

Weakly Informative Prior for Point Estimation of Covariance Matrices in Hierarchical Models

Yeojin Chung
Kookmin University

Andrew Gelman
Columbia University

Sophia Rabe-Hesketh
University of California

Jingchen Liu
Columbia University

Vincent Dorie
New York University

When fitting hierarchical regression models, maximum likelihood (ML) estimation has computational (and, for some users, philosophical) advantages compared to full Bayesian inference, but when the number of groups is small, estimates of the covariance matrix (Σ) of group-level varying coefficients are often degenerate. One can do better, even from a purely point estimation perspective, by using a prior distribution or penalty function. In this article, we use Bayes modal estimation to obtain positive definite covariance matrix estimates. We recommend a class of Wishart (not inverse-Wishart) priors for Σ with a default choice of hyperparameters, that is, the degrees of freedom are set equal to the number of varying coefficients plus 2, and the scale matrix is the identity matrix multiplied by a value that is large relative to the scale of the problem. This prior is equivalent to independent gamma priors for the eigenvalues of Σ with shape parameter 1.5 and rate parameter close to 0. It is also equivalent to independent gamma priors for the variances with the same hyperparameters multiplied by a function of the correlation coefficients. With this default prior, the posterior mode for Σ is always strictly positive definite. Furthermore, the resulting uncertainty for the fixed coefficients is less underestimated than under classical ML or restricted maximum likelihood estimation. We also suggest an extension of our method that can be used when stronger prior information is available for some of the variances or correlations.

Keywords: *Bayes modal estimation; penalized likelihood estimation; variance estimation; Heywood case; mixed-effects model; multilevel model*

Hierarchical or mixed-effects regression models are increasingly popular in applied statistics and can be viewed as Bayesian at the following two levels: A prior distribution is assigned to the varying coefficients, and the parameters of that prior distribution themselves are given a hyperprior. The family of models can be written in general terms as follows: Data are in groups $j = 1, \dots, J$. For each group j , there is a response vector \mathbf{y}_j and two data matrices, X_j and Z_j , that have fixed and varying coefficients, respectively. The data model is $p(\mathbf{y}_j | X_j \boldsymbol{\beta} + Z_j \mathbf{b}_j)$, where $\boldsymbol{\beta}$ is the vector of fixed coefficients and \mathbf{b}_j is the vector of regression coefficients that varies by group. The vectors \mathbf{b}_j are modeled as independent draws from a prior distribution, $p(\mathbf{b}_j)$, given some hyperparameters. We shall assume a normal model for the varying coefficients, so that $\mathbf{b}_j \sim N(\mathbf{0}, \Sigma)$. The model could also include a nonzero mean vector or a group-level regression structure for the hyperprior distribution, but these can be folded into the fixed coefficients in the data model without loss of generality.

There is a rich literature on full Bayesian inference for hierarchical regressions. There is also an empirical Bayes version in which the hyperparameters (in this case, Σ) are estimated via maximum likelihood (ML) and then inference for the coefficients is performed conditional on the estimated Σ . From the Bayesian perspective, the empirical Bayes approach is suboptimal, both because it avoids the use of any prior information on Σ and because it understates posterior uncertainty. From a pragmatic perspective, however, we recognize that the point estimation approach has two advantages that give it great appeal to many users. First, existing software such as `lme4` in R and various commands in Stata allow such models to be fit fast and reliably for moderate-sized data sets, whereas software for Markov chain Monte Carlo simulation for full Bayes inference is not yet so immediately practical. Second, the non-Bayesian motivation behind point estimation is attractive to practitioners who want the benefits of partial pooling and hierarchical modeling without needing to specify prior information or fully buy into the Bayesian paradigm.

The subject of this article is the use of Bayesian ideas and methods to produce better inferences for hierarchical models via better point estimates of the hyperparameters. In that sense, this work falls into a long tradition of Bayesian tools used for practical non-Bayesian inferences (e.g., Agresti & Coull, 1998). Bayes modal (BM) estimation (or penalized likelihood) has also been used to obtain more stable estimates in item response theory (e.g., Mislevy, 1986; Swaminathan & Gifford, 1985; Tsutakawa & Lin, 1986) and to avoid boundary estimates (or logit parameters tending to $\pm\infty$) in log-linear models (Galindo-Garre, Vermunt, & Bergsma, 2004), logistic regression (Gelman, Jakulin, Pittau, & Su, 2008), varying-intercept models with constant coefficients (Chung, Rabe-Hesketh, Dorie, Gelman, & Liu, 2013), random-effects

meta-analysis models (Chung, Rabe-Hesketh, & Choi, 2013), and latent class analysis (Galindo-Garre & Vermunt, 2006; Maris, 1999). Such an approach has also been used to obtain nondegenerate covariance matrices in factor analysis (Martin & McDonald, 1975), in finite mixtures of normal densities (Ciuperca, Ridolfi, & Idier, 2003; Vermunt & Magidson, 2005), and in multivariate regression (Warton, 2008). In varying intercept models, the Stein loss function (Srivastava & Kubokawa, 1999) and an extension of MANOVA estimation (Amemiya, 1985) have been used for obtaining nonnegative definite covariance estimators.

The key problem solved by our method is the tendency of ML estimates of Σ to be degenerate, that is, on the border of positive definiteness, which corresponds to zero variance or perfect correlation among some linear combinations of the parameters. When the ML estimate of a hierarchical covariance matrix is degenerate, this often arises from a likelihood that is nearly flat in the relevant dimension and just happens to have a maximum at the boundary.

Our solution is a class of weakly informative prior densities for Σ that go to zero on the boundary as Σ becomes degenerate, thus ensuring that the posterior mode (i.e., the maximum penalized likelihood estimate) is always nondegenerate. We recommend a class of Wishart priors with a default choice of hyperparameters, that is, the degrees of freedom is the dimension of \mathbf{b}_j plus 2 and the scale matrix is the identity matrix multiplied by a large enough number. This prior can be expressed as a product of gamma(1.5, θ) priors on the eigenvalues of Σ or as a product of gamma(1.5, θ) priors on variances of the varying effects with rate parameter $\theta \rightarrow 0$ and a function of the correlations (a beta prior in the two-dimensional case). In the varying-intercept model (Chung, Rabe-Hesketh, Dorie, et al., 2013) and random-effects meta-analysis model (Chung, Rabe-Hesketh, & Choi, 2013), the gamma(1.5, θ) prior successfully avoids boundary estimates while producing estimates that are consistent with the data. We show that this is also true for the default Wishart prior proposed in this article for general varying coefficient models.

In a simulation study and an education example presented later, the default Wishart prior always gives nondegenerate estimates of Σ (in particular, nonperfect correlation coefficients) without decreasing the log likelihood substantially. The BM estimators of the standard deviations and correlations using the default Wishart prior have better statistical properties than the (restricted) ML estimators.

When prior information is available for specific standard deviations or correlations, additional penalty functions may be included. Specifically, if the prior most plausible value for a standard deviation or correlation parameter is σ^* or ρ^* , respectively, then we propose multiplying the Wishart prior by the gamma(2, $2/\sigma^*$) or $N(\rho^*, .25^2)$ densities. This assigns more prior probability around the preferred values while exploiting the property of the Wishart prior that it ensures that the estimates remain positive definite.

The outline of the article is as follows. First, we illustrate the boundary estimation problems encountered in ML estimation of hierarchical variance and covariance parameters. Then, we introduce the default Wishart prior for Σ and

Chung et al.

investigate its properties. Next, additional penalty functions are proposed that incorporate further prior knowledge for some of the parameters. Finally, our method is applied to an example from education research and simulated data.

Boundary Estimation Problem

Consider the varying-coefficients model,

$$y_{ij} = \mathbf{x}_{ij}^T \boldsymbol{\beta} + \mathbf{z}_{ij}^T \mathbf{b}_j + \varepsilon_{ij}, i = 1, \dots, n_j, j = 1, \dots, J, \quad (1)$$

where y_{ij} is the response variable for unit i in group j , \mathbf{x}_{ij} is a p -dimensional covariate vector with constant (or fixed) coefficients $\boldsymbol{\beta}$, \mathbf{z}_{ij} is a d -dimensional covariate vector with varying coefficients $\mathbf{b}_j \sim N(0, \Sigma)$, and $\varepsilon_{ij} \sim N(0, \sigma_\varepsilon^2)$ is a residual for each observation. We further assume that \mathbf{b}_j and ε_{ij} are independent of each other and of the covariates (and suppress conditioning on covariates throughout the paper).

Non-Bayesian Point Estimation

For each j , $\mathbf{y}_j = (y_{1j}, \dots, y_{n_j j}) \sim N(X_j \boldsymbol{\beta}, V_j)'$, where X_j is a $n_j \times p$ matrix with \mathbf{x}_{ij}^T in the i th row, $V_j = Z_j \Sigma Z_j^T + \sigma_\varepsilon^2 I$, and Z_j is a $n_j \times d$ matrix with \mathbf{z}_{ij}^T in the i th row. The log-likelihood function is given by:

$$\log p(\mathbf{y} | \boldsymbol{\beta}, \Sigma, \sigma_\varepsilon^2) = -\frac{1}{2} \left[\sum_{j=1}^J \log |V_j| + \sum_j (\mathbf{y}_j - X_j \boldsymbol{\beta})^T V_j^{-1} (\mathbf{y}_j - X_j \boldsymbol{\beta}) \right], \quad (2)$$

where the constant term, $-(N/2)\log(2\pi)$, has been dropped. The ML estimator is obtained by maximizing the log-likelihood function.

It is known that the ML estimator of the covariance matrix is biased for finite samples (Lehmann & Casella, 1998), and an often-preferred option is restricted maximum likelihood (REML; Patterson & Thompson, 1971), as it takes into account the degrees of freedom for the fixed coefficients $\boldsymbol{\beta}$. Harville (1974) showed that the REML estimator can be derived by specifying flat prior distributions for $\boldsymbol{\beta}$, marginalizing over $\boldsymbol{\beta}$, and maximizing the marginal (or restricted) likelihood with respect to Σ and σ_ε^2 . The restricted log-likelihood function is given by:

$$\log p_R(\mathbf{y} | \Sigma, \sigma_\varepsilon^2) = -\frac{1}{2} \left[\log \left| \sum_{j=1}^J X_j^T V_j^{-1} X_j \right| + \sum_{j=1}^J \log |V_j| + \sum_{j=1}^J (\mathbf{y}_j - X_j \hat{\boldsymbol{\beta}})^T V_j^{-1} (\mathbf{y}_j - X_j \hat{\boldsymbol{\beta}}) \right], \quad (3)$$

up to a constant, where

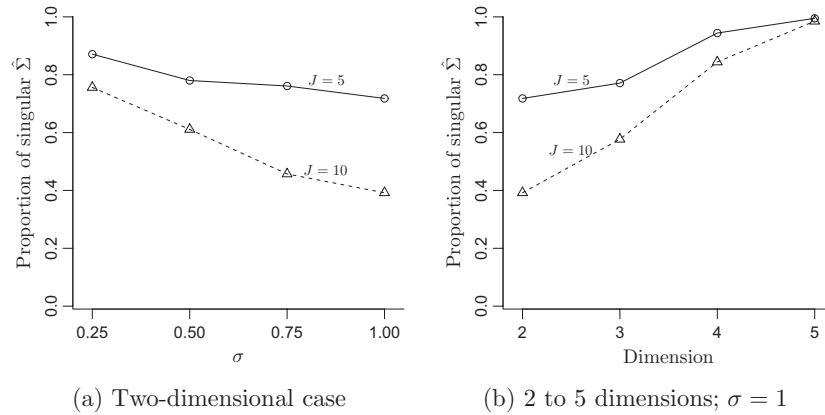


FIGURE 1. Proportion of data sets, out of 1,000, where the maximum likelihood (ML) estimate of the covariance matrix is singular. (a) Two-dimensional case: When $\sigma = .25$, 87% of the ML estimates are singular for $J = 5$. As σ and J increase, the proportion decreases but is greater than 40% for the conditions considered. (b) Two to five dimensions $\sigma = 1$: As the dimension of Σ increases, there is a rapid increase in the probability of the estimate being degenerate.

$$\hat{\beta} = \left(\sum_{j=1}^J X_j^T V_j^{-1} X_j \right)^{-1} \left(\sum_{j=1}^J X_j^T V_j^{-1} \mathbf{y}_j \right).$$

Singular Estimates of Σ using ML and REML

ML and REML often yield singular (i.e., nonpositive definite) estimates of Σ . This boundary includes the cases where some varying coefficients have zero variance or a varying coefficient is a linear combination of the other varying coefficients.

We present two simulation studies to demonstrate how often singular estimates of Σ occur in the varying-coefficients model. In the first study, we consider a model with two-dimensional varying coefficients, that is, a varying intercept b_{0j} and a varying slope b_{1j} . We set the group size to $n = 10$ and the number of groups to $J = 5$ or 10. A covariate that varies within group only was generated from $N(0, 1)$ and group-mean centered. The varying coefficients (b_{0j}, b_{1j}) were generated from $N(\mathbf{0}, \sigma^2 I_2)$ with $\sigma = 0.25, 0.5, 0.75, 1$. Setting the correlation to 0 corresponds to the best-case scenario in the sense of being furthest from the boundary. The within-group variance σ_ε^2 was set to 1 and the fixed coefficients β_0 and β_1 were set to 0. For each of 1,000 random samples of data from the model, we obtained ML and REML estimates using lmer (Bates & Maechler, 2010) in R.

Figure 1a shows the proportion of ML estimates of Σ on the boundary for the two-dimensional case. For $J = 5$ groups, 87% of the ML estimates are singular

when $\sigma = 0.25$ and the proportion decreases as σ increases but remains as high as 72% when $\sigma = 1$. For $J = 10$ groups, the proportions are smaller than those for $J = 5$ but still, in more than 40% of the simulations, the likelihood is maximized at a singular $\hat{\Sigma}$. The REML estimator yields smaller proportions of singular estimates with a similar trend (not shown). For $J = 10$, 79% and 64% of the REML estimates are singular when $\sigma = 0.25$ and $\sigma = 1$, respectively. For $J = 10$, the proportion is reduced to 69% and 35% when $\sigma = 0.25$ and $\sigma = 1$, respectively.

Our second simulation study considers various dimensions, from $d = 2$ to $d = 5$, each time with a varying intercept and $d - 1$ varying slopes for $n = 10$ and $J = 5$ or 10. The $d - 1$ covariates were independently drawn from $N(0, 1)$ and centered at their group means as in the previous simulation. The varying coefficients \mathbf{b}_j were drawn from $N(\mathbf{0}, I_d)$ and σ_ε^2 was set to 1. Figure 1b presents the proportion of replicates where the ML estimate $\hat{\Sigma}$ is singular. As the number of dimensions increases, this proportion increases rapidly, exceeding 95% with five varying coefficients for both $J = 5$ and $J = 10$. For REML, the proportions of singular estimates are slightly lower than for ML but follow a similar pattern and exceed 35% across all simulation conditions.

In some contexts, singular estimates of the covariance matrix are acceptable or considered as an indication of structural misspecification of the model. In the varying-intercept model, a negative group-level variance estimate is sometimes permitted if the model is viewed as a marginal model for the responses, given the covariates where only the sum of the group-level and within-group variance must be positive (Verbeke & Molenberghs, 2000, pp. 52–53). In factor analysis and structural equation models, a negative variance estimate, called a Heywood case, is sometimes interpreted as model misspecification, especially if the null hypothesis that the variance is nonnegative can be rejected (Kolenikov & Bollen, 2012). However, this article takes a hierarchical perspective of the multilevel linear model, where the intercepts and slopes vary due to omitted group-level variables. Therefore, the variances of the varying coefficients must be nonnegative, and perfect correlations among linear combinations of varying coefficients are regarded as unrealistic.

Weakly Informative Wishart Prior for Σ

We propose posterior modal estimation with a prior on Σ , implicitly assuming uniform priors for the other parameters. With a prior $p(\Sigma)$, the log-posterior function can be written as follows:

$$\log p(\boldsymbol{\beta}, \Sigma, \sigma_\varepsilon | \mathbf{y}) = \log p(\mathbf{y} | \boldsymbol{\beta}, \Sigma, \sigma_\varepsilon) + \log p(\Sigma) + c, \quad (4)$$

and we find the mode of $\log p(\boldsymbol{\beta}, \Sigma, \sigma_\varepsilon | \mathbf{y})$. This approach can also be viewed as maximum penalized likelihood estimation where $\log p(\Sigma)$ is a penalty function. We consider a family of Wishart (*not* inverse-Wishart) densities for the prior on Σ . The Wishart density function on Σ with hyperparameters ν and Ψ is defined by:

$$p(\Sigma) = \frac{|\Sigma|^{(\nu-d-1)/2} \exp\left[-\frac{1}{2} \text{tr}(\Psi^{-1}\Sigma)\right]}{2^{\nu d/2} |\Psi|^{\nu/2} \Gamma_d(\nu/2)}, \nu > d - 1, \Psi > 0, \quad (5)$$

where $\Gamma_d(\nu/2) = \pi^{d(d-1)/4} \prod_{j=1}^d \Gamma(\nu/2 + (1-j)/2)$, ν is the degrees of freedom, and Ψ is a scale matrix with $E(\Sigma) = \nu\Psi$.

If we set Ψ to be a diagonal matrix $(1/2\theta)I_d$, the Wishart density of Σ in Equation 5 can be written as:

$$\begin{aligned} p(\Sigma) &= \frac{\theta^{d\nu/2}}{\Gamma_d(\nu/2)} |\Sigma|^{(\nu-d-1)/2} \exp(-\theta \text{tr}(\Sigma)) \\ &= \frac{\theta^{d\nu/2}}{\Gamma_d(\nu/2)} \prod_{r=1}^d \lambda_r^{(\nu-d-1)/2} \exp(-\theta \lambda_r) \\ &\propto \prod_{r=1}^d g\left(\lambda_r \mid \frac{\nu-d+1}{2}, \theta\right), \end{aligned} \quad (6)$$

where $\lambda_1, \dots, \lambda_d$ are the eigenvalues of Σ and $g(x|\alpha, \theta)$ is the gamma(α, θ) density with shape parameter α and rate parameter θ , $g(x|\alpha, \theta) = \frac{\theta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} \exp(-x\theta)$. In the previous equations, note that we do not transform the density of Σ to the density of eigenvalues, but just rewrite Equation 5 as a function of eigenvalues without including a Jacobian term.

As a default choice, we propose $\nu = d + 2$ and $\theta \rightarrow 0$. In practice, we can choose a sufficiently small number for θ , for example, $\theta = 10^{-4}$ or 10^{-5} . If these two values of θ lead to almost the same parameter estimates, we can consider the choice of θ to be sufficiently close to the limit 0. In order to avoid dependency on the scale of the response variable, we can also use an improper prior $|\Sigma|^{(\nu-d-1)/2}$, which is the same as the Wishart prior up to constant in the limit $\theta \rightarrow 0$ (Chung, Rabe-Hesketh, Dorie, et al., 2013). This prior is proportional to independent gamma(1.5, θ) densities of the eigenvalues as observed in Equation 6. If Σ is a diagonal matrix, this prior implies gamma(1.5, θ) priors on the diagonal elements of Σ , which is equivalent to gamma(2, θ) priors on the standard deviations when $\theta \rightarrow 0$. If Σ is not diagonal, we obtain gamma(1.5, θ) priors on the variances and a function of the correlations.

The advantage of this family of density functions is that they equal zero at the boundary—thus, the BM or penalized likelihood estimate for Σ will never be degenerate—but the densities move away from zero when Σ moves off the boundary, so that the posterior mode can be arbitrarily close to degeneracy if this is what the data demand. In contrast, various other families of models do not have these properties, making them less desirable when used for the purpose of BM point estimation. The inverse-Wishart family of density, one of the most commonly used priors for Σ in the full Bayesian inference, is also zero at the

boundary. However, it tends to assign an excessive penalty near the boundary because it is a function of Σ^{-1} and $|\Sigma|^{-1}$ while the Wishart density is a function of Σ and $|\Sigma|$.

Alternative choices of ν and θ can be considered but ν and θ larger than the default choice will make the prior more informative. This behavior might be preferable if a plausible value of Σ is available. In the next section, we suggest including additional prior information about any specific standard deviation by multiplying the default prior by an additional penalty function, which can be viewed as a special case of the Wishart prior with larger ν and θ .

Priors on the covariance matrix in the varying-coefficients model have been investigated by several authors in the context of full Bayesian modeling. Daniels and Kass (1999) investigated nonconjugate Bayesian estimation of covariance matrices in hierarchical models including an inverse-Wishart prior on covariance matrices with unknown scale and degrees of freedom and a normal prior on Fisher's z -transformed correlations. Barnard, McCulloch and Meng (2000) decomposed $\Sigma = \text{Diag}(\mathbf{s})R\text{Diag}(\mathbf{s})$ where \mathbf{s} is a vector of standard deviations and R is the correlation matrix, which is assigned marginal or jointly uniform priors. O'Malley and Zaslavsky (2005) propose a scaled inverse Wishart, a decomposition similar to that of Barnard, McCulloch, and Meng (2000) except that the central matrix R itself has an inverse-Wishart distribution rather than being constrained to be a correlation matrix. Our approach is different from these others in being explicitly intended not for full Bayes inference but as a tool to obtain positive definite posterior modal estimates. As such, our concerns are different from those involved in constructing traditional Bayesian priors.

Unlike posterior mean estimation, BM estimation does not involve simulation and is computationally as efficient as ML estimation. By modifying existing ML estimation procedures, `gllamm` (Rabe-Hesketh, Skrondal, & Pickles, 2005) in Stata and `lmer` (Bates & Maechler, 2010) in R, we have developed software to find the maximum of the penalized likelihood. The modified `gllamm` is available from www.gllamm.org and `blmer`, the modified `lmer` function, can be found in the `blme` package available from the Comprehensive R Archive Network.

Varying-Intercept Models: $d = 1$

The varying-intercept model is a special case of the model in Equation 1 with $d = 1$, given by:

$$y_{ij} = \mathbf{x}_{ij}^T \boldsymbol{\beta} + b_j + \varepsilon_{ij},$$

where $b_j \sim N(0, \sigma_b^2)$ and $\varepsilon_{ij} \sim N(0, \sigma_\varepsilon^2)$. The Wishart prior in Equation 6 is equivalent to a $\text{gamma}(\nu/2, \theta)$ prior on σ_b^2 . With the default choice of hyperparameters, $\nu = 3 (= d + 2)$ and $\theta \rightarrow 0$, the Wishart prior coincides with a $\text{gamma}(1.5, \theta)$ prior on σ_b^2 .

When $\theta \rightarrow 0$, the $\text{gamma}(1.5, \theta)$ prior on σ_b^2 has a density function proportional to σ_b , which is also proportional to the $\text{gamma}(2, \theta)$ prior on σ_b . The $\text{gamma}(2, \theta)$ prior on σ_b is recommended as a weakly informative prior for avoiding estimates of σ_b equal to zero in the varying-intercept model (Chung, Rabe-Hesketh, Dorie, et al., 2013) and in random-effects meta-analysis models (Chung, Rabe-Hesketh, & Choi, 2013). Since the $\text{gamma}(2, \theta)$ prior is 0 at $\sigma_b = 0$, the posterior density is also 0 at $\sigma_b = 0$ and thus the posterior mode of σ_b is always strictly positive. In addition, since the gamma density has a positive constant derivative at $\sigma_b = 0$, the $\text{gamma}(2, \theta)$ density increases linearly at zero. It follows that the profile likelihood of σ_b (maximized over all the other parameters) dominates the posterior density of σ_b if the likelihood is strongly curved near $\sigma_b = 0$. That is, the prior does not rule out positive values near zero if they are supported by the likelihood. Chung, Rabe-Hesketh, Dorie, et al. (2013) show that the posterior mode is approximately one standard error away from zero when the ML estimate of σ_b is zero. Finally, the estimator behaves reasonably well in terms of mean squared error of parameter estimates and coverage of confidence intervals for fixed parameters.

In the context of small area estimation, strictly positive group-level variance estimators have been proposed for the Fay and Herriot model (1979), a varying-intercept model for aggregated group-level data and known heterogeneous within-group variances. *Adjustment for density maximization* (Li & Lahiri, 2010; Morris, 2006; Morris & Tang, 2011) applies a penalty term $\pi(\sigma_b^2) = (\sigma_b^2)^{c-1}$ to the likelihood, and this approach turns out to be equivalent to posterior modal estimation with a $\text{gamma}(\alpha, \theta)$ prior on σ_b with $\alpha = 2c + 1$ and $\theta \rightarrow 1$. Therefore, for this specific varying-intercept model, our estimator shares the properties of adjustment for density maximization, such as predictions of the group means being minimax for mean squared-error loss when the within-group variances are equal and $c \leq 1$ (Morris & Tang, 2011).

Varying-Intercept and Varying-Slope Models: d = 2

When $d = 2$, the model includes a varying intercept and a varying slope of one covariate, written as:

$$y_{ij} = \mathbf{x}_{ij}^T \boldsymbol{\beta} + b_{0j} + b_{1j} z_{ij} + \varepsilon_{ij},$$

where $(b_{0j}, b_{1j}) \sim N(\mathbf{0}, \boldsymbol{\Sigma})$ and $\varepsilon_{ij} \sim N(0, \sigma_\varepsilon^2)$.

As shown in Equation 6, with the default choice $\nu = d + 2$, the Wishart density can be written as a product of $\text{gamma}(1.5, \theta)$ densities on the eigenvalues λ_1 and λ_2 . For the bivariate case, we can also express the default prior as a function of the variances (σ_1^2 and σ_2^2) and the correlation (ρ) between the two varying effects b_{0j} and b_{1j} , given by:

$$p(\boldsymbol{\Sigma}) \propto |\boldsymbol{\Sigma}|^{1/2} = \sigma_1 \sigma_2 \sqrt{1 - \rho^2}. \quad (7)$$

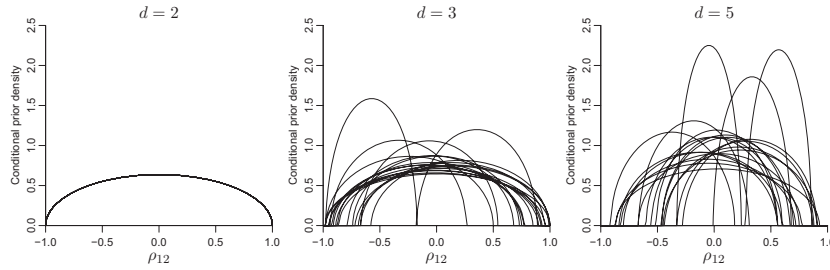


FIGURE 2. Conditional density of ρ_{ij} with Wishart $(d + 2, (1/2\theta)I)$ on Σ , $\theta = 10^{-4}$, where the other parameters are randomly generated from the Wishart distribution for 20 replicates. When $d = 2$, the conditional density is $\text{beta}(1.5, 1.5)$, but for larger d , the curves are more scattered and the supports of the densities become narrower.

This expression implies that $\text{Wishart}(4, (1/2\theta)I_d)$ with $\theta \rightarrow 0$ is equivalent to the joint density of independent $\text{gamma}(1.5, \theta)$ priors on both σ_1^2 and σ_2^2 , and a $\text{beta}(1.5, 1.5)$ prior on $(\rho + 1)/2$.

Since the $\text{beta}(1.5, 1.5)$ prior on $(\rho + 1)/2$ is zero at the boundaries $\rho = \pm 1$, the posterior mode of Σ cannot be attained at any matrices with perfect correlation. In addition, the $\text{beta}(1.5, 1.5)$ density function increases rapidly as ρ approaches 0 from ± 1 and so does not rule out values close to ± 1 . The left panel of Figure 2 shows the $\text{beta}(1.5, 1.5)$ density on $(\rho + 1)/2$. Whereas $\text{gamma}(2, \theta)$ increases linearly at 0, the slopes of $\text{beta}(1.5, 1.5)$ at ± 1 are $\pm \infty$. Therefore, compared to the $\text{gamma}(2, \theta)$ prior for σ_1 and σ_2 , the $\text{beta}(1.5, 1.5)$ for ρ is less informative with lower penalties on the values around the boundaries.

The beta priors have been used to avoid boundary estimates of the probability parameter p of the binomial distribution. When the sample proportion is 0 or 1, the traditional Wald confidence interval for p degenerates to the point estimate. To avoid such boundary estimates, Agresti and Coull (1998) specified a $\text{beta}(2, 2)$ prior on p . The posterior mean of p then is the sample proportion after adding two successes and two failures to the data. Compared with the $\text{beta}(2, 2)$, the $\text{beta}(1.5, 1.5)$ tends to assign less penalty at the boundaries and so is less informative.

Higher Dimensional Case: $d \geq 3$

Similar to the case $d = 2$, the default prior for $d \geq 3$ can be written as a product of σ_r , $r = 1, \dots, d$ and a function of ρ_{rs} , the correlation between the r th and s th varying effects ($0 < r < s, s = 2, \dots, d$). For example with $d = 3$, the $\text{Wishart}(5, (1/2\theta)I_3)$ prior with $\theta \rightarrow 0$ can be written as:

$$p(\Sigma) \propto |\Sigma|^{1/2} \propto \sigma_1 \sigma_2 \sigma_3 \sqrt{1 - \rho_{12}^2 - \rho_{23}^2 - \rho_{13}^2 + 2\rho_{12}\rho_{23}\rho_{13}}. \quad (8)$$

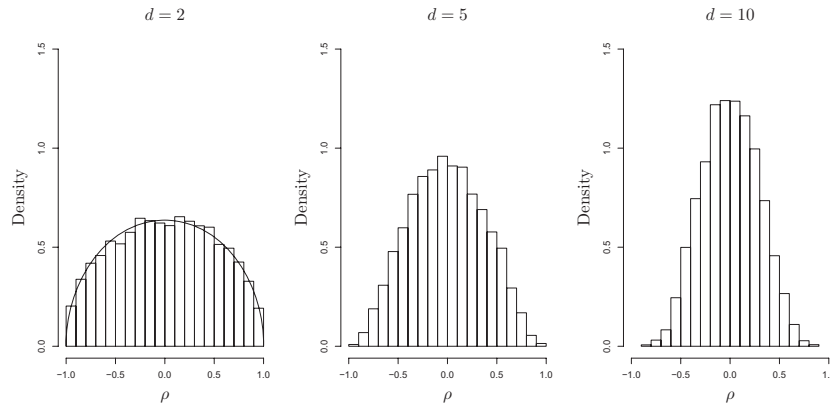


FIGURE 3. Marginal density of ρ_{rs} with Wishart $(d + 2, (1/20)I)$, $\theta = 10^{-4}$. When $d = 2$, the marginal density of ρ is equivalent to $\text{beta}(1.5, 1.5)$ on $(\rho + 1)/2$ (solid curve). As d increases, the marginal density has more mass around 0 due to the positive semidefinite constraint of the covariance matrix.

This is a product of $\text{gamma}(1.5, \theta)$ priors on the variances and a function of the correlations. This function depends on the squares of the correlations, as in the two-dimensional case (Equation 7), but also contains the product of three correlations, which comes from the constraint $|\Sigma| > 0$ that defines the support of Wishart distributions. Because of this constraint, the Wishart prior automatically restricts the posterior mode of Σ to be strictly positive definite.

The graphs in Figure 2 show the conditional densities of ρ_{12} when Σ follows the $\text{Wishart}(d + 2, (1/20)I_d)$, $\theta = 10^{-4}$. The curves are the density of ρ_{12} conditional on the other parameter values (standard deviations and the other correlations) that are randomly generated from $\text{Wishart}(d + 2, (1/20)I_d)$ with 20 replicates. When $d = 2$, the correlation follows $\text{beta}(1.5, 1.5)$ as discussed previously. When $d = 3$, the curves have distinct supports, defined by $1 - (\rho_{12})^2 - (\rho_{23}^0)^2 - (\rho_{13}^0)^2 + 2\rho_{12}\rho_{23}^0\rho_{13}^0 > 0$ where ρ_{13}^0 and ρ_{23}^0 for each replicate are given by randomly generated Σ . The curves for $d = 5$ are more scattered and the supports of the densities tend to be narrower than for $d = 2$ and 3 due to more restrictions required for the higher dimensional Σ to be positive definite.

The marginal prior densities of ρ_{rs} are displayed in Figure 3 for $d = 2, 5$, and 10. With 10,000 replicates, d -dimensional matrices were randomly generated from the $\text{Wishart}(d + 2, (1/20)I)$ with $\theta = 10^{-4}$ and 10,000 $(d - 1)(d - 2)/2$ correlation coefficients were used to construct the histograms. For $d = 2$ (left), the distribution of the correlation coefficient matches the $\text{beta}(1.5, 1.5)$ density, shown as a solid curve. As d increases, the marginal prior density

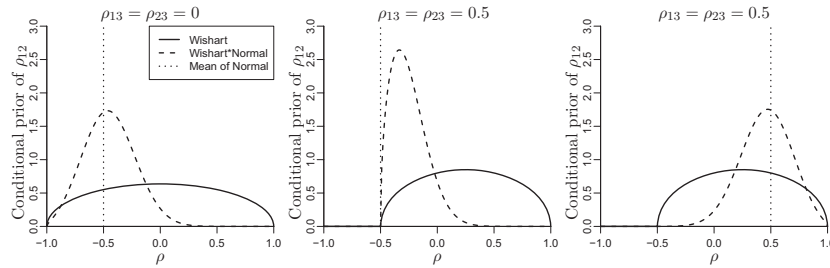


FIGURE 4. Conditional prior density of ρ_{12} with additional $N(-.5, .25^2)$ (left and middle) and $N(.5, .25^2)$ (right) densities multiplying the default Wishart prior. The Wishart prior is on three-dimensional Σ and ρ_{13} and ρ_{23} are fixed as 0 (left) and .5 (middle and right). The additional normal penalty makes the prior density skewed toward the prior value, but still enforces positive definiteness.

of ρ_{rs} becomes more concentrated around zero because of the positive definiteness of Σ .

Incorporating Additional Prior Information

In the previous section, we suggested the $\text{Wishart}(d + 2, (1/2\theta)I)$ with $\theta \rightarrow 0$ as a default prior when no other information is available. If a researcher has additional prior knowledge about any specific standard deviations or correlations, he or she might want to adjust the prior to incorporate such information. In this section, we suggest multiplying the Wishart prior by functions of the parameters on which we have information. Because the Wishart density ensures that Σ is positive definite, we can choose the functions for the other parameters to be intuitive and easy to specify without regard for the parameter space.

If σ^* is a plausible value for σ_r , then the $\text{gamma}(2, 2/\sigma^*)$ density is recommended as a penalty. Recall that the default Wishart prior is proportional to $\text{gamma}(2, \theta)$ priors with $\theta \rightarrow 0$ on each standard deviation, multiplied by a function of the correlations. When the $\text{gamma}(2, 2/\sigma^*)$ density of σ is multiplied by the Wishart, the part including σ_r becomes $\sigma_r^2 \exp(-2\sigma_r/\sigma^*)$. This is proportional to the $\text{gamma}(3, 2/\sigma^*)$ density that has its mode at $\sigma_r = \sigma^*$. The gamma prior with shape parameter greater than two assigns more penalty near zero than for shape parameter equal to two. Therefore, we have a more informative prior with mode at σ^* .

If any specific correlation ρ_{rs} is believed to be close to ρ^* , we can incorporate this prior information by multiplying the default Wishart prior by a $N(\rho^*, \tau^2)$ density. As usual, the scale parameter τ can be chosen depending on the prior uncertainty regarding ρ_{rs} . A possible default choice is $\tau = .25$ because it is the standard deviation of the $\text{beta}(1.5, 1.5)$ distribution. Figure 4 displays the shape of conditional prior densities of ρ_{12} with additional normal

priors in the three-dimensional case. When ρ_{13} and ρ_{23} are fixed at zero (left), the default Wishart($5, (1/20)I_3$) prior (solid curve) is pretty flat. In order to incorporate the prior information, for example, $\rho^* = -.5$, the Wishart is multiplied by the $N(-.5, .25^2)$ density, and then the prior mode moves toward $-.5$ (dashed curve). When ρ_{13} and ρ_{23} are $.5$ (middle and right), the support of the Wishart for ρ_{12} is on $[-0.5, 1]$ because of the constraint of positive definiteness. When our prior value is on the boundary $\rho^* = -.5$ (middle), the Wishart multiplied by $N(-.5, .25^2)$ density is skewed toward $-.5$, but still enforces positive definiteness. When the prior value is inside the support, $\rho^* = -.5$, the resulting density is less skewed (right).

The default prior for ρ in the two-dimensional case is $\text{beta}(1.5, 1.5)$, and so it would seem natural to use the beta family for ρ_{rs} when constructing an additional penalty. However, the parameters of the normal distribution are more intuitive because they represent the prior mean (and mode) and variance. In addition, since the positive definiteness of $\hat{\Sigma}$ is already guaranteed by the Wishart prior, estimates of Σ remain positive definite regardless of the type of additional penalties that multiply the Wishart prior. Furthermore, computation is no problem in any case; including any closed-form prior density adds essentially no cost to the optimization.

Example: A Varying Intercept, Varying Slope Model in Education Research

We illustrate our approach using a study of Heller et al. (2007) on the effects of the Mathematics Pathways and Pitfalls (MPP) teacher professional development program on mathematics learning for students at different levels of English language proficiency. Half of the 36 teachers were randomized to MPP and the other half to the control condition. Teachers randomized to the MPP condition were taught how to use the materials and then substituted MPP for part of their mathematics curriculum during the 2003–2004 school year, while control teachers used their regular mathematics curriculum. All students received an MPP test as a pretest before the lessons and took the same test after the lessons as a posttest.

Posttest scores are regressed on the mean-centered pretest scores, an indicator for treatment group (1 for MPP and 0 for control), English language learner (ELL) status (1 for ELL and 0 for non-ELL), and the Treatment \times ELL Interaction Term. A varying intercept and a varying slope for ELL status are included to allow for the cluster-randomized design. The model can be written as follows:

$$y_{ij} = \mathbf{x}_{ij}^T \boldsymbol{\beta} + b_{0j} + b_{1j} z_{ij} + \varepsilon_{ij},$$

where y_{ij} is the posttest score for the i th student of the j th teacher, \mathbf{x}_{ij} is the covariate vector that includes the mean-centered pretest score, the treatment group indicator, ELL status, and the interaction between ELL status and treatment, and z_{ij} is ELL status. As usual, we assume $(b_{0j}, b_{1j}) \sim N(0, \Sigma)$ and $\varepsilon_{ij} \sim N(0, \sigma_\varepsilon^2)$. After dropping observations with missing values on any of the variables, data were available on 755 students and $J = 36$ teachers, with between 12 and 27

TABLE 1.
Parameter estimates for education example

	ML	REML	BM
Fixed effect			
Intercept	32.39 (2.01)	32.40 (2.07)	32.31 (2.11)
Pretest	0.56 (0.06)	0.56 (0.06)	0.56 (0.06)
Treatment	12.84 (3.15)	12.81 (3.24)	13.01 (3.30)
ELL	-2.46 (2.73)	-2.54 (2.77)	-2.66 (3.17)
ELL \times Treatment	1.00 (4.13)	1.24 (4.19)	1.56 (4.84)
Varying effect (group: teacher)			
Intercept <i>SD</i>	8.31 (1.18)	8.62 (1.25)	8.50 (1.22)
ELL <i>SD</i>	0.71 (2.09)	0.48 (2.18)	3.64 (2.50)
Correlation	-1.00 (2.93)	-1.00 (0.00)	-0.32 (0.22)
Residual <i>SD</i>	226.5	227.3	226.3
Log likelihood	-3,153.7	-3,153.8	-3,154.2

Note. ELL = English language learner; *SD* = standard deviation; ML = maximum likelihood; REML = restricted maximum likelihood; BM = Bayes Modal. The ML and REML estimates imply perfect correlation between the varying intercept and varying slope, whereas BM produces more reasonable estimates. The log likelihood stays almost the same among the three methods. We present results here to more decimal places than would be recommended in practice in order to display the sometimes-small differences between the different estimates.

students per teacher. We fit the models by ML and REML using `lmer` in the `lme4` package and by BM using `blmer` in the `blme` package.

Table 1 presents ML, REML, and BM estimates with the default $\text{Wishart}(4, (1/2\theta)I_d)$ prior with $\theta = 10^{-4}$. Both ML and REML estimates of the correlation between b_{0i} and b_{1j} are -1 . This implies an unrealistic perfect correlation between the teacher-level slopes and intercepts. The BM estimate of ρ is -0.32 and the standard deviation estimate of the varying slope for ELL status increases from 0.71 for ML and 0.48 for REML to 3.64, a change that is within the uncertainty implied by the asymptotic standard error of 2.1 (ML) or 2.2 (REML) for that parameter. The standard deviation of the varying intercept stays similar for ML, REML, and BM.

The fixed coefficient estimates are similar across estimation methods. The coefficient for the interaction term between ELL and treatment changes the most among all the fixed coefficients, but the differences are negligible considering that the standard errors of the interaction term are greater than 4. The standard errors of the fixed coefficient estimates of Treatment, ELL, and Treatment by ELL are larger for BM than for ML or REML, suggesting that ML and REML underestimate the uncertainty.

The log likelihood at the BM estimates differs from the maximum by less than 1. Figure 5 shows the profile likelihood of ρ (profiling out all the other parameters) divided by its maximum. Although the ML is attained at $\rho = -1$, the

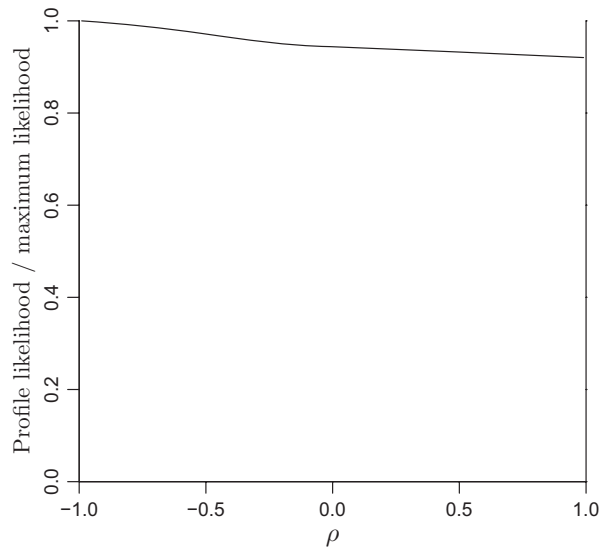


FIGURE 5. Profile likelihood of ρ . The maximum likelihood estimate of ρ is -1 but the likelihood has very little information. Therefore, the Bayes modal estimate of -0.3 is also well supported by the data.

profile likelihood is very flat and so the minimum (at $\rho = 1$) is attained with only an 8% decrement from the maximum. Therefore, all the values of ρ including $\rho = -0.32$ are well supported by the data. As is typical in such settings, there is nothing special about the point estimate on the boundary, and it would be inappropriate for a researcher to use that estimate. Our BM approach gives a default procedure that allows a classical statistician to avoid the inappropriate degenerate estimate. A full Bayes approach using real prior information would do better, but our BM approach takes us a bit in the right direction and has the advantage of being fast and easy to implement.

When a researcher is interested in comparing teacher-specific effects, b_{0j} and b_{1j} can be predicted by their conditional posterior means (or modes), given the estimates of the model parameters and the data (called empirical Bayes prediction or best linear unbiased prediction).

In Figure 6, scatter plots of empirical Bayes predictions of b_{1j} versus b_{0j} are displayed with the proportion of ELL students of each teacher represented by the gray scale, that is, black indicates all the students are ELL and white indicates none are ELL. The sizes of the squares are proportional to the numbers of students for each teacher. For ML (left), due to the estimate $\hat{\rho} = -1$, the slopes b_{1j} are predicted perfectly linearly by the intercepts b_{0j} . In contrast, BM (right)

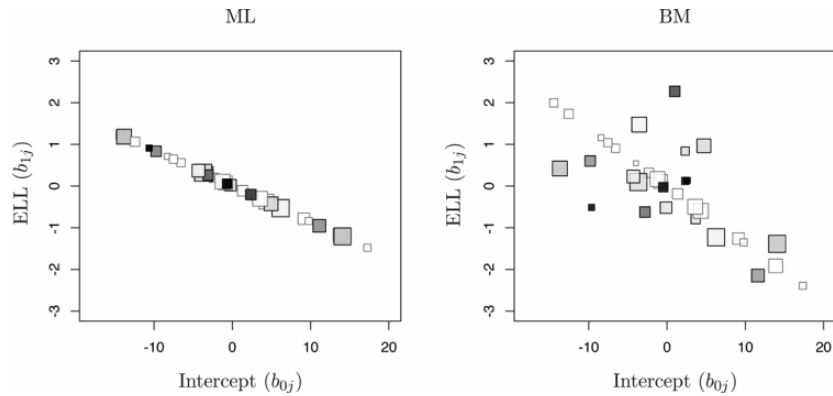


FIGURE 6. Empirical Bayes predictions of varying effects. The size of each square represents n_j for the j th teacher. The ratio of English language learner (ELL) students for each teacher is shown on a gray scale, that is, black indicates all the students are ELL and white indicates none are ELL. For maximum likelihood, b_{1j} are predicted perfectly linearly in b_{0j} . On the right graph, Bayes modal estimation shows more reasonable predictions for the varying slopes and intercepts.

shows more reasonable predictions for the varying slopes and intercepts. In addition, we can observe that 18 (out of 36) white squares with a gray border fall perfectly on a line—these are teachers without any ELL students in their classes. Four black squares (only three visible due to overlap) correspond to teachers with only ELL students. The 18 groups without ELL students and the 4 groups with only ELL students do not provide any information about the slope variance and intercept variance, respectively, and none of the 22 groups provide information about the correlation between the varying slope and the intercept. This lack of information could be one of the reasons we obtain the boundary estimates using ML and REML. As the group size increases (i.e., the square increases) and the proportion of ELL students increases (i.e., the square gets darker), the empirical Bayes predictions tend to be less shrunken toward the line formed by the white squares.

Using the fitted covariance matrix, we can calculate the marginal variances and correlations of the posttest score given ELL status. The variance of the posttest scores for ELL students is $\text{Var}(y_{ij}|z_{ij} = 1) = \sigma_1^2 + \sigma_2^2 + 2\sigma_{12} + \sigma_\epsilon^2$ and, similarly, the variance for non-ELL student is $\text{Var}(y_{ij}|z_{ij} = 0) = \sigma_1^2 + \sigma_\epsilon^2$. The covariance between the posttest scores of two students of the same teacher is $\text{Cov}(y_{ij}, y_{i'j}|z_{ij} = 1, z_{i'j} = 1) = \sigma_1^2 + \sigma_2^2 + 2\sigma_{12}$ if both students are ELL, $\text{Cov}(y_{ij}, y_{i'j}|z_{ij} = 1, z_{i'j} = 0) = \sigma_1^2 + \sigma_{12}$ if one student is ELL, and $\text{Cov}(y_{ij}, y_{i'j}|z_{ij} = 0, z_{i'j} = 0) = \sigma_1^2$ if neither student is ELL.

TABLE 2.
Marginal standard deviations and correlations of posttest scores given ELL status

	ML	BM
SD of ELL student	16.86	17.06
SD of non-ELL student	17.19	17.26
Correlation of (ELL, ELL)	0.20	0.22
Correlation of (ELL, non-ELL)	0.22	0.21
Correlation of (non-ELL, non-ELL)	0.23	0.24

Note. ELL = English language learner; BM = Bayes modal; SD = standard deviation; ML = maximum likelihood. These values do not differ much between ML and BM although the slope standard deviation estimate and correlation estimate increased notably from ML to BM.

Table 2 shows these model-implied marginal standard deviations and correlations with estimates from ML and BM substituted for the parameters. These standard deviation and correlation estimates are remarkably similar, which also explains why the log likelihood evaluated at the BM estimates is not much smaller than that evaluated at the ML estimates.

Simulation

We simulated data from the varying coefficient model as described in the preliminary simulation for Figure 1 but with only one covariate. We explored different values of the correlation ρ (0, 0.225, 0.450, 0.675, and 0.900), setting σ to be a moderate value of 0.5. With 1,000 replicated samples generated with $J = 5$ and $n = 30$, we estimated the bias and root mean squared error (RMSE) for σ_1 , σ_2 , and ρ . For ML and REML, the bias and RMSE of $\hat{\rho}$ are based on the replicates that generate legitimate estimates (i.e., when neither $\hat{\sigma}_1$ nor $\hat{\sigma}_2$ is zero which happened in 1.2% of the replicates for ML and 0.9% of the replicates for REML). For BM estimation, we assigned a Wishart(4, (1/20)I) prior on Σ with $\theta = 10^{-4}$.

Figure 7 shows the proportion of boundary estimates of ρ (where $1 - |\hat{\rho}| < 10^{-5}$). When ρ is 0, 21% of the ML estimates and 17% of the REML estimates have perfect correlations. As ρ increases, the proportion of $\hat{\rho}$ on the boundary also increases and reaches 60% for ML and 51% for REML. The BM method does not produce any boundary estimates of ρ for any of the simulation conditions.

In spite of the absence of boundary estimates, the log likelihood is not reduced substantially by using BM estimation. Investigating the difference in deviances ($= 2[\log L(\hat{\Sigma}_{ML}) - \log L(\hat{\Sigma}_{BM})]$) for all the replicates, the BM method never reduces the log likelihood by more than 2.2 from the maximum.

Figure 8 summarizes the estimated bias and RMSE of $\hat{\rho}$, $\hat{\sigma}_1$, and $\hat{\sigma}_2$. When $\rho = 0$, the estimated bias of $\hat{\rho}$ is almost zero for all three methods.

ML, REML, and BM all have some bias in estimating ρ , with BM having the most bias (i.e., the most shrinkage toward 0), as would be expected given the

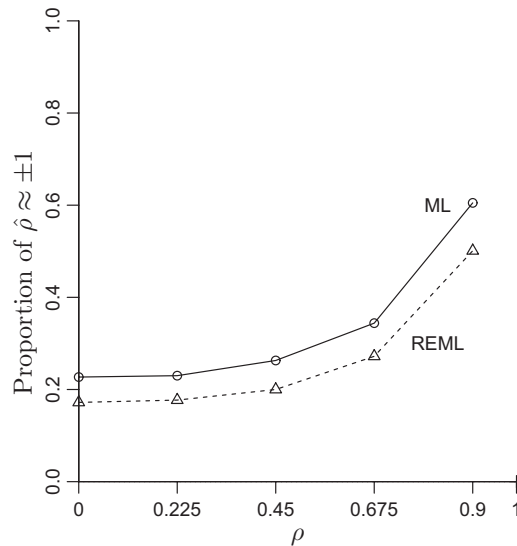


FIGURE 7. Proportion of maximum likelihood (ML) and restricted maximum likelihood (REML) estimates of ρ that are on the boundary. When $\rho = 0$, 21% of the ML estimates and 17% of the REML estimates are ± 1 . As ρ increases, the proportion of estimates on the boundary, $\hat{\rho}$ equal to ± 1 , also increases and reaches 60% for ML and 51% for REML when $\rho = .9$.

regularization from the Wishart prior that squeezes $\hat{\rho}$ toward zero as seen in the shape of the prior density for ρ with $d = 2$ in Figure 3. However, BM gives the smallest estimated RMSE of $\hat{\rho}$. The estimated bias of $\hat{\sigma}_1$ and $\hat{\sigma}_2$ is similar across the different values of ρ for all the estimation methods. The BM estimates of the standard deviations are less biased and have smaller estimated RMSE than ML and REML.

The coverage of 95% confidence intervals for β_0 and β_1 does not change much with ρ . The average coverage of the BM confidence intervals is .940 for β_0 and .943 for β_1 . The coverage for REML is about the same as that for BM, whereas ML shows slightly lower coverage with averages of .935 for β_0 and .937 for β_1 .

Conclusion

For the hierarchical regression model, particularly with several varying coefficients, degenerate covariance matrix estimates do not have a practical interpretation. Unfortunately, such boundary estimates commonly arise in ML estimation because there is often little information on these parameters when there is only a moderate number of groups. In addition, when $\hat{\Sigma}$ is singular, underestimated standard errors of the fixed coefficients make the researcher overconfident about the

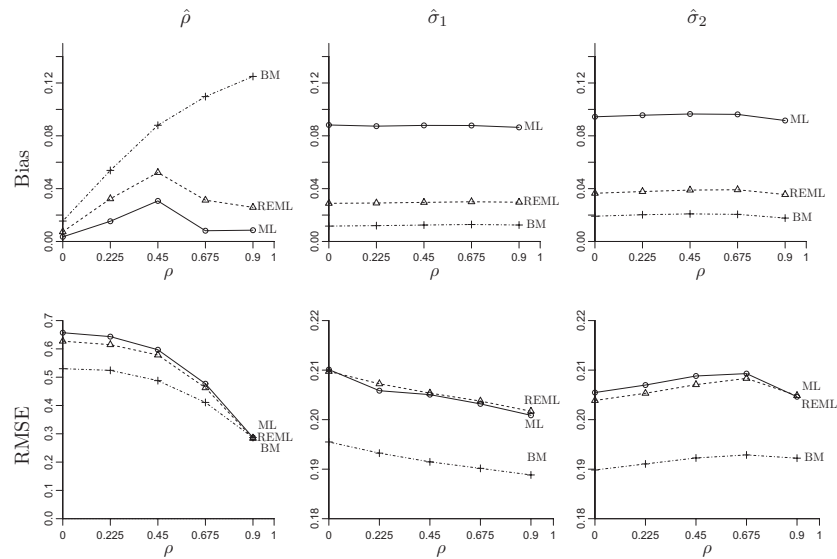


FIGURE 8. Bias and root mean squared error of $\hat{\sigma}_1$, $\hat{\sigma}_2$, and $\hat{\rho}$ with $J = 5$ and $n = 30$ of the varying-coefficient model with $\sigma_1 = \sigma_2 = .5$ and ρ in the grid. In our simulation, with ρ set to various positive values, the bias values are all negative, so we display absolute values to make the graphs easier to read given the convention that high values of bias are bad. Bayes modal (BM) has higher bias for ρ (i.e., shrinking the estimate toward 0) compared to maximum likelihood (ML) and restricted maximum likelihood (REML), but the RMSE is smaller for BM. For both σ_1 and σ_2 , BM has smaller bias and RMSE than ML and REML.

effect of the covariates. When a boundary estimate is attained but no prior information is available for Σ , the BM estimator using the default Wishart prior is recommended because it ensures strictly positive definite $\hat{\Sigma}$ and is weakly informative at the same time. The modified gllamm from www.gllamm.org for Stata and blme package for R allow straightforward application of our method for practitioners.

In varying-slope models, changing the location and scale of the covariates that have varying slopes implies that Σ must change to produce an equivalent model. For example, for longitudinal data, we might want to transform the time variable to have a value 0 at the initial time point. In this case, subtracting a constant from the covariate changes the variance of the varying intercepts and the correlation between intercepts and slopes. Although ML and REML will yield equivalent models after linearly transforming the covariate, this is no longer true for BM estimation, which pulls the correlation toward 0. When using Bayesian regularization in this setting, it therefore becomes more

important to choose meaningful centering points for the covariates with varying coefficients.

Authors' Note

The data from Math Pathways and Pitfalls Lessons on Students Mathematics Achievement study are Copyright ©2011 by WestEd. All rights reserved. Any opinion, findings, and conclusions or recommendations expressed in this material are those of the author and may not reflect the views, findings, or opinions of the National Science Foundation or WestEd.

Declaration of Conflicting Interests

The authors declared no potential conflicts of interest with respect to the research, authorship, and/or publication of this article.

Funding

The authors disclosed receipt of the following financial support for the research, authorship, and/or publication of this article: The research reported here was supported by the Institute of Education Sciences (R305D100017), the National Science Foundation (SES-1323977), and the Army Grant (W911NF-14-1-0020). This data set is based upon work supported by the National Science Foundation under Grant No. 9911374 along with materials developed by WestEd.

References

- Agresti, A., & Coull, B. A. (1998). Approximate is better than “exact” for interval estimation of binomial proportions. *The American Statistician*, *52*, 119–126.
- Amemiya, Y. (1985). What should be done when an estimated between-group covariance matrix is not nonnegative definite? *The American Statistician*, *39*, 112–117.
- Barnard, J., McCulloch, R., & Meng, X. (2000). Modeling covariance matrices in terms of standard deviations and correlations, with application to shrinkage. *Statistica Sinica*, *10*, 1281–1312.
- Bates, D., & Maechler, M. (2010). *lme4: Linear mixed-effects models using s4 classes* [Computer software manual]. Retrieved from <http://CRAN.R-project.org/package=lme4> (R package version 0.999375-37)
- Chung, Y., Rabe-Hesketh, S., & Choi, I.-H. (2013). Avoiding zero between-study variance estimates in random-effects meta-analysis. *Statistics in Medicine*, *32*, 4071–4089.
- Chung, Y., Rabe-Hesketh, S., Dorie, V., Gelman, A., & Liu, J. (2013). A nondegenerate penalized likelihood estimator for variance parameters in multilevel models. *Psychometrika*, *78*, 685–709.
- Ciuperca, G., Ridolfi, A., & Idier, J. (2003). Penalized maximum likelihood estimator for normal mixtures. *Scandinavian Journal of Statistics*, *30*, 45–59.
- Daniels, M. J., & Kass, R. E. (1999). Nonconjugate Bayesian estimation of covariance matrices and its use in hierarchical models. *Journal of the American Statistical Association*, *94*, 1254–1263.

- Fay, R. E., & Herriot, R. A. (1979). Estimates of income for small places: An application of James-Stein procedures to census data. *Journal of the American Statistical Association*, *74*, 269–277.
- Galindo-Garre, F., & Vermunt, J. (2006). Avoiding boundary estimates in latent class analysis by Bayesian posterior mode estimation. *Behaviormetrika*, *33*, 43–59.
- Galindo-Garre, F., Vermunt, J., & Bergsma, W. (2004). Bayesian posterior mode estimation of logit parameters with small samples. *Sociological Methods & Research*, *33*, 88–117.
- Gelman, A., Jakulin, A., Pittau, M. G., & Su, Y. S. (2008). A weakly informative default prior distribution for logistic and other regression models. *The Annals of Applied Statistics*, *2*, 1360–1383.
- Harville, D. A. (1974). Bayesian inference for variance components using only error contrasts. *Biometrika*, *61*, 383–385.
- Kolenikov, S., & Bollen, K. A. (2012). Testing negative error variances: Is a Heywood case a symptom of misspecification? *Sociological Methods & Research*, *41*, 124–167.
- Lehmann, E. L., & Casella, G. (1998). *Theory of point estimation*. New York, NY: Springer.
- Li, H., & Lahiri, P. (2010). An adjusted maximum likelihood method for solving small area estimation problems. *Journal of Multivariate Analysis*, *101*, 882–892.
- Maris, E. (1999). Estimating multiple classification latent class models. *Psychometrika*, *64*, 187–212.
- Martin, J. K., & McDonald, R. P. (1975). Bayesian estimation in unrestricted factor analysis: A treatment for Heywood cases. *Psychometrika*, *40*, 505–517.
- Mislevy, R. J. (1986). Bayes modal estimation in item response models. *Psychometrika*, *51*, 177–195.
- Morris, C. (2006). Mixed model prediction and small area estimation (with discussions). *Test*, *15*, 72–76.
- Morris, C., & Tang, R. (2011). Estimating random effects via adjustment for density maximization. *Statistical Science*, *26*, 271–287.
- O'Malley, A. J., & Zaslavsky, A. M. (2005). *Cluster-level covariance analysis for survey data with structured nonresponse* (Tech. Rep.). Boston, MA: Department of Health Care Policy, Harvard Medical School.
- Patterson, H. D., & Thompson, R. (1971). Recovery of inter-block information when block sizes are unequal. *Biometrika*, *58*, 545–554.
- Rabe-Hesketh, S., Skrondal, A., & Pickles, A. (2005). Maximum likelihood estimation of limited and discrete dependent variable models with nested random effects. *Journal of Econometrics*, *128*, 301–323.
- Srivastava, M., & Kubokawa, T. (1999). Improved nonnegative estimation of multivariate components of variance. *Annals of Statistics*, *27*, 2008–2032.
- Swaminathan, H., & Gifford, J. A. (1985). Bayesian estimation in the two-parameter logistic model. *Psychometrika*, *50*, 349–364.
- Tsutakawa, R. K., & Lin, H. Y. (1986). Bayesian estimation of item response curves. *Psychometrika*, *51*, 251–267.
- Verbeke, G., & Molenberghs, G. (2000). *Linear mixed models for longitudinal data*. New York, NY: Springer.
- Vermunt, J., & Magidson, J. (2005). *Technical guide for Latent Gold 4.0: Basic and advanced* (Tech. Rep.). Belmont, MA: Statistical Innovations.

Chung et al.

Warton, D. I. (2008). Penalized normal likelihood and ridge regularization of correlation and covariance matrices. *Journal of the American Statistical Association*, 103, 340–349.

Authors

YEOJIN CHUNG is an assistant professor at the School of Business Administration, Kookmin University, Seoul 136-702, South Korea; e-mail: ychung@kookmin.ac.kr. Her research interest is in multilevel modeling, nonparametric density estimation, and model-based clustering.

ANDREW GELMAN is a professor at the Department of Statistics, Columbia University, New York, NY 10027, USA; e-mail: gelman@stat.columbia.edu. He specializes in Bayesian data analysis.

SOPHIA RABE-HESKETH is a professor at the Graduate School of Education and Graduate Group in Biostatistics, University of California, Berkeley, CA 94720, USA; e-mail: sophiarh@berkeley.edu. Her primary research interests include multilevel and latent variable modeling.

JINGCHEN LIU is an associate professor at the Department of Statistics, Columbia University, New York, NY 10027, USA; e-mail: jeliu@stat.columbia.edu. His research interests are in importance sampling, extremes of Gaussian random fields, and Bayesian modelling for missing data problems.

VINCENT DORIE is a postdoctoral research fellow at the Center for the Promotion of Research Involving Innovative Statistical Methodology at New York University, New York, NY, USA; e-mail: vjd4@nyu.edu. He is currently researching Bayesian nonparametrics and causal inference.

Manuscript received April 2, 2014

Revision received August 6, 2014

Accepted November 14, 2014