, ‘Faux Fur’, ‘Leather’, ‘Modal’, ‘Pleather’, ‘PU’, ‘PVC’, ‘Silk’, ‘Spandex’, and ‘Tencel’ are significantly different from the baseline, ‘Polyester’ in our crossed random effects model. ‘PU’ and ‘PVC’ are associated with an increase in rating of at least half a point. Those materials are often used to make shoes and specialty clothing, which may be related to their association with high ratings.

The computations in this section were performed in Python; the code is available at <https://github.com/kxgao/scalable-crossed-mixed-effects>.

5. Asymptotic Behavior

Here we give sufficient conditions to ensure that the parameter estimates $\hat{\beta}$, $\hat{\sigma}_A^2$, $\hat{\sigma}_B^2$, and $\hat{\sigma}_E^2$ obtained from Algorithm 1 are consistent. We also give a CLT for $\hat{\beta}$. We use the sample size growth conditions from Section 2 and some additional conditions on x_{ij} . Our results are conditional on the observed predictors x_{ij} for which $Z_{ij} = 1$.

As in ordinary i.i.d. error regression problems, our CLT requires the information in the observed x_{ij} to grow quickly in every projection while also imposing a limit on the largest x_{ij} . For each i with $N_{i\bullet} > 0$, let $\bar{x}_{i\bullet}$ be the average of those x_{ij} with $Z_{ij} = 1$ and similarly define column averages $\bar{x}_{\bullet j}$.

For a symmetric positive semi-definite matrix V , let $\mathcal{I}(V)$ be the smallest eigenvalue of V . We will need lower bounds on $\mathcal{I}(V)$ for various V to rule out singular or nearly singular designs. Some of those V involve centered variables. In most applications x_{ij} will include an intercept term, and so we assume that the first component of every x_{ij} is equal to one. That term raises some technical difficulties as centering that component always yields zero. Thus, we will treat

this term specially in some of our proofs. For a symmetric matrix $V \in \mathbb{R}^{p \times p}$, we let

$$\mathcal{I}_0(V) = \mathcal{I}((V_{ij})_{2 \leq i, j \leq p})$$

be the smallest eigenvalue of the lower $(p - 1) \times (p - 1)$ submatrix of V .

In our motivating applications, it is reasonable to assume that $\|x_{ij}\|$ are uniformly bounded. We let

$$M_N \equiv \max_{ij} Z_{ij} \|x_{ij}\|^2 \tag{5.1}$$

quantify the largest x_{ij} in the data so far. Some of our results would still hold if we were to let M_N grow slowly with N . To focus on the essential ideas, we simply take $M_N \leq M_\infty < \infty$ for all N .

5.1. Consistency

First, we give conditions under which $\hat{\beta}_{OLS}$ from step 1 is consistent.

Theorem 4. *Let $\max(\epsilon_R, \epsilon_C) \rightarrow 0$ and $\mathcal{I}(X^T X) \geq cN$ for some $c > 0$, as $N \rightarrow \infty$. Then $E(\|\hat{\beta}_{OLS} - \beta\|^2) = O((\epsilon_R + \epsilon_C)/\mathcal{I}(X^T X/N)) \rightarrow 0$ and $\hat{\beta}_{OLS} \xrightarrow{p} \beta$.*

Proof. See Section S3.1 in the Supplementary Material.

Second, we show that the variance component estimates computed in step 2 are consistent. Recall that we compute the U -statistics (3.8) on data $(i, j, \hat{\eta}_{ij} = Y_{ij} - x_{ij}^T \hat{\beta})$ and use them to obtain estimates $\hat{\sigma}_A^2$, $\hat{\sigma}_B^2$, and $\hat{\sigma}_E^2$ via (3.9).

Theorem 5. *Suppose that as $N \rightarrow \infty$ with $\max(\epsilon_R, \epsilon_C) \rightarrow 0$, $\max(R, C)/N \leq \theta \in (0, 1)$, $\hat{\beta} \xrightarrow{p} \beta$, and that M_N is bounded. Then $\hat{\sigma}_A^2 \xrightarrow{p} \sigma_A^2$, $\hat{\sigma}_B^2 \xrightarrow{p} \sigma_B^2$, and $\hat{\sigma}_E^2 \xrightarrow{p} \sigma_E^2$.*

Proof. See Section S3.2 in the Supplementary Material.

From Theorem 4, the estimate of β obtained in step 1 of Algorithm 1 is consistent. Therefore, from Theorem 5, the variance component estimates obtained in step 2 are consistent, given the combined assumptions of those two theorems. The proof of Theorem 4 shows that the estimated variance components differ by $O(\|\hat{\beta} - \beta\|^2 + \epsilon\|\hat{\beta} - \beta\|)$ from what we would get replacing $\hat{\beta}$ by an oracle value β and computing variance components of $Y_{ij} - x_{ij}^T \beta$. Such an estimate would have mean squared error $O(\epsilon)$ by Theorem 1. As a result the mean squared error for all parameters of interest is $O(\epsilon)$.

Our third result shows that the estimate of β obtained in step 3 is consistent. We do so by showing that the estimators $\hat{\beta}_{RLS}$ and $\hat{\beta}_{CLS}$ are consistent when

constructed using consistent variance component estimates. We give the version for $\hat{\beta}_{\text{RLS}}$.

Theorem 6. *Let $\hat{\beta}_{\text{RLS}}$ be computed with $\hat{\sigma}_A^2 \xrightarrow{P} \sigma_A^2$ and $\hat{\sigma}_E^2 \xrightarrow{P} \sigma_E^2$ as $N \rightarrow \infty$, where $\sigma_E^2 > 0$. If $\max(\epsilon_R, \epsilon_C) \rightarrow 0$ and,*

$$\mathcal{I}_0 \left(\sum_{ij} \frac{Z_{ij}(x_{ij} - \bar{x}_{i\bullet})(x_{ij} - \bar{x}_{i\bullet})^\top}{N} \right) \geq c > 0 \quad (5.2)$$

and

$$\frac{1}{R^2} \sum_{ir} (ZZ^\top)_{ir} N_{i\bullet}^{-1} N_{r\bullet}^{-1} \rightarrow 0, \quad (5.3)$$

then $\hat{\beta}_{\text{RLS}} \xrightarrow{P} \beta$.

Proof. See Section S3.3 in the Supplementary Material.

The most complicated part of the proof of Theorem 6 involves handling the contribution of b_j to $\hat{\beta}_{\text{RLS}}$. In a row-weighted GLS it is quite standard to have random errors a_i and e_{ij} but here we must also contend with errors b_j that do not appear in the model for which $\hat{\beta}_{\text{RLS}}$ is the MLE. Condition (5.3) is used to control the variance contribution of the column random effects to the intercept in $\hat{\beta}_{\text{RLS}}$. For balanced data it reduces to $1/C \rightarrow 0$ and so it has an effective number of columns interpretation. Recalling that $(ZZ^\top)_{ir}$ is the number of columns sampled in both rows i and r , we have $(ZZ^\top)_{ir} \leq N_{r\bullet}$ and so a sufficient condition for (5.3) is that $(1/R) \sum_i N_{i\bullet}^{-1} \rightarrow 0$. For sparsely observed data we expect $(ZZ^\top)_{ir} \ll \max(N_{i\bullet}, N_{r\bullet})$ to be typical, in which case, these bounds are conservative.

Any realistic setting will have $\sigma_E^2 > 0$, which we need for $\hat{\beta}_{\text{RLS}}$ to be well defined, and so that condition in Theorem 6 is not restrictive.

It remains to show that the variance component estimates from step 4 are consistent. We can just apply Theorem 5 again. Therefore the final estimates returned by Algorithm 1 are consistent given only weak conditions on the behavior of Z_{ij} and x_{ij} .

5.2. Asymptotic normality of $\hat{\beta}_{\text{RLS}}$

Here we show that the estimator $\hat{\beta}_{\text{RLS}}$ constructed using consistent estimates of σ_A^2 , σ_B^2 , and σ_E^2 is asymptotically Gaussian. The same result applies to $\hat{\beta}_{\text{CLS}}$ after transposing the conditions. We need stronger conditions than we needed

for consistency.

These conditions are expressed in terms of some weighted means of the predictors. First, let

$$\tilde{x}_{\bullet,j} = \frac{1}{N_{\bullet,j}} \sum_i Z_{ij} \frac{\sigma_A^2}{\sigma_A^2 + \sigma_E^2/N_{i\bullet}} \bar{x}_{i\bullet}. \tag{5.4}$$

This is a ‘second-order’ average of x for column j : it is the average over rows i that intersect j , of averages $\bar{x}_{i\bullet}$ shrunken towards zero. For a balanced design with $Z_{ij} = 1_{i \leq R} 1_{j \leq C}$ we would have $\tilde{x}_{\bullet,j} = \bar{x}_{\bullet\bullet} \sigma_A^2 / (\sigma_A^2 + \sigma_E^2/C)$, in which case, the second-order means would all be very close to $\bar{x}_{\bullet\bullet}$ for large C . Apart from the shrinkage, we can think of $\tilde{x}_{\bullet,j}$ as a local version of $\bar{x}_{\bullet\bullet}$ appropriate to column j . Next let

$$k = \sum_j \frac{N_{\bullet,j}^2 (\bar{x}_{\bullet,j} - \tilde{x}_{\bullet,j})}{\sum_j N_{\bullet,j}^2} \in \mathbb{R}^p. \tag{5.5}$$

This is a weighted sum of adjusted column means, weighted by the squared column size. The intercept component of this k will not be used.

Theorem 7. *Let $\hat{\beta}_{\text{RLS}}$ be computed with $\hat{\sigma}_A^2 \xrightarrow{p} \sigma_A^2$, $\hat{\sigma}_B^2 \xrightarrow{p} \sigma_B^2$, and $\hat{\sigma}_E^2 \xrightarrow{p} \sigma_E^2 > 0$ as $N \rightarrow \infty$. Suppose also that*

$$\begin{aligned} & \mathcal{I} \left(\sum_i \bar{x}_{i\bullet} \bar{x}_{i\bullet}^\top \right), \quad \mathcal{I}_0 \left(\sum_{ij} Z_{ij} (x_{ij} - \bar{x}_{i\bullet}) (x_{ij} - \bar{x}_{i\bullet})^\top \right), \quad \text{and} \\ & \frac{\mathcal{I}_0 \left(\sum_j N_{\bullet,j}^2 (\bar{x}_{\bullet,j} - \tilde{x}_{\bullet,j} - k) (\bar{x}_{\bullet,j} - \tilde{x}_{\bullet,j} - k)^\top \right)}{\max_j N_{\bullet,j}^2} \end{aligned}$$

all tend to infinity, where $\tilde{x}_{\bullet,j}$ is given by (5.4) and k is given by (5.5). Next for $c_j = \sum_i Z_{ij} \sigma_E^2 / (\sigma_E^2 + \sigma_A^2 N_{i\bullet})$ and $c_{ij} = \sigma_E^2 / (\sigma_E^2 + \sigma_A^2 N_{i\bullet})$ assume that both $\max_j c_j^2 / \sum_j c_j^2$ and $\max_{ij} c_{ij}^2 / \sum_{ij} c_{ij}^2$ converge to zero. Then $\hat{\beta}_{\text{RLS}}$ is asymptotically distributed as

$$\mathcal{N}(\beta, (X^\top V_A^{-1} X)^{-1} X^\top V_A^{-1} V_R V_A^{-1} X (X^\top V_A^{-1} X)^{-1}). \tag{5.6}$$

Proof. See Section S3.4 in the Supplementary Material.

The statement that $\hat{\beta}_{\text{RLS}}$ has asymptotic distribution $\mathcal{N}(\beta, V)$ is shorthand for $V^{-1/2}(\hat{\beta} - \beta) \xrightarrow{p} \mathcal{N}(0, I_p)$.

Theorem 7 imposes three information criteria. First, the R rows i with

$N_{i\bullet} > 0$ must have sample average $\bar{x}_{i\bullet}$ vectors with information tending to infinity. It would be reasonable to expect that information to be proportional to R and also reasonable to require $R \rightarrow \infty$ for a CLT. Next, the sum of within row sums of squares and cross-products of row-centered x_{ij} must have growing information, apart from the intercept term. Finally, thinking of $\bar{x}_{\bullet j} - \tilde{x}_{\bullet j}$ as the locally centered mean for column j , those quantities centered on the vector k must have a weighted sum of squares that is not dominated by any single column when weights proportional to $N_{\bullet j}^2$ are applied.

The conditions on c_j and c_{ij} are used to show that the CLT will apply to the intercept in the regression. The condition on $\max_j c_j^2 / \sum_j c_j^2$ will fail if for example column $j = 1$ has half of the N observations, all in rows of size $N_{i\bullet} = 1$. In the case of an $R \times C$ grid $\max_j c_j^2 / \sum_j c_j^2 = 1/C$ and so we can interpret this condition as requiring a large enough effective number $\sum_j c_j^2 / \max_j c_j^2$ of columns in the data.

The condition on $\max_{ij} c_{ij}^2 / \sum_{ij} c_{ij}^2$ will fail if for example the data contain a full $R \times C$ grid of values plus a single observation with $i = R+1$ and $j = C+1$. The problem is that in a row based regression, a single small row can have outsized leverage. This can be controlled by excluding relatively small rows. This pruning of rows is only used to apply the CLT to the intercept term. It is not needed for other components of β nor is it needed for consistency. We do not know if it is necessary for the CLT.

5.3. Computing $\text{Var}(\hat{\beta}_{\text{RLS}})$

Here we show how to compute the estimate of the asymptotic variance of $\hat{\beta}_{\text{RLS}}$ from Theorem 7. First,

$$\begin{aligned} & (X^T V_A^{-1} X)^{-1} X^T V_A^{-1} V_R V_A^{-1} X (X^T V_A^{-1} X)^{-1} \\ &= (X^T V_A^{-1} X)^{-1} X^T V_A^{-1} (V_A + \sigma_B^2 B_R) V_A^{-1} X (X^T V_A^{-1} X)^{-1} \\ &= (X^T V_A^{-1} X)^{-1} + (X^T V_A^{-1} X)^{-1} X^T V_A^{-1} \sigma_B^2 B_R V_A^{-1} X (X^T V_A^{-1} X)^{-1} \\ &= (X^T V_A^{-1} X)^{-1} + (X^T V_A^{-1} X)^{-1} \text{Var}(X^T V_A^{-1} b) (X^T V_A^{-1} X)^{-1}, \end{aligned} \tag{5.7}$$

where b is the length- N vector of column random effects for each observation. That is b_j appears $N_{\bullet j}$ times in b .

Using the Woodbury formula we find that $\text{Var}(X^T V_A^{-1} b)$ is equal to

$$\frac{\sigma_B^2}{\sigma_E^4} \sum_j \left(X_{\bullet j} - \sigma_A^2 \sum_i Z_{ij} \frac{X_{i\bullet}}{\sigma_E^2 + \sigma_A^2 N_{i\bullet}} \right) \left(X_{\bullet j} - \sigma_A^2 \sum_i Z_{ij} \frac{X_{i\bullet}}{\sigma_E^2 + \sigma_A^2 N_{i\bullet}} \right)^T. \tag{5.8}$$

Recall that $X_{i\bullet}$ and $X_{\bullet j}$ are row and column totals, not means.

In practice, we substitute consistent estimates $\hat{\sigma}_A^2$, $\hat{\sigma}_B^2$, and $\hat{\sigma}_E^2$ in for σ_A^2 , σ_B^2 , and σ_E^2 , respectively, in (5.7) and (5.8). We already have $(X^\top \hat{V}_A^{-1} X)^{-1}$ as well as $N_{i\bullet}$ and $X_{i\bullet}$ for $i = 1, \dots, R$ available from computing $\hat{\beta}_{\text{RLS}}$. In a new pass over the data, we compute $X_{\bullet j}$ and $\sum_i Z_{ij} X_{i\bullet}$ for $j = 1, \dots, C$, incurring $O(Np)$ computational and $O(Cp)$ storage costs. Then, (5.8) can be found in $O(Cp^2)$ time; a final step finds (5.7) in $O(p^3)$ time. Overall, estimating the variance of $\hat{\beta}_{\text{RLS}}$ requires $O(Np + Cp^2 + p^3)$ additional computation time and $O(Cp)$ additional space.

6. Comparisons with the MLE

Here we compare Algorithm 1 to maximum likelihood for a linear mixed effects model, considering both computational efficiency and statistical efficiency. We use a state-of-the-art code for linear mixed models called MixedModels (Bates (2016)). This is written in Julia (Bezanson et al. (2017)) and is much faster than other linear mixed model code we have tried.

Our examples use $R = C = 2\sqrt{N}$ for various N . We create an $R \times C$ matrix of Z_{ij} and randomly choose exactly $RC/4$ components to be one. We have an intercept and p other x 's with $x_{ij,t} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$, for $2 \leq t \leq p + 1$. We use all $p \in \{1, 5, 10, 20\}$. We take $\sigma_A^2 = 2$, $\sigma_B^2 = 1/2$, $\sigma_E^2 = 1$ and all $\beta_j = 1$. Our simulated random effects and our noise are all Gaussian because we are comparing to code that computes a Gaussian MLE.

6.1. Computational cost

The Julia package MixedModels uses a derivative-free optimization method from the BOBYQA package (Powell (2009)). At each iteration it evaluates the log-likelihood at a set of points, fits a quadratic function to those points and minimizes the quadratic. The number of likelihood evaluations per iteration is fixed, but we are unable to model the number of iterations required. We consider the cost per likelihood evaluation next.

The log-likelihood is

$$(Y - X\beta)^\top (\sigma_A^2 A_R + \sigma_B^2 B_R + \sigma_E^2 I_N)^{-1} (Y - X\beta) + \ln |\sigma_A^2 A_R + \sigma_B^2 B_R + \sigma_E^2 I_N|.$$

In an analysis using the Woodbury formula we find that the log-likelihood can be computed in $O(R^3 + \sum_i N_{i\bullet}^2)$ time. Because $1 \leq R \leq N$ we can write $R = N^\alpha$

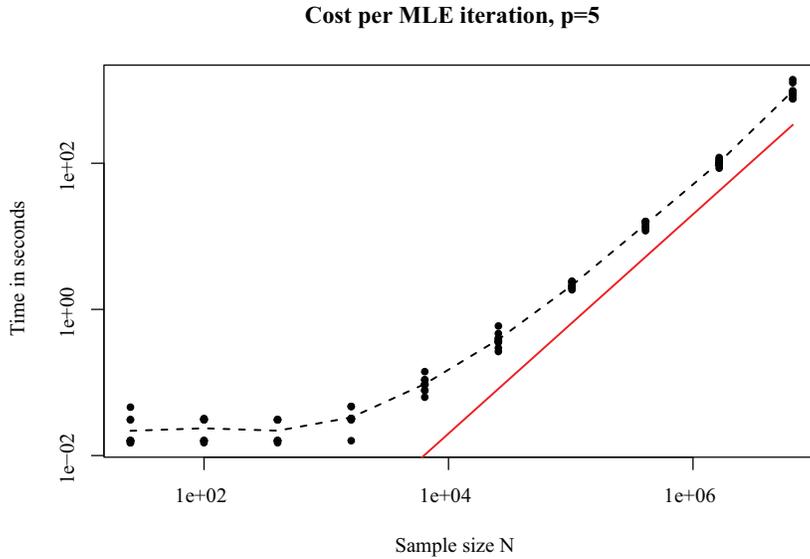


Figure 2. For $p = 5$, the cost per iteration of the MLE. The dashed curve is the average of 10 replicates. The solid reference line is parallel to $N^{3/2}$.

for some $0 \leq \alpha \leq 1$. Then

$$R^3 + \sum_i N_{i\bullet}^2 = R^3 + R \left(\frac{1}{R} \sum_i N_{i\bullet} \right)^2 \geq R^3 + N^2 R^{-1} = N^{3\alpha} + N^{2-\alpha}.$$

Now $N^{3\alpha} + N^{2-\alpha} > \max(N^{3\alpha}, N^{2-\alpha})$ and $\alpha = 1/2$ minimizes $\max(3\alpha, 2 - \alpha)$. Therefore $R^3 + \sum_i N_{i\bullet}^2 > N^{3/2}$.

This is the same estimate that we would obtain by considering the cost of solving a system of $R + C$ equations in $R + C$ unknowns. There are faster ways to solve the equations in special cases (e.g., for nested models), and there is the possibility that sparsity patterns in the data can be exploited for speed, as is done in MixedModels (Bates et al. (2015)). However, we are interested in arbitrarily complicated sampling plans where these special cases cannot be assumed.

Figure 2 shows the computed cost per iteration for 10 replicates at each of 11 different sample sizes $R^2/4$ given by $R = 10, 20, 40, \dots, 2^9 \times 10$, with $p = 5$. The cost per iteration is flat for small N presumably due to some overhead. It grows slowly until about $N = 10^4$ and then it appears to increase parallel to a reference line with slope $3/2$.

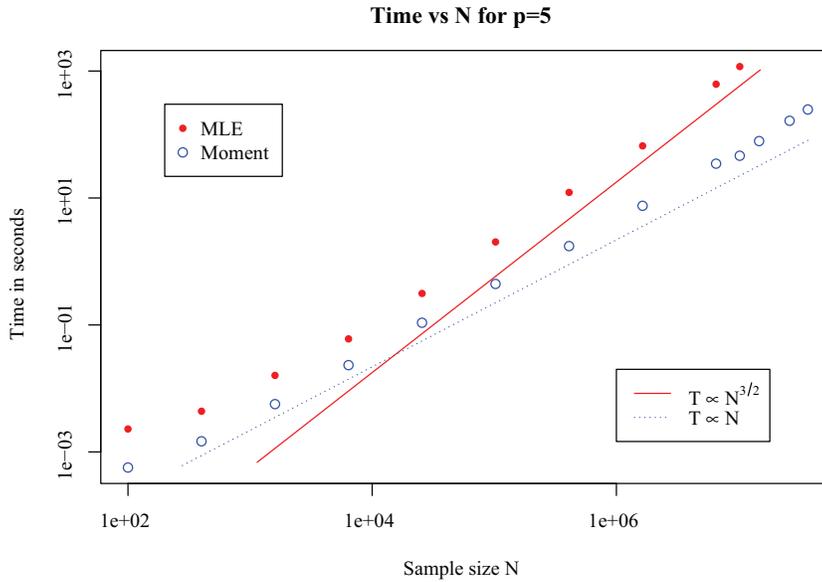


Figure 3. Computational cost for MLE and moments versus sample size N . There are reference lines parallel to $N^{3/2}$ and N^1 .

Figure 3 shows total cost versus N in a setting with $p = 5$ averaged over 100 data sets. The cost curve for the MLE computation looks different to that in Figure 2. It does not start out flat for small N . We found that the number of iterations required to find the MLE generally rose over the range $64 \leq N \leq 6,400$ and then declined gently thereafter, making it more difficult to discern the $O(N^{3/2})$ rate. Because the number of iterations cannot be below one, we can be sure that the MLE cost is at least a multiple of $N^{3/2}$.

From the analysis and empirical results, we find that a cost per iteration of $O(N^{3/2})$ is a realistic lower bound for the MLE code. The method of moments cost is $O(N)$ theoretically and appears to be proportional to N empirically.

Our computations were performed using data generated in memory. In commercial applications, there could be a much larger time cost proportional to N related to reading the data from external storage. However, the $N^{3/2}$ cost component would be considerably larger at commercial scale, where N is much larger than in our examples. For the method of moments it is straightforward to read and use the data in parallel even for large N .

A second computational issue arises with the linear mixed effects MLE. The code crashes on large enough data sets because the algorithm requires

$O((R + C)^2)$ memory. For $p = 5$ we were unable to take the next step past $N = 5,120^2/4 \doteq 6.5 \times 10^6$. The program runs out of memory on our cluster where we had one virtual machine with 4Gb memory. For $p = 1$, it crashes for N near 18 million observations. The method of moments in Algorithm 1 has linear cost both theoretically and empirically and can be implemented in $O(R + C)$ memory. The difference is minor for our CPU time simulations that also keep all N observations in memory, but it will be critical in large commercial applications. Commercial computing resources are much greater than ours but their data are vastly greater than our sample.

6.2. Statistical efficiency

For statistical efficiency we considered $p = 1, 5, 10$ and 20 . Sample sizes $N = 100 \times 4^j$ for $j = 0, 1, \dots, 8$ were replicated 100 times each. A few larger values of N were replicated 10 times each, though the MLE code would not run on all of the largest sample sizes we tried. The pattern in the results was the same for all of those p . We display results for $p = 5$ in Figure 4. The MSEs for β decay proportionally to $1/N$. The reference curves for variance components in Figure 4 are what we would expect from i.i.d. sampling of a_i , b_j and e_{ij} , respectively: $2\sigma_A^4/R$, $2\sigma_B^4/C$ and $2\sigma_E^4/N$ where $R = C = 2\sqrt{N}$.

The parameter of greatest interest will ordinarily be β . The MLE has greater accuracy for β , as it must by the Gauss-Markov theorem. In this instance the MLE has about half the MSE that the method of moments does. For the variance components, the method of moments attains essentially the same MSE as the MLE does for σ_A^2 and σ_B^2 . The MLE has greater efficiency for σ_E^2 . In ordinary use we would want to know ratios of variance components and the uncertainty in such ratios is dominated by that in σ_A^2 and σ_B^2 , where the two methods have comparable accuracy.

In this example, we saw a modest loss in statistical efficiency of $\hat{\beta}$ and $\hat{\sigma}_E^2$ and comparable accuracy for $\hat{\sigma}_A^2$ and $\hat{\sigma}_B^2$. These comparisons were run on data simulated from the Gaussian model that the MLE assumes. The method of moments does not require that assumption. Likelihood based variance estimates for variance components, such as $\widehat{\text{Var}}(\hat{\sigma}_A^2)$, can fail to be even asymptotically correct when the Gaussian model does not hold.

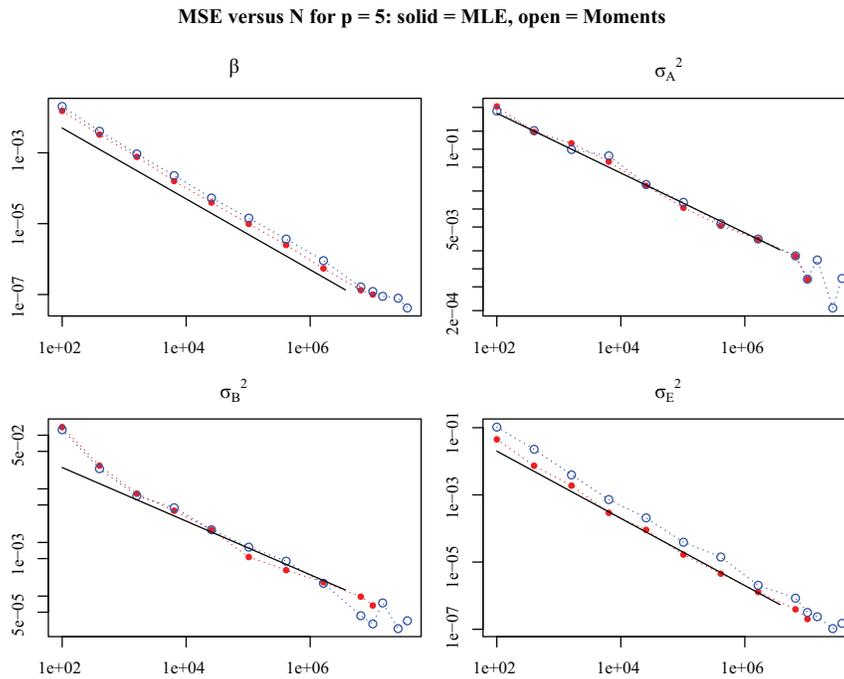


Figure 4. Mean squared errors for β , σ_A^2 , σ_B^2 and σ_E^2 versus N . Reference lines for β and σ_E^2 are parallel to N^{-1} . Reference lines for σ_A^2 and σ_B^2 are parallel to $N^{-1/2}$.

7. Discussion

We have proposed an algorithm for the two-factor linear mixed effects model with a crossed covariance structure that provides consistent and asymptotically normal parameter estimates. It alternates twice between estimating the regression coefficients and estimating the variance components via the method of moments.

Unlike available methods based on Bayes' theorem or maximum likelihood, the moment estimates cost $O(N)$ time and $O(R+C)$ space. The variance estimate for $\hat{\beta}$ is obtained by substituting consistent estimates of σ_A^2 , σ_B^2 , and σ_E^2 into exact finite sample formulae for that variance matrix. The variance estimates for $\hat{\sigma}_A^2$, $\hat{\sigma}_B^2$, and $\hat{\sigma}_E^2$ are obtained by such a substitution in mildly conservative formulae from Gao and Owen (2017). Here the usual root- n consistency from i.i.d. settings is replaced by a $1/\sqrt{\epsilon}$ -consistency for $\epsilon = \max(\epsilon_R, \epsilon_C)$. Interpreting $1/\sqrt{\epsilon}$ as an effective sample size might be somewhat conservative because in theorems such as Theorem 1 the value of ϵ appears in upper bounds.

We exchange higher MSEs for an algorithm with cost only linear in the

number of observations. We do not know how bad the efficiency loss might be in general, but we expect that when the pure error term σ_E^2 is meaningfully large, the loss will not be extreme. Also, if one of σ_A^2 and σ_B^2 very much dominates the other one, we can get a GLS estimate that accounts for the dominant source of correlation.

Gao (2017) proves a martingale central limit theorem for the variance component estimates $\hat{\sigma}_A^2$, $\hat{\sigma}_B^2$, and $\hat{\sigma}_E^2$. We do not anticipate those variance components to be uncorrelated with $\hat{\beta}$ because the random variables a_i , b_j , and e_{ij} might not have symmetric distributions.

This study is a second step in developing big-data versions of mixed model procedures such as the Henderson estimators. One follow-up step is to incorporate interactions between fixed and random steps, as the Henderson III model allows. Another is to incorporate interactions among latent variables. At present both kinds of interactions would serve to inflate σ_E^2 . A third step is to adapt to binary responses, for instance by replacing the identity link in Model 1, with a logit or probit link. This third step is of value because many responses in e-commerce are categorical, e.g., for Stitch Fix, whether the client keeps the item of clothing.

The computation for GLMMs is daunting, especially for large ones. Referring to penalized quasi-likelihood, (McCulloch, Searle and Neuhaus (2008, p.341)) write

In the “derivation” of the PQL equations quite a few approximations of undetermined accuracy are bandied about and the development has an air of ad hocery. How well do these methods work in practice? Unfortunately, not very.

The latest version of the `lme4` R package Bates (2014) does not include their previous `mcmcSamp` method because it was deemed to be unreliable. Jiang Jiang (1998) proposes a general method of simulated moments estimator for generalized linear mixed models by deriving sufficient statistics, but they have superlinear computational complexity at equation (16) in our context. Even the theory of GLMMs is difficult. The consistency of the maximum likelihood estimate of μ in a balanced data set for a binary response Y_{ij} with $\text{logit}(\Pr(Y_{ij} = 1 | a_i, b_j)) = \mu + a_i + b_j$ was only established in Jiang (2013).

Finally, Papaspiliopoulos, Roberts and Zanella (2018) has recently shown that for sparse observation patterns, the convergence time of a collapsed Gibbs sampler for an intercept-only version of Model 1, alternating between sampling

μ and the a_i 's and μ and the b_j 's, would not grow with the size of the data set. This suggests that a reparameterization could enable MCMC to have good performance on our model as well. That paper does however require a strong balance assumption in which all rows are equally commonly represented in the data and similarly for the columns. That assumption is very unrealistic for e-commerce.

Supplementary Material

The proofs of our results are provided in the an online Supplementary Material, at: <http://statweb.stanford.edu/~owen/reports/vllmemsupp.pdf>.

Acknowledgments

This work was supported by the US NSF under grants DMS-1407397 and DMS-1521145. KG was supported by US NSF Graduate Research Fellowship under grant DGE-114747. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation.

We would like to thank Stitch Fix and in particular Brad Klingenberg for providing us with the data used in our real-world experiment and motivation and encouragement during the project.

References

- Bates, D. (2014). Computational methods for mixed models. Technical report, Department of Statistics, University of Wisconsin–Madison. <https://cran.r-project.org/web/packages/lme4/vignettes/Theory.pdf>.
- Bates, D. (2016). Linear mixed-effects models in Julia. <https://github.com/dmbates/MixedModels.jl>.
- Bates, D., Mächler, M., Bolker, B. and Walker, S. (2015). Fitting linear mixed-effects models using lme4. *Journal of Statistical Software* **67**, 1–48.
- Bennett, J. and Lanning, S. (2007). The Netflix prize. In *Proceedings of KDD Cup and Workshop*.
- Bezanson, J., Edelman, A., Karpinski, S. and Shah, V. B. (2017). Julia: A fresh approach to numerical computing. *SIAM Review* **59**, 65–98.
- Blei, D. M., Kucukelbir, A. and McAuliffe, J. D. (2017). Variational inference: A review for statisticians. *Journal of the American Statistical Association* **112**, 859–877.
- Demidenko, E. (2013). *Mixed Models: Theory and Applications with R*. John Wiley & Sons.
- 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 (Statistical Methodology)* **39**, 1–38.

- Fox, C. (2013). Polynomial accelerated MCMC and other sampling algorithms inspired by computational optimization in Dick, J., Kuo, F. Y., Peters, G. W. and Sloan, I. H., editors, *Monte Carlo and Quasi-Monte Carlo Methods 2012*, 349–366. Springer, Berlin.
- Gao, K. (2017). *Scalable Estimation and Inference for Massive Linear Mixed Models with Crossed Random Effects*. PhD thesis, Stanford University.
- Gao, K. and Owen, A. B. (2017). Efficient moment calculations for variance components in large unbalanced crossed random effects models. *Electronic Journal of Statistics* **11**, 1235–1296.
- Gelman, A., Van Dyk, D. A., Huang, Z. and Boscardin, J. W. (2012). Using redundant parameterizations to fit hierarchical models. *Journal of Computational and Graphical Statistics* **17**.
- Hager, W. W. (1989). Updating the inverse of a matrix. *SIAM Review* **31**, 221–239.
- Hartley, H. O. and Rao, J. N. K. (1967). Maximum-likelihood estimation for the mixed analysis of variance model. *Biometrika* **54**, 93–108.
- Henderson, C. R. (1953). Estimation of variance and covariance components. *Biometrics* **9**, 226–252.
- Imbens, G. W. and Rubin, D. B. (2015). *Causal Inference in Statistics, Social, and Biomedical Sciences*. Cambridge University Press, Cambridge.
- Ioannidis, J. P. A. (2005). Why most published research findings are false. *PLoS Medicine* **2**, e124.
- Jiang, J. (1998). Consistent estimators in generalized linear mixed models. *Journal of the American Statistical Association* **93**, 720–729.
- Jiang, J. (2007). *Linear and Generalized Linear Mixed Models and Their Applications*. Springer Science & Business Media.
- Jiang, J. (2013). The subset argument and consistency of MLE in GLMM: Answer to an open problem and beyond. *The Annals of Statistics* **41**, 177–195.
- Marshall, A. W. and Olkin, I. (1990). Matrix versions of the Cauchy and Kantorovich inequalities. *Aequationes Mathematicae* **40**, 89–93.
- McCulloch, C. E., Searle, S. R. and Neuhaus, J. M. (2008). *Generalized, Linear, and Mixed Models*. John Wiley & Sons, Hoboken, NJ.
- Papaspiliopoulos, O., Roberts, G. O. and Zanella, G. (2018). Scalable inference for crossed random effects models. *arXiv:1803.09460*.
- Perry, P. O. (2017). Fast moment-based estimation for hierarchical models. *Journal of the Royal Statistical Society, Series B (Statistical Methodology)* **79**, 267–291.
- Powell, M. J. D. (2009). The BOBYQA Algorithm for Bound Constrained Optimization without Derivatives. Technical Report NA2009/06, University of Cambridge.
- Provost, F. and Fawcett, T. (2013). Data science and its relationship to big data and data-driven decision making. *Big Data* **1**, 51–59.
- Raudenbush, S. W. (1993). A crossed random effects model for unbalanced data with applications in cross-sectional and longitudinal research. *Journal of Educational and Behavioral Statistics* **18**, 321–349.
- Roberts, G. O. and Sahu, S. K. (1997). Updating schemes, correlation structure, blocking and parameterization for the Gibbs sampler. *Journal of the Royal Statistical Society, Series B (Statistical Methodology)* **59**, 291–317.
- Scott, S. L., Blocker, A. W., Bonassi, F. V., Chipman, H. A., George, E. I. and McCulloch, R. E. (2016). Bayes and big data: The consensus Monte Carlo algorithm. *International*

- Journal of Management Science and Engineering Management* **11**, 78–88.
- Searle, S. R., Casella, G. and McCulloch, C. E. (1992). *Variance Components*. John Wiley & Sons, New York.
- Varian, H. R. (2014). Big data: New tricks for econometrics. *The Journal of Economic Perspectives* **28**, 3–27.
- Wasserstein, R. L. and Lazar, N. A. (2016). The ASA’s statement on p -values: context, process, and purpose. *American Statistician* **70**, 129–133.
- Wu, P., Stute, W. and Zhu, L.-X. (2012). Efficient estimation of moments in linear mixed models. *Bernoulli* **18**, 206–228.
- Wu, P. and Zhu, L. X. (2010). An orthogonality-based estimation of moments for linear mixed models. *Scandinavian Journal of Statistics* **37**, 253–263.
- Yu, Y. and Meng, X.-L. (2011). To center or not to center: That is not the question – an ancillarity–sufficiency interweaving strategy (ASIS) for boosting MCMC efficiency. *Journal of Computational and Graphical Statistics* **20**, 531–570.

Intel Labs, Santa Clara, CA, 95054, USA.

E-mail: katelyng3@gmail.com

Dept of Statistics, Sequoia Hall, Stanford University, Stanford, CA 94305, USA.

E-mail: owen@stanford.edu

(Received January 2018; accepted October 2018)