

Paper 253-28

An Introduction to the Analysis of Mixed Models

Dallas E. Johnson, Kansas State University, Manhattan, KS

ABSTRACT

This paper introduces the General Linear Mixed Model (GLMM) and compares various alternatives for estimating estimable functions of the model parameters and provides some of the advantages and disadvantages of each of the alternative methods of estimation. The paper also considers estimating the standard errors of the various estimates of estimable functions and how the estimates and their estimated standard errors can be used for statistical inference. Special cases of a GLMM include all types of split-plot experiments, repeated measures experiments, and various combinations of these. Examples will be given to illustrate how these kinds of experiments fit into a GLMM framework.

INTRODUCTION

The General Linear Mixed Model (GLMM) is defined by

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\epsilon}$$

where \mathbf{y} is an $n \times 1$ observable data vector, $\boldsymbol{\beta}$ is a $p \times 1$ vector of unknown parameters, \mathbf{u} is a $q \times 1$ vector of unobservable random variables, \mathbf{X} and \mathbf{Z} are design matrices corresponding to the fixed and random effects, respectively, and $\boldsymbol{\epsilon}$ is a vector of random errors. The basic assumptions are that

$$\begin{aligned} \boldsymbol{\epsilon} &\sim N(\mathbf{0}, \mathbf{R}_{n \times n}), \\ \mathbf{u} &\sim N(\mathbf{0}, \mathbf{G}_{q \times q}), \end{aligned}$$

and

$$\text{COV}(\boldsymbol{\epsilon}, \mathbf{u}) = \mathbf{0}_{n \times q}.$$

Note that

$$E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$$

and

$$\text{Cov}(\mathbf{y}) = \mathbf{Z}\mathbf{G}\mathbf{Z}' + \mathbf{R} = \mathbf{V} \text{ (say).}$$

Thus

$$\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \mathbf{V}).$$

ESTIMATORS OF $\ell'\boldsymbol{\beta}$

Four different estimators of $\ell'\boldsymbol{\beta}$ are introduced where $\ell'\boldsymbol{\beta}$ is an estimable function of the parameters in $\boldsymbol{\beta}$. That is, ℓ belongs to the column space of \mathbf{X}' .

Generalized Least Squares Estimators

If \mathbf{G} and \mathbf{R} are both known, then \mathbf{V} is also known, and in this special case, the uniformly minimum variance unbiased

(UMVU) estimate of any estimable function $\ell'\boldsymbol{\beta}$ is $\ell'\hat{\boldsymbol{\beta}}_{\mathbf{V}}$ where

$$\hat{\boldsymbol{\beta}}_{\mathbf{V}} = (\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1} \mathbf{X}'\mathbf{V}^{-1}\mathbf{y}.$$

Furthermore, under the basic assumptions given above,

$$\ell'\hat{\boldsymbol{\beta}}_{\mathbf{V}} \sim N(\ell'\boldsymbol{\beta}, \ell'(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\ell).$$

This estimator of $\ell'\boldsymbol{\beta}$ is called the *Generalized Least Squares Estimator*.

Estimated Generalized Least Squares Estimators

In practice, \mathbf{G} and \mathbf{R} are rarely known! However, suppose $\hat{\mathbf{G}}$ and $\hat{\mathbf{R}}$ can be estimated in some way by

$$\hat{\mathbf{G}} \text{ and } \hat{\mathbf{R}}, \text{ respectively. Then } \ell'\boldsymbol{\beta} \text{ could be estimated}$$

by $\ell'\hat{\boldsymbol{\beta}}_{\hat{\mathbf{V}}}$ where

$$\hat{\boldsymbol{\beta}}_{\hat{\mathbf{V}}} = (\mathbf{X}'\hat{\mathbf{V}}^{-1}\mathbf{X})^{-1} \mathbf{X}'\hat{\mathbf{V}}^{-1}\mathbf{y} \text{ and } \hat{\mathbf{V}} = \mathbf{Z}\hat{\mathbf{G}}\mathbf{Z}' + \hat{\mathbf{R}}.$$

The estimator $\ell'\hat{\boldsymbol{\beta}}_{\hat{\mathbf{V}}}$ is called the *Estimated Generalized Least Squares Estimator* (EGLS) of $\ell'\boldsymbol{\beta}$.

In most balanced experiments $\ell'\hat{\boldsymbol{\beta}}_{\hat{\mathbf{V}}}$ is an unbiased estimator of $\ell'\boldsymbol{\beta}$, but in unbalanced experiments $\ell'\hat{\boldsymbol{\beta}}_{\hat{\mathbf{V}}}$ cannot be shown to be an unbiased estimator of $\ell'\boldsymbol{\beta}$.

Ordinary Least Squares Estimators

Another estimator of $\ell'\boldsymbol{\beta}$, called the *Ordinary Least Squares Estimator* of $\ell'\boldsymbol{\beta}$, is $\ell'\hat{\boldsymbol{\beta}}_{\mathbf{O}}$ where $\hat{\boldsymbol{\beta}}_{\mathbf{O}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}$.

This estimator is always an unbiased estimator of $\ell'\boldsymbol{\beta}$. In many balanced experiments

$$\ell'\hat{\boldsymbol{\beta}}_{\hat{\mathbf{V}}} = \ell'\hat{\boldsymbol{\beta}}_{\mathbf{O}}.$$

When this is true, then $\ell'\hat{\boldsymbol{\beta}}_{\hat{\mathbf{V}}}$ is an unbiased estimator of $\ell'\boldsymbol{\beta}$. It can also be shown that,

$$\ell'\hat{\boldsymbol{\beta}}_{\mathbf{O}} \sim N(\ell'\boldsymbol{\beta}, \ell'(\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{V}\mathbf{X})(\mathbf{X}'\mathbf{X})^{-1}\ell).$$

Even though $\ell'\hat{\boldsymbol{\beta}}_{\hat{\mathbf{V}}}$ is not unbiased for $\ell'\boldsymbol{\beta}$ in unbalanced cases, one might expect that $\ell'\hat{\boldsymbol{\beta}}_{\hat{\mathbf{V}}}$ should be a reasonable estimator of $\ell'\boldsymbol{\beta}$.

In most unbalanced cases one cannot determine the variance of $\ell'\hat{\boldsymbol{\beta}}_{\hat{\mathbf{V}}}$. However, note that

$$\text{Var}(\ell' \hat{\beta}_V) = \ell'(X'V^{-1}X)\ell,$$

and so it would seem reasonable to approximate the variance of

$$\ell' \hat{\beta}_V \text{ with } \ell'(X'\hat{V}^{-1}X)\ell \text{ so that so that approximate}$$

estimated standard error of $\ell' \hat{\beta}_V$ is

$$\text{s.e.}(\ell' \hat{\beta}_V) = \sqrt{\ell'(X'\hat{V}^{-1}X)\ell}.$$

GLM Estimators

Another possible estimator of $\ell'\beta$ is

$$\ell' \hat{\beta}_{GLM} + a'u$$

where

$$\begin{bmatrix} \hat{\beta}_{GLM} \\ u \end{bmatrix} = \begin{bmatrix} X'X & X'Z \\ Z'X & Z'Z \end{bmatrix}^{-1} \begin{bmatrix} X'y \\ Z'y \end{bmatrix}$$

and where a is chosen so $\begin{bmatrix} \ell \\ a \end{bmatrix}$ belongs to the column

space of $\begin{bmatrix} X'X & X'Z \\ Z'X & Z'Z \end{bmatrix}$. This estimator is called the GLM

estimator as it is the estimator that one gets from SAS-GLM when the random effects are treated as fixed effects. The GLM estimator can be shown to be an unbiased estimator of $\ell'\beta$. One

would like to choose a so that $\text{VAR}(\ell' \hat{\beta}_{GLM} + a'u)$ is

minimized. The variance of $\ell' \hat{\beta}_{GLM} + a'u$ is equal to

$$\begin{bmatrix} \ell \\ a \end{bmatrix}' \begin{bmatrix} X'X & X'Z \\ Z'X & Z'Z \end{bmatrix}^{-1} \begin{bmatrix} X'y \\ Z'y \end{bmatrix} + \begin{bmatrix} \ell \\ a \end{bmatrix}' \begin{bmatrix} X'X & X'Z \\ Z'X & Z'Z \end{bmatrix}^{-1} \begin{bmatrix} X'y \\ Z'y \end{bmatrix}.$$

Choosing an appropriate a is easy in balanced cases, and in these cases SAS-GLM does the work for you. Choosing an appropriate a in unbalanced cases is much harder, and SAS-GLM rarely chooses an appropriate a for you. Usually, SAS-GLM gives you a message indicating that the LSMEANS are not estimable or if you are using ESTIMATE and/or CONTRAST options, you often get a message that the contrast is not estimable. In all cases, these happen because the criteria that SAS-GLM uses to select a does not provide an a such that

$$\begin{bmatrix} \ell \\ a \end{bmatrix} \text{ belongs to the column space of } \begin{bmatrix} X'X & X'Z \\ Z'X & Z'Z \end{bmatrix}. \text{ The}$$

GLM estimator, $\ell' \hat{\beta}_{GLM} + a'u$ of $\ell'\beta$ should also be a

reasonable estimator of $\ell'\beta$. Like the OLS estimator of $\ell'\beta$, it is unbiased.

To estimate the variance of the GLM estimator, one would need to estimate R and G , then $\text{Var}(\ell' \hat{\beta}_{GLM} + a'u)$ is estimated

$$\text{by } \begin{bmatrix} \ell \\ a \end{bmatrix}' \begin{bmatrix} X'X & X'Z \\ Z'X & Z'Z \end{bmatrix}^{-1} \begin{bmatrix} X'y \\ Z'y \end{bmatrix} + \begin{bmatrix} \ell \\ a \end{bmatrix}' \begin{bmatrix} X'X & X'Z \\ Z'X & Z'Z \end{bmatrix}^{-1} \begin{bmatrix} X'y \\ Z'y \end{bmatrix}.$$

In all cases, with any of these estimators, in order to make inferences about $\ell'\beta$, it is usually assumed that

$$\frac{\ell' \hat{\beta} - \ell' \beta}{\text{s.e.}(\ell' \hat{\beta})} \sim t(v)$$

for some appropriate value of v which usually has to be estimated in some way.

In order to compute the EGLS estimate of $\ell'\beta$, one must first be able to estimate R and G . While general estimates of R and G do not exist, there are many special cases in which one can estimate R and G .

Split-Plot Type Cases

In many cases the GLMM,

$$y_{n \times 1} = X_{n \times p} \beta_{p \times 1} + Z_{n \times q} u_{q \times 1} + \epsilon_{n \times 1}$$

can be written as

$$y = X\beta + Z_1u_1 + Z_2u_2 + \dots + Z_ru_r + \epsilon \quad (1)$$

where it can be assumed that $u_i \sim$ independent $N(0, \sigma_i^2 I)$ for $i=1,2,\dots,r$, $\epsilon \sim N(0, \sigma_\epsilon^2 I)$, and that $\text{COV}(u_i, \epsilon) = 0$ for all i . In this case, $V = \sigma_1^2 Z_1Z_1' + \sigma_2^2 Z_2Z_2' + \dots + \sigma_r^2 Z_rZ_r' + \sigma_\epsilon^2 I$. Note that each u_i is a $q_i \times 1$ random vector with corresponding design matrix Z_i which is an $n \times q_i$ matrix. Here estimating V is accomplished by estimating $\sigma_1^2, \sigma_2^2, \dots, \sigma_r^2$, and σ_ϵ^2 , called variance components.

All split-plot type experiments fit into this framework. Repeated measures experiments also fit into this general framework when the repeated measures satisfy compound symmetry assumptions, and all repeated measures experiments can be written in the framework of model (1) simply by placing fewer restrictions on $R=\text{COV}(\epsilon)$.

Example 1: A Simple Split-Plot Model Suppose

$$y_{ijk} = \mu + T_i + \beta_j + Y_{ij} + \delta_{k(i)} + \epsilon_{ijk}$$

for $i=1,2,\dots,t$; $j=1,2,\dots,b$; and $k=1,2,\dots,n_i$; where $\delta_{k(i)} \sim$ i.i.d. $N(0, \sigma_\delta^2)$, $\epsilon_{ijk} \sim$ i.i.d. $N(0, \sigma_\epsilon^2)$, and where all the $\delta_{k(i)}$'s and ϵ_{ijk} 's are independent. Then $y = X\beta + Z_1u_1 + \epsilon$ where $y' = [y_{111} \ y_{112} \ \dots \ y_{11n_1} \ y_{121} \ y_{122} \ \dots \ y_{12n_2} \ \dots \ y_{tbn} \ y_{tbn} \ \dots \ y_{tbn}]$, $\beta' = [\mu \ T_1 \ T_2 \ \dots \ T_t \ \beta_1 \ \beta_2 \ \dots \ \beta_b \ Y_{11} \ Y_{12} \ \dots \ Y_{tb}]$, $u_1' = [\delta_{1(1)} \ \delta_{2(1)} \ \dots \ \delta_{n_1(1)} \ \dots \ \delta_{1(t)} \ \delta_{2(t)} \ \dots \ \delta_{n_t(t)}]$, $\epsilon' = [\epsilon_{111} \ \epsilon_{112} \ \dots \ \epsilon_{11n_1} \ \epsilon_{121} \ \epsilon_{122} \ \dots \ \epsilon_{12n_2} \ \dots \ \epsilon_{tbn} \ \epsilon_{tbn} \ \dots \ \epsilon_{tbn}]$, and X and Z_1 are design matrices determined by the pattern of observed data. The ideal conditions for the errors in the split-plot model can be restated as $u_1 \sim N(0, \sigma_\delta^2 I_N)$ where $N=n_1+n_2+\dots+n_t$, $\epsilon \sim N(0, \sigma_\epsilon^2 I_{nb})$, and that u_1 and ϵ are independent. For the simple split plot model, one must estimate σ_δ^2 and σ_ϵ^2 in order to estimate V .

Example 2: A Simple Repeated Measures Model. Suppose

$$y_{ijk} = \mu + T_i + \beta_j + Y_{ij} + \epsilon_{ijk}$$

for $i=1,2,\dots,t$; $j=1,2,\dots,b$; and $k=1,2,\dots,n_i$. Let $\epsilon_{ik}' = [\epsilon_{i1k} \ \epsilon_{i2k} \ \dots \ \epsilon_{ibk}]$. Assume that the ϵ_{ik} 's are i.i.d. $N(0, \Sigma)$ for $i=1,2,\dots,t$; $k=1,2,\dots,n_i$. Thus ϵ_{ik} represents the $b \times 1$ vector of errors corresponding to the b repeated measures on the k th subject in the i th treatment group.

For the simple repeated measures experiment $Z_1=0$ which implies that $G=0$ and $V=R$. Here

$$R = \text{COV}(\epsilon) = \text{COV} \left(\begin{bmatrix} \epsilon_{11} \\ \epsilon_{12} \\ \vdots \\ \epsilon_{tn} \end{bmatrix} \right) = \begin{bmatrix} \Sigma & 0 & \dots & 0 \\ 0 & \Sigma & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \Sigma \end{bmatrix}, \text{ for}$$

$$Y' = [y_{111} \ y_{121} \ \dots \ y_{1b1} \ y_{112} \ y_{122} \ \dots \ y_{1b2} \ \dots \ y_{11n} \ y_{12n} \ \dots \ y_{1bn}]$$

ESTIMATING G and R

Methods for estimating G and R are considered in this section.

Estimates of G and R in Split-plot Type Experiments.

For the general case, one has

$$y = X\beta + Z_1u_1 + Z_2u_2 + \dots + Z_ru_r + \epsilon$$

where it can be assumed that

$$u_i \sim \text{independent } N(0, \sigma_i^2 I) \text{ for } i=1,2,\dots,r,$$

$$\epsilon \sim N(0, \sigma_e^2 I),$$

and that

$$\text{COV}(u_i, \epsilon) = 0 \text{ for all } i.$$

In this case,

$$V = \sigma_1^2 Z_1 Z_1' + \sigma_2^2 Z_2 Z_2' + \dots + \sigma_r^2 Z_r Z_r' + \sigma_e^2 I$$

and estimating V is accomplished by estimating $\sigma_1^2, \sigma_2^2, \dots, \sigma_r^2$, and σ_e^2 .

Methods of Estimating the Variance Components

1. ANOVA (Method of Moment) Estimates

- (a) Fit the GLMM by assuming that the random effects in the model are fixed effects.
- (b) Next obtain the corresponding ANOVA table and compute the expected mean squares of the observed mean squares in the ANOVA table under the true assumptions about the u_i 's and ϵ .
- (c) Equate the observed mean squares to their expected mean squares and solve the resulting system of equations for each of the variance components.

- (d) Use the resulting solutions as the estimates of the variance components.

Advantages:

- (a) The estimators of the variance components are unbiased.
- (b) One can often approximate the degrees of freedom corresponding to the estimated standard errors of estimators of estimable functions of the fixed effects by using Satterthwaite's Method.
- (c) SAS - GLM can produce the necessary information to perform these analyses.

Disadvantages:

- (a) There is no unique way in which to form an ANOVA table when the data are not balanced.
- (b) The procedure can produce negative estimates of the variance components which do not make sense. For the simple repeated measures model, one must estimate Σ in order to estimate V .
- (c) If some of the expected mean squares of the random effects in the ANOVA table depend on fixed effects, the method cannot be applied. This problem can be avoided by placing all of the fixed effects in the model first followed by the random effects.

2. Maximum Likelihood Method

Find the values of $\beta, \sigma_1^2, \sigma_2^2, \dots, \sigma_r^2$, and σ_e^2 that maximize the likelihood function over the parameter space.

Advantages

- (a) Avoids negative estimates of the variance components.

Disadvantages

- (a) Numerically intensive
- (b) Resulting estimators are not unbiased.
- (c) Solving the likelihood equations requires an iterative process which may or may not converge. Even when it converges, it may converge at a local maxima rather than a global maximum.
- (d) Tends to underestimate the variance components.
- (e) Distributional properties are not known except asymptotically.

3. Restricted Maximum Likelihood (REML)

REML estimators of the variance components are found by maximizing that part of the likelihood function that is invariant to fixed effects in the model.

The GLMM is given by

$$y = X\beta + Zu + \epsilon$$

Let L be a full row rank matrix that satisfies $LX=0$ and such that $\text{rank}(L) = n - \text{rank}(X)$ where n is the dimension of y . Let $y^* = Ly$. Then

$$y^* \sim N(0, \sigma_1^2 LZ_1Z_1'L' + \sigma_2^2 LZ_2Z_2'L' + \dots + \sigma_r^2 LZ_rZ_r'L' + \sigma_e^2 I).$$

The likelihood function formed from y^* depends only on the variance components; and the REML estimates of the variance components are those values of $\sigma_1^2, \sigma_2^2, \dots, \sigma_r^2$, and σ_e^2 that maximize the restricted likelihood function formed from y^* .

Advantages:

- (a) Less numerically intensive than the Maximum Likelihood Method.
- (b) The REML estimates and the ANOVA estimates agree when the data are balanced and all MM estimates of the variance components are non-negative.
- (c) REML estimates tend to be less biased than the Maximum Likelihood Estimates.

Disadvantage:

- (a) The distributional properties of these estimators are not known, except asymptotically.

Repeated Measures (Continued)

In the simple repeated measures experiment,

$$y = X\beta + \epsilon \text{ where}$$

$$\epsilon \sim N(0, \begin{bmatrix} \Sigma & 0 & \dots & 0 \\ 0 & \Sigma & \dots & 0 \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & \dots & \Sigma \end{bmatrix}).$$

If there are p repeated measures on each subject, then each Σ is a $p \times p$ matrix, provided that there are no missing data for any subject.

Here, in order to estimate V , one must estimate Σ . The estimate of Σ depends on whether one assumes any structure on Σ .

Reasonable Structures on Σ .

1. Unstructured Σ . Here it is assumed that

$$\Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2p} \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \sigma_{p1} & \sigma_{p2} & \dots & \sigma_{pp} \end{bmatrix}.$$

2. Compound Symmetry. Here it is assumed that

$$\Sigma = \sigma^2 \begin{bmatrix} 1 & \rho & \dots & \rho \\ \rho & 1 & \dots & \rho \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \rho & \rho & \dots & 1 \end{bmatrix}.$$

3. Huyhn-Feldt Conditions. Here it is assumed that $\Sigma = \lambda I + \eta j' + j\eta'$ for some λ and some η . This implies

$$\Sigma = \begin{bmatrix} \lambda + 2\eta_1 & \eta_1 + \eta_2 & \dots & \eta_1 + \eta_p \\ \eta_2 + \eta_1 & \lambda + 2\eta_2 & \dots & \eta_2 + \eta_p \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \eta_p + \eta_1 & \eta_p + \eta_2 & \dots & \lambda + 2\eta_p \end{bmatrix}.$$

4. AR(1). Here it is assumed that

$$\Sigma = \sigma^2 \begin{bmatrix} 1 & \rho & \rho^2 & \dots & \rho^{p-1} \\ \rho & 1 & \rho & \dots & \rho^{p-2} \\ \rho^2 & \rho & 1 & \dots & \rho^{p-3} \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \rho^{p-1} & \rho^{p-2} & \rho^{p-3} & \dots & 1 \end{bmatrix}.$$

In each of the above structures for Σ as well as other possible structures, estimating V requires one to estimate the parameters in Σ . For the unstructured case and the compound symmetry case, Σ can be estimated either by method of moments (MoM), maximum likelihood (ML), or restricted maximum likelihood (REML). If there are no missing data on any subject, the estimates will agree. For all other cases, the parameters in Σ must be estimated by ML or REML methods.

Finally, in each of these cases, once estimates of the variance parameters are found, V can be estimated, and then $l'\beta$ is estimated by $l'\hat{\beta}_V$ where $\hat{\beta}_V = (X'V^{-1}X)^{-1} X'V^{-1}y$. The

standard error of $l'\hat{\beta}_V$ is estimated by $\sqrt{l'(X'V^{-1}X)^{-1}l}$.

CONTACT INFORMATION

Dallas E. Johnson
 Department of Statistics
 101 Dickens Hall
 Kansas State University
 Manhattan, KS 66506-0802
 Work Phone: 785-532-0510
 Fax: 785-532-7736
 Email: dallas@stat.ksu.edu
 Web: <http://www.ksu.edu/stats/facultypages/johnson.htm>