

# Design Document: Monte Carlo Maximum Likelihood for Generalized Linear Mixed Models

Christina Knudson

January 13, 2015

## Abstract

This design document describes the process of performing Monte Carlo maximum likelihood (MCML) for generalized linear mixed models. First, penalized quasi-likelihood (PQL) estimates are calculated, which help generate simulated random effects. Then, the Monte Carlo likelihood approximation (MCLA) is calculated using the simulated random effects. Next, the MCLA is maximized to find Monte Carlo maximum likelihood estimates, the corresponding Fisher Information, and other statistics. Additional inference is then possible, including confidence intervals for the parameters and likelihood ratio tests for comparing nested models.

## 1 Theory

Let  $y = (y_1, \dots, y_n)'$  be a vector of observed data. Let  $u = (u_1, \dots, u_q)'$  be a vector of unobserved random effects centered at 0 with variance matrix  $D$ . Let  $\beta$  be a vector of  $p$  fixed effect parameters and let  $\nu$  be a vector of  $T$  variance components for the random effects so that  $D$  depends on  $\nu$ . Let  $\theta = (\beta, \nu)'$  be a vector containing all unknown parameters. Then the data  $y$  are distributed conditionally on the random effects according to  $f_\theta(y|u)$  and the random effects are distributed according to  $f_\theta(u)$ . Although  $f_\theta(u)$  does not actually depend on  $\beta$  and  $f_\theta(y|u)$  does not depend on  $\nu$ , we write both densities with  $\theta$  to keep notation simple in future equations.

Since  $u$  is unobservable, the log likelihood must be expressed by integrating out the random effects:

$$l(\theta) = \log \int f_\theta(y|u) f_\theta(u) du \quad (1)$$

For most datasets, this integral is intractable. In these cases, performing even basic inference on the likelihood is not possible. Rather than evaluating the integral, Geyer and Thompson (1992) suggest using a Monte Carlo approximation to the likelihood. Monte Carlo likelihood approximation

(MCLA) uses an importance sampling distribution  $\tilde{f}(u)$  to generate random effects  $u_k, k = 1, \dots, m$  where  $m$  is the Monte Carlo sample size. MCLA theoretically works for any  $\tilde{f}(u)$ , but the  $\tilde{f}(u)$  chosen for this package is the mixture distribution specified in section 5.3.

Then the Monte Carlo log likelihood approximation is

$$l_m(\theta) = \log \frac{1}{m} \sum_{k=1}^m f_\theta(y|u_k) \frac{f_\theta(u_k)}{\tilde{f}(u_k)} \quad (2)$$

$$= \log \frac{1}{m} \sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}. \quad (3)$$

When  $\tilde{f}$  does not depend on  $\theta$ , the gradient vector of the MCLA with respect to  $\theta$  is

$$\nabla l_m(\theta) = \frac{\sum_{k=1}^m \nabla \log f_\theta(u_k, y) \frac{f_\theta(u_k, y)}{\tilde{f}(u_k)}}{\sum_{k=1}^m \frac{f_\theta(u_k, y)}{\tilde{f}(u_k)}}, \quad (4)$$

and the Hessian matrix of the MCLA is

$$\begin{aligned} \nabla^2 l_m(\theta) &= \frac{\sum_{k=1}^m \nabla^2 \log f_\theta(y, u_k) \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}}{\sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}} \\ &+ \frac{\sum_{k=1}^m [\nabla \log f_\theta(y, u_k) - \nabla l_m(\theta)] [\nabla \log f_\theta(y, u_k) - \nabla l_m(\theta)]' \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}}{\sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}} \end{aligned} \quad (5)$$

Details for these derivatives can be found in section A.1. Calculation for the gradient and hessian when  $\tilde{f}$  is dependent on  $\theta$  can be found in section A.2.

Now any likelihood-based inference, such as maximum likelihood, can be performed on  $l_m(\theta)$  and its derivatives.

## 2 Model fitting function

The model fitting function is be the primary function the user uses. The user specifies the response and the predictors using the R formula mini-language as interpreted by `model.matrix`. Let  $n$  be the observed sample size and recall that  $p$  is the number of fixed effects. As a result of the user specifying the fixed effects, the model fitting function creates matrix  $X$ , which has dimensions  $n \times p$ .

The user specifies the random effects in the same way as the fixed effects. This is also how users of **reaster** (of R package **aster** (Geyer, 2014)) specify the random effects. That is, random effects are expressed using the R formula mini-language. Let  $q_t$  be the number of random effects associated with variance component  $\nu_t$ . When the random effects are specified, a list of  $T$  model matrices are created. The  $t$ th model matrix  $Z_t$  has dimensions  $n \times q_t$ .

The user also specifies the exponential family (details in section 3), the name of the data set, and the names of the variance components.

Thus, the following is a sample command with fixed predictors  $x_1$  and  $x_2$  and with random effects created by categorical variables *school* and *classroom*:

```
glmm(y ~ x1+ x2, list(~0+school,~0+classroom), family.glmm="bernoulli.glmm",
data=schooldat,varcomps.names=c("school","classroom"),varcomps.equal=c(1,2),
debug=FALSE)
```

It is possible that the user could want some variance components to be set equal. For example, Coull and Agresti's (2000) influenza dataset contains four years with random effects for each year. The authors want the within-year variance components to be equal. There are also variance components for subject-specific intercepts and for the decreased susceptibility to illness in year 4 (since the strain of flu during year 4 was a repeat of a previous year). Suppose *year* is a categorical variable with four levels, and *year1* through *year4* are dummy variables. Thus, the call contains these arguments:

```
glmm(y~year,list(~0+subject,~0+year1,~0+year2,~0+repeat,~0+year3,~0+year4),
varcomps.equal=c(1,2,2,3,2,2), varcomps.names=c("subject","year","repeat"),
data=flu,family.glmm=bernoulli.glmm)
```

Initially, `model.matrix` makes  $Z$  into a list of 6 design matrices. Since we have 3 distinct variance components, we want 3 design matrices. We take the design matrices that share a variance component and use `cbind` to bind them together. The result is 3 model matrices in the  $Z$  list, one for each variance component. We want to put them in order (1,2,3 according to `varcomps.equal`) so that the names for each model matrix ("subject", "year", "repeat") are correctly assigned. The variance estimates are eventually listed in the model summary in this same order as well.

After interpreting the model the user has specified, the next step is to find penalized quasi-likelihood (PQL) estimates. The process of finding these estimates is detailed in section 8. The PQL estimates parameterize the importance sampling distribution  $\tilde{f}(u_k)$  which generate the random effects. More information on generating the random effects is in section 5.3.

Next, `trust` is implemented to maximize the MCLA objective function (details on the objective function are in section 6). Finally, the function returns parameter estimates, the log likelihood evaluated at those estimates, the gradient vector of the log likelihood, the Hessian of the log likelihood at those estimates, information from the `trust` optimization, and other information.

### 3 Families

Let  $g$  be the canonical link function and  $\mu$  be a vector of length  $n$  such that

$$g(\mu) = X\beta + Zu \tag{6}$$

The choice of the link function is related to the distribution of the data,  $f_\theta(y|u)$ . If the data have a Bernoulli distribution, the link is  $\text{logit}(\mu)$ . If the data have a Poisson distribution, the link is  $\log(\mu)$ . Currently, `glmm` offers only these two choices for the family, but any exponential family could work and can be easily added later.

For simplicity of future notation, let  $\eta = g(\mu) = X\beta + Zu$ . Let  $c(\eta)$  denote the cumulant function such that the log of the data density can be written as

$$y'\eta - c(\eta) = \sum_i [y_i\eta_i - c(\eta_i)] \tag{7}$$

The user is required to specify the family in the model-fitting function. Once the family is specified, many family-specific functions are called. They are contained in an S3 class called “`glmm.family`”. Each family function outputs a list including the family name (a character string such as “`bernoulli.glmm`”), a function that calculates the value of the cumulant function  $c(\eta)$ , a function that calculates the cumulant’s first derivative  $c'(\eta)$  with the derivative taken with respect to  $\eta$ , and a function that calculates the cumulant’s second derivative  $c''(\eta)$ .

The users provide the family in the model-fitting function by either enter the character string (“`bernoulli.glmm`”), the function (`bernoulli.glmm()`), or the result of invoking the function. The following code for using the input to determine the family is adapted from `glm`.

```
logDensity<-function(family.glmm)
{
  if(is.character(family.glmm))
  family.glmm<-get(family.glmm,mode="function",envir=parent.frame())
  if(is.function(family.glmm))
  family.glmm<-family.glmm()
  if(!inherits(family.glmm,"glmm.family"))
```

```

stop(" 'family.glmm' not recognized")
return(family.glmm)
}

```

We interpret this as follows. If the user has entered the family as a string, go get the R object with that family name, either from the immediate environment or the parent environment. If this has happened, `family.glmm` is now a function. If `family.glmm` is a function (either because the user entered it as a function or because of the preceding step), invoke that function. At this point, `family.glmm` should have class “`glmm.family`.” If this is not the case (maybe because of a typo or maybe because they entered “`poisson`” rather than “`poisson.glmm`”), then stop and return an error.

With this family of functions, calculating  $c(\eta_i)$ ,  $c'(\eta_i)$ , and  $c''(\eta_i)$  is as simple as:

```

family.glmm$cum(args)
family.glmm$cp(args)
family.glmm$cpp(args)

```

For the Bernoulli distribution, we calculate these values (`cum`, `cp`, and `cpp`) as follows. In order to be careful with our computer arithmetic, we use the `log1p` function and the second set of equalities.

$$c(\eta_i) = \log(1 + e^{\eta_i}) = \begin{cases} \log(1 + e^{\eta_i}) & \text{if } \eta_i \leq 0, \\ \eta_i + \log(e^{-\eta_i} + 1) & \text{if } \eta_i > 0, \end{cases} \quad (8)$$

$$c'(\eta_i) = \frac{e^{\eta_i}}{1 + e^{\eta_i}} = \frac{1}{1 + e^{-\eta_i}} \quad (9)$$

$$c''(\eta_i) = \frac{e^{\eta_i}}{1 + e^{\eta_i}} - \frac{e^{2\eta_i}}{(1 + e^{\eta_i})^2} = \frac{1}{1 + e^{-\eta_i}} \cdot \frac{1}{1 + e^{\eta_i}} \quad (10)$$

For Poisson, these values (`cum`, `cp`, and `cpp`) are

$$c(\eta_i) = e^{\eta_i} \quad (11)$$

$$c'(\eta_i) = e^{\eta_i} \quad (12)$$

$$c''(\eta_i) = e^{\eta_i}. \quad (13)$$

Then we use these pieces to create the scalar  $c(\eta)$ , the vector  $c'(\eta)$  and the matrix  $c''(\eta)$ . We calculate

$$c(\eta) = \sum_i c(\eta_i). \quad (14)$$

The vector  $c'(\eta)$  has components  $c'(\eta_i)$ . The matrix  $c''(\eta)$  is diagonal with diagonal elements  $c''(\eta_i)$ .

In the R function `glm`, the user can choose the link. The canonical link must be used in the `glmm` package so that we have an exponential family. The canonical link is included in `family.glmm` in case the user does not know it already.

Also, this family of functions contains a check to ensure the data are valid given the family type. If `family.glmm` is `bernoulli.glmm`, the data should contain only 0 and 1. If `family.glmm` is `poisson.glmm`, then the data should be nonnegative. If the data set does not pass the check, the check returns an error message.

### 3.1 Redone in C

Redone in C, I have separate functions for calculating  $c(\eta)$ ,  $c'(\eta)$ ,  $c''(\eta)$ : `cum3.c`, `cp3.c`, and `cpp3.c`. The inputs for each function are identical: `eta` (an array of doubles), `neta` (the length of `eta`), the type (to denote the family: 1 indicates Bernoulli and 2 indicates Poisson), and an array of doubles to contain the result. These are passed in as pointers. The functions calculate the cumulant function or one of its first two derivatives. Each function contains a switch statement for the glmm family. The calculations for each of these functions have been shown earlier in this section. This function is type void: rather than returning the cumulant or its derivatives, the pointers are changed to contain the results. This function is invoked by `e1`, described in section 4.

## 4 Log density of the data (e1)

This section provides details for the log density of the data and two of its derivatives.

### 4.1 Equations

Recall the log of the data density is

$$\log f_{\theta}(y|u) = y'\eta + c(\eta) \quad (15)$$

$$= \sum_i y_i \eta_i - c(\eta_i) \quad (16)$$

where

$$\eta = X\beta + Zu. \quad (17)$$

The derivative of this with respect to one component,  $\eta_j$ , is

$$\frac{\partial}{\partial \eta_j} \log f_{\theta}(y|u) = y_j - c'(\eta_j). \quad (18)$$

The derivative of the component  $\eta_j$  with respect to one of the fixed effect predictors,  $\beta_l$ , is

$$\frac{\partial \eta_j}{\partial \beta_l} = X_{jl} \quad (19)$$

We'd like the derivative of the log of the data density with respect to  $\beta$ . The first step is using the chain rule as follows:

$$\frac{\partial}{\partial \beta_l} \log f_\theta(y|u) = \frac{\partial \eta_j}{\partial \beta_l} \frac{\partial}{\partial \eta_j} \log f_\theta(y|u) \quad (20)$$

$$= [y_j - c'(\eta_j)] X_{jl} \quad (21)$$

The mixed partial derivative (with respect to  $\beta_{l_1}$  and  $\beta_{l_2}$ ) of the log data density can be written similarly:

$$\frac{\partial^2}{\partial \beta_{l_1} \partial \beta_{l_2}} \log f_\theta(y|u) = \frac{\partial}{\partial \beta_{l_2}} ([y_j - c'(\eta_j)] X_{jl_1}) \quad (22)$$

$$= \frac{\partial \eta_j}{\partial \beta_{l_2}} \frac{\partial}{\partial \eta_j} ([y_j - c'(\eta_j)] X_{jl_1}) \quad (23)$$

$$= -X_{jl_1} X_{jl_2} c''(\eta_j) \quad (24)$$

Letting  $c'(\eta)$  be a vector with components  $c'(\eta_j)$ , the first derivative of the log data density can be written in matrix form as:

$$\frac{\partial}{\partial \beta} \log f_\theta(y|u) = X' [y - c'(\eta)] . \quad (25)$$

Letting  $c''(\eta)$  be a diagonal matrix with components  $c''(\eta_j)$ , the second derivative of the log data density can be written in matrix form as:

$$\frac{\partial^2}{\partial \beta^2} \log f_\theta(y|u) = X' [-c''(\eta)] X \quad (26)$$

## 4.2 Redone in C

The C function to calculate the value of the log data density and its two derivatives are called by reference. The following pointers are passed in: double **Y**, double **X**, int **nrowX**, int **ncolX**, double **eta**, int **family**, double **elval**, double **elgradient**, double **elhessian**. The pointers **elval**, **elgradient** and **elhessian** are zeros before **el.C** is invoked. Invoking **el.C** then places the calculated value of the log data density and two derivatives into **elval**, **elgradient** and **elhessian**.

The function **el.C** calls the following C functions: **cum3.C** to calculate the cumulant given a value of  $\eta$ , **cp3.C** to calculate the derivative of the cumulant given a value of  $\eta$ , **cpp3.C** to calculate the hessian of the cumulant given a value of  $\eta$ , and functions to perform matrix multiplication.

## 5 Random effect generation and calculations

This section is focused on a vector of random effects  $u$  with length  $q$ . Section 5.1 details the relationship between  $\nu$  and  $D$ . Section 5.2 states the importance sampling distribution. Section 5.3 describes the process of generating the random effects. Section 5.4 explains how to evaluate the distribution of the random effects and its first two derivatives.

### 5.1 Constructing $D$ (`getEk`)

Recall that, for this version of the package, we assume that  $D$  is diagonal. Let  $\nu$  be length  $T$  with components  $\nu_t$ , and let  $E_t, t = 1, \dots, T$  be diagonal matrices with indicators on the diagonal so that  $\sum_{t=1}^T E_t = I$ . That is, the diagonal entries of  $E_t$  indicate whether that random effect has  $\nu_t$  as its variance component. Then  $D = \sum_{t=1}^T \nu_t E_t$ .

Recall  $q_t$  is the number of random effects associated with variance component  $\nu_t$ . Then  $q_t$  is also the number of nonzero entries in  $E_t$  and  $q = \sum_{t=1}^T q_t$  is the total number of random effects in the model. The model fitting function of `glmm` has made  $Z$  into a list with  $T$  design matrices (one for each distinct variance component). In order to create  $E_t$ , we need to go through and count the number of columns ( $q_t$ ) for each design matrix  $t$  in the list. Then  $q = \sum_{t=1}^T q_t$  is the total number of random effects in the model, and  $D$  is  $q \times q$ . Thus, each  $E_t$  is  $q \times q$ .  $E_1$  has  $q_1$  ones on the diagonal, followed by zeros to fill the rest of the diagonal.  $E_2$  has  $q_1$  zeros, then  $q_2$  ones, then zeros for the rest of the diagonal. We continue this process to complete the construction of  $E_t, t = 1, \dots, T$ .

### 5.2 Importance sampling distribution $\tilde{f}(u_k)$

Let  $s$  be a vector of length  $q$  that represents the random effects on the standard normal scale. That is  $s$  is defined such that,  $u = D^{1/2}s$ . Let  $\sigma$  be a vector of length  $T$  with components  $\sqrt{\nu_t}$ . Prior to generating random effects, the model fitting function of `glmm` has performed PQL and recorded the PQL estimates for  $\beta$ ,  $\sigma$  and  $s$ . The PQL estimates are denoted by  $\beta^*$ ,  $s^*$  and  $\sigma^*$ . Let

$$A^* = \sum_{t=1}^T E_t \sigma_t^* \quad (27)$$

and

$$D^* = A^* A^* \quad (28)$$

be matrices based on PQL estimates. We can “unstandardize” our PQL-based random effects:

$$u^* = A^* s^*. \quad (29)$$



Let  $p_1, p_2, p_3$  be the proportions of the mixture distribution such that  $p_1 + p_2 + p_3 = 1$ . Let  $f(u|\mu, \Sigma)$  denote the pdf for  $u \sim N(\mu, \Sigma)$ . Let  $\dot{f}(u)$  denote the pdf for  $u \sim t_\zeta$ , a  $q$ -dimensional multivariate standard  $t$  with  $\zeta$  degrees of freedom. Then the importance sampling distribution is:

$$\tilde{f}(u) = p_1 \dot{f}(u) + p_2 f(u|u^*, D^*) + p_3 f(u|u^*, (Z'c''(X\beta^* + Zu^*)Z + (D^*)^{-1})^{-1}). \quad (30)$$

Then

$$\log \tilde{f}(u) = \log \left( p_1 \dot{f}(u) + p_2 f(u|u^*, D^*) + p_3 f(u|u^*, (Z'c''(X\beta^* + Zu^*)Z + (D^*)^{-1})^{-1}) \right) \quad (31)$$

The first component of the mixture distribution  $t_\zeta$  is chosen to ensure the gradient of the MCLA has a central limit theorem, which is proven in section B. The second component is chosen because it is centered at the PQL best guess of the random effect values  $u^*$  and has the PQL guess of the variance. The last component is centered at the PQL guess  $u^*$  and has a variance based on the Hessian of the PQL penalized likelihood. The idea of this last distribution is to generate random effects from a distribution whose Hessian matches that of the target distribution  $f_\theta(y|u)f_\theta(u)$ .

### 5.3 Generating random effects (genRand)

Recall that  $m$  is the overall Monte Carlo sample size. Let  $p_1, p_2, p_3$  be the proportions of the mixture distribution such that  $p_1 + p_2 + p_3 = 1$ . Let  $m_1 = mp_1$ ,  $m_2 = mp_2$ ,  $m_3 = mp_3$ . Then we draw

$$u_k \sim t_\zeta, k = 1, \dots, m_1 \quad (32)$$

$$u_k \sim N(u|u^*, D^*), k = m_1 + 1, \dots, m_1 + m_2 \quad (33)$$

$$u_k \sim N(u|u^*, (Z'c''(X\beta^* + Zu^*)Z + (D^*)^{-1})^{-1}), k = m_1 + m_2 + 1, \dots, m. \quad (34)$$

Details for drawing from a nonstandard normal are in section 5.3.2.

#### 5.3.1 Finding the square root of a variance matrix

In order to generate draws from a nonstandard normal distribution, we need to calculate the square root of a variance matrix. When the variance matrix is diagonal (for example,  $D^*$ ), we simply take the root of the diagonal elements. When the variance matrix is not diagonal (for example, in the second component of the mixture distribution), we can use eigendecomposition. Eigendecompositions take a little bit more time than Cholesky decompositions but are more stable.

Let  $\Sigma$  denote the nondiagonal variance matrix. Eigendecomposition of  $\Sigma$  provides orthogonal matrix  $O$  (containing the eigenvectors) and diagonal matrix  $\Lambda$  (with diagonal entries  $\lambda$  being the

eigenvalues of  $\Sigma$ ). Then

$$\Sigma^{1/2} = O\Lambda^{1/2}O'. \quad (35)$$

Finding  $\Lambda^{1/2}$  is as easy as taking the root of the diagonal entries. We know that  $\Sigma^{1/2}$  is correct because

$$\begin{aligned} \Sigma^{1/2}\Sigma^{1/2} &= O\Lambda^{1/2}O'O\Lambda^{1/2}O' \\ &= O\Lambda^{1/2}\Lambda^{1/2}O' \\ &= O\Lambda O' \\ &= \Sigma. \end{aligned} \quad (36)$$

### 5.3.2 Generating from nonstandard normal distributions

Construct  $A^* = (D^*)^{1/2}$  and  $((Z'c''(X\beta^* + Zu^*)Z + (D^*)^{-1})^{-1})^{1/2}$  as described in section 5.3.1.

Let  $\check{u}_k, k = m_1 + 1, \dots, m$  be vectors of random effects that are drawn from the standard normal distribution. Then we center and scale  $\check{u}_k$  as follows:

$$u_k = u^* + A^* \check{u}_k, k = m_1 + 1, \dots, m_1 + m_2 \quad (37)$$

$$u_k = u^* + ((Z'c''(X\beta^* + Zu^*)Z + (D^*)^{-1})^{-1})^{1/2} \check{u}_k, k = m_1 + m_2 + 1, \dots, m. \quad (38)$$

The result of centering and scaling is that

$$u_k \sim N(u^*, D^*), k = m_1 + 1, \dots, m_1 + m_2 \quad (39)$$

$$u_k \sim N(u^*, (Z'c''(X\beta^* + Zu^*)Z + (D^*)^{-1})^{-1}), k = m_1 + m_2 + 1, \dots, m. \quad (40)$$

## 5.4 Evaluating the distribution of random effects (distRand)

This section discusses evaluating the distribution of the random effects, the gradient of the random effects, and the Hessian of the random effects. That is, the equations in this section provide  $\log f_\theta(u)$ ,  $\nabla \log f_\theta(u)$ , and  $\nabla^2 \log f_\theta(u)$  for equation 58.

We can express the log unnormalized pdf of  $N(0, D)$  as:

$$\log f_\theta(u) = \log f(u|0, D) = (-1/2) \log |D| - (1/2)u'D^{-1}u \quad (41)$$

To find derivatives of this, we are going to use the diagonal form of  $D$ . Recall that for every  $\nu_t$ , we can construct a matrix  $E_t$  (with dimensions the same as matrix  $D$ ) that has 1s on the diagonal elements corresponding to the elements of  $D$  that contain  $\nu_t$  and 0s elsewhere.

We can partition the random effects according to their variance components:  $u = (U'_1, \dots, U'_T)'$ . Let  $D_t$  be the variance matrix for  $U_t$ .  $D_t$  has  $q_t$  rows and  $q_t$  columns. Thus  $D$  can be expressed as:

$$D = \begin{bmatrix} D_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & D_T \end{bmatrix} \quad (42)$$

Since  $D$  is diagonal, it follows that  $D^{-1}$  is also diagonal with diagonal entries  $\frac{1}{\nu_1}, \dots, \frac{1}{\nu_T}$ . Also, the assumption that  $D$  is diagonal makes calculating the determinant of  $D$  easy:

$$|D| = \nu_1^{q_1} \dots \nu_T^{q_T} \quad (43)$$

Taking these two pieces of information into account allows us to write the log density for  $U$  as follows:

$$\log f_\theta(u) = -\frac{1}{2} \log |D| - \frac{1}{2} u' D^{-1} u \quad (44)$$

$$= -\frac{1}{2} \left[ \sum_{t=1}^T q_t \log \nu_t \right] - \frac{1}{2} \sum_{t=1}^T \left[ \frac{1}{\nu_t} U'_t U_t \right] \quad (45)$$

The first and second derivatives of each summand with respect to its associated  $\nu_t$  are:

$$\frac{\partial}{\partial \nu_t} \log f_\theta(u_t) = -\frac{q_t}{2\nu_t} + \frac{1}{2\nu_t^2} U'_t U_t \quad (46)$$

and

$$\frac{\partial^2}{\partial \nu_t^2} \log f_\theta(u_t) = \frac{q_t}{2\nu_t^2} - \frac{1}{\nu_t^3} U'_t U_t. \quad (47)$$

Any other derivative is equal to 0. That is, for all  $t_1 \neq t_2$ ,

$$\frac{\partial}{\partial \nu_{t_1}} \log f_\theta(u_{t_2}) = 0. \quad (48)$$

Also,

$$\frac{\partial}{\partial \beta} \log f_\theta(u_{t_2}) = 0. \quad (49)$$

Thus, if  $\nu = (\nu_1, \dots, \nu_T)$ , the gradient of the random effects distribution is the following vector of length  $T$ :

$$\frac{\partial}{\partial \nu} \log f_\theta(u) = \left[ -\frac{q_1}{2\nu_1} + \frac{1}{2\nu_1^2} U'_1 U_1 \quad \dots \quad -\frac{q_T}{2\nu_T} + \frac{1}{2\nu_T^2} U'_T U_T \right] \quad (50)$$

The Hessian matrix is the following diagonal matrix:

$$\frac{\partial^2}{\partial \nu^2} \log f_\theta(u) = \begin{bmatrix} \frac{q_1}{2\nu_1^2} - \frac{1}{\nu_1^3} U_1' U_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{q_T}{2\nu_T^2} - \frac{1}{\nu_T^3} U_T' U_T \end{bmatrix} \quad (51)$$

To calculate the value, gradient, and Hessian, we need to provide  $\nu$ ,  $u$ , and the list  $\mathbf{z}$  (from `mod.mcm1`). The list  $\mathbf{z}$  has  $T$  matrices, each with the number of columns equal to  $q_t$ . We need  $q_t, t = 1, \dots, T$  to calculate the log density and its derivatives.

The last thing we need to discuss is how to split  $U$  into  $U_1, \dots, U_T$ . We know that the first  $q_1$  items of  $U$  are  $U_1$ , the next  $q_2$  items are  $U_2$ , etc. In other words, entries 1 through  $q_1$  are  $U_1$ . Items  $q_1 + 1$  through  $q_1 + q_2$  are  $U_2$ , etc. We find these numbers ( $q_1, q_2$ , etc) from the number of columns of the items in the list  $\mathbf{z}$ . This information of how to split  $u$  up is contained in vector `meow` and used by function `distRand3`. This is faster than recalculating `meow` for every call of `distRand3`.

Rewritten in C, this function takes the following as pointers: the double array `nu` that contains the variance components, the int `T` to specify the length of `nu`, the double array `mu` that contains the means of the random effects, the int array `nrandom` of length `T` that contains the number random effects from that variance component, the double array `Uvec` that contains one vector of generated random effects, the int array `meow` that specifies how to split  $u$  up based on the variance components, the double array `drgradient` that contains the resulting gradient, and the double array `drhessian` that contains the resulting hessian.

The result of invoking the C function `distRand3` is the calculation of the gradient and hessian for the distribution of the random effects. The value is evaluated by the C function `distRandGeneral`, described in section 5.5.1.

## 5.5 Evaluating the importance sampling distribution

This section discusses evaluating the importance sampling distribution shown in 30. Section 5.5.1 describes how to evaluate the pdf for the normal components of  $\tilde{f}(u)$ . Section 5.5.2 describes how to evaluate the pdf for the t component of  $\tilde{f}(u)$ . The  $\tilde{f}(u)$  calculation is used in equation 58.

### 5.5.1 Normal distribution for general variance matrix $\Sigma$ (`distRandGeneral`)

Let  $\Sigma$  be a variance matrix and let  $\mu$  be a mean vector. Consider  $N(\mu, \Sigma)$ . We can write the log pdf of this distribution as:

$$\log f(u|\mu, \Sigma) = (1/2) \log |\Sigma^{-1}| - (1/2)(u - \mu)' \Sigma^{-1} (u - \mu) \quad (52)$$

The only part of this to discuss is  $|\Sigma^{-1}|$ . We can use eigendecomposition to make  $\Sigma^{-1} = O\Lambda O'$  where  $O$  is orthogonal and  $\Lambda$  is the diagonal matrix with eigenvalues. Since orthogonal matrices have determinant  $\pm 1$ , then  $|O||O'| = 1$ . Thus

$$|\Sigma^{-1}| = |O'| |\Lambda| |O| \quad (53)$$

$$= |O'| |O| |\Lambda| \quad (54)$$

$$= |\Lambda|, \quad (55)$$

which is just the product of the eigenvalues. The log of the determinant is calculated beforehand to save time, since this function is called  $3m$  times throughout each trust optimization iteration, where  $m$  is again the Monte Carlo sample size.

This function is also rewritten in C with the following passed in as pointers: double `Sigma.inv`  $\Sigma^{-1}$ , double `logdet`  $\log |\Sigma^{-1}|$ , int `nrow`, double `uvec` a vector of random effects, double `mu`  $\mu$ , and double `distRandGenVal`.

### 5.5.2 t distribution

Consider  $t_\zeta$ . We can write the log unnormalized pdf of this distribution as:

$$\log \hat{f}(u) = [1 + u'u/\zeta]^{-(\zeta+q)/2} \quad (56)$$

## 6 Objective function: approximated log likelihood

The objective function is optimized by `trust` within the the model fitting function. In order to evaluate and maximize the approximated log likelihood, I need an objective function that returns the value, the gradient vector, and the Hessian matrix of the approximated log likelihood. Recall equations 3, 4, ?? for calculating these quantities. This objective function needs  $\log f_\theta(u_k)$ ,  $c(\eta_i)$ ,  $\tilde{f}(u_k)$ , each of their gradient vectors, and each of their Hessian matrices to plug into these expressions.

Denote

$$b_k = \log f_\theta(u_k, y) - \log \tilde{f}(u_k) \quad (57)$$

$$= \log f_\theta(u_k) + \log f_\theta(y|u_k) - \log \tilde{f}(u_k) \quad (58)$$

For computational stability, we'll set  $a = \max(b_k)$ . Then the value of the MCLA is expressed as:

$$l_m(\theta) = a + \log \left[ \frac{1}{m} \sum_{k=1}^m e^{b_k - a} \right] \quad (59)$$

Define the weights as:

$$v_\theta(u_k, y) = \frac{e^{b_k - a}}{\sum_{k=1}^m e^{b_k - a}} \quad (60)$$

Rewriting equation 4 using the notation for the weights, the gradient vector is:

$$\nabla l_m(\theta) = \sum_{k=1}^m \left( \nabla \log f_\theta(u_k, y) - \nabla \log \tilde{f}(u_k) \right) v_\theta(u_k, y) \quad (61)$$

$$= \sum_{k=1}^m \nabla \left[ \log f_\theta(y|u_k) + \log f_\theta(u_k) - \nabla \log \tilde{f}(u_k) \right] v_\theta(u_k, y) \quad (62)$$

$$= \sum_{k=1}^m \left[ \nabla \log f_\theta(y|u_k) + \nabla \log f_\theta(u_k) - \nabla \log \tilde{f}(u_k) \right] v_\theta(u_k, y) \quad (63)$$

$$= \sum_{k=1}^m \left[ \frac{\partial}{\partial \beta} \log f_\theta(y|u_k) \quad \frac{\partial}{\partial \nu} \log f_\theta(y|u_k) \right] v_\theta(u_k, y) + \left[ \frac{\partial}{\partial \beta} \log f_\theta(u_k) \quad \frac{\partial}{\partial \nu} \log f_\theta(u_k) \right] v_\theta(u_k, y) \quad (64)$$

$$+ \sum_{k=1}^m \left[ \frac{\partial}{\partial \beta} \log \tilde{f}(u_k) \quad \frac{\partial}{\partial \nu} \log \tilde{f}(u_k) \right] v_\theta(u_k, y) \quad (65)$$

$$= \sum_{k=1}^m \left( \left[ \frac{\partial}{\partial \beta} \log f_\theta(y|u_k) \quad 0 \right] + \left[ 0 \quad \frac{\partial}{\partial \nu} \log f_\theta(u_k) \right] + \left[ 0 \quad \frac{\partial}{\partial \nu} \log \tilde{f}(u_k) \right] \right) v_\theta(u_k, y) \quad (66)$$

$$= \sum_{k=1}^m \left[ \frac{\partial}{\partial \beta} \log f_\theta(y|u_k) \quad \left( \frac{\partial}{\partial \nu} \log f_\theta(u_k) \right) + \left( \frac{\partial}{\partial \nu} \log \tilde{f}(u_k) \right) \right] v_\theta(u_k, y) \quad (67)$$

Then the Hessian matrix from equation ?? is expressed as:

$$\nabla^2 l_m(\theta) = \sum_{k=1}^m \left[ \nabla^2 \log f_\theta(y, u_k) - \nabla^2 \log \tilde{f}(u_k) \right] v_\theta(u_k, y) \quad (68)$$

$$+ \sum_{k=1}^m \left[ \nabla \log f_\theta(y, u_k) - \nabla \log \tilde{f}(u_k) - \nabla l_m(\theta) \right] \left[ \nabla \log f_\theta(y, u_k) - \nabla \log \tilde{f}(u_k) - \nabla l_m(\theta) \right]' v_\theta(u_k, y).$$

Everything in this has already been defined except:

$$\nabla^2 \log f_\theta(u_k, y) = \begin{bmatrix} \frac{\partial^2}{\partial \beta^2} \log f_\theta(u_k, y) & \frac{\partial^2}{\partial \beta \partial \nu} \log f_\theta(u_k, y) \\ \frac{\partial^2}{\partial \beta \partial \nu} \log f_\theta(u_k, y) & \frac{\partial^2}{\partial \nu^2} \log f_\theta(u_k, y) \end{bmatrix} \quad (69)$$

$$= \begin{bmatrix} \frac{\partial^2}{\partial \beta^2} \log f_\theta(u_k) & \frac{\partial^2}{\partial \beta \partial \nu} \log f_\theta(u_k) \\ \frac{\partial^2}{\partial \beta \partial \nu} \log f_\theta(u_k) & \frac{\partial^2}{\partial \nu^2} \log f_\theta(u_k) \end{bmatrix} + \begin{bmatrix} \frac{\partial^2}{\partial \beta^2} \log f_\theta(y|u_k) & \frac{\partial^2}{\partial \beta \partial \nu} \log f_\theta(y|u_k) \\ \frac{\partial^2}{\partial \beta \partial \nu} \log f_\theta(y|u_k) & \frac{\partial^2}{\partial \nu^2} \log f_\theta(y|u_k) \end{bmatrix} \quad (70)$$

$$= \begin{bmatrix} 0 & 0 \\ 0 & \frac{\partial^2}{\partial \nu^2} \log f_\theta(u_k) \end{bmatrix} + \begin{bmatrix} \frac{\partial^2}{\partial \beta^2} \log f_\theta(y|u_k) & 0 \\ 0 & 0 \end{bmatrix} \quad (71)$$

$$= \begin{bmatrix} \frac{\partial^2}{\partial \beta^2} \log f_\theta(y|u_k) & 0 \\ 0 & \frac{\partial^2}{\partial \nu^2} \log f_\theta(u_k) \end{bmatrix} \quad (72)$$

## 7 Summary of model

Typing `summary(mod)` gives more detailed info of the model. This is broken into two pieces (as is usually done for summaries): `summary.glmm` and `print.summary.glmm`. We split this up because we have a lot of information in the summary that we do not necessarily want printed each time. When a user types `summary(mod)`, only the basic information is automatically printed. More information can be found in the summary list.

The summary performs all calculations and its value is a list of the following:

- call
- the variance estimate(s)
- a matrix with the predictor in the first column, the estimated coefficient in the second column, the standard error in the third column, the z value in the fourth column, and the two-sided pvalue in the fifth column. All inference is asymptotic, so we use the standard normal distribution to calculate the p-values.
- the evaluated Monte Carlo log likelihood along with its first and second derivative
- maybe other things that I haven't thought of

A note on the standard errors: to calculate the standard errors, we take the MCLA Hessian matrix, invert it, take the diagonal elements, and square root them. It's possible that the Hessian

will be noninvertible if it is close enough to singular that the computer thinks it's singular. Then the standard errors will all be infinite. We also need to warn the user why this is happening.

Then `print.summary.glmm` prints the following:

- call
- the variance estimate(s)
- a matrix similar to the output of `summary.glm`, which will be put through `printCoefmat` in order to get the significance stars that we're used to. We'll have the predictor in the first column, the estimated coefficient in the second column, the standard error in the third column, the z value in the fourth column, and the two-sided pvalue in the fifth column, the significance stars, and the significance legend.
- maybe other things that I haven't thought of. I'll probably realize more once I get coding.

I think for the `printCoefmat` to work, the fixed effects matrix needs to look like

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	2	0.50	4	6.334248e-05
x1	2	0.25	8	1.244192e-15

(or at least it worked when I had these as the column names and didn't work before that).

Note that the `summary` and `print.summary` functions will be S3 generic functions. This means is the user types `summary(mod)` and `summary` is a generic function. R checks the class of the model `mod` and then automatically uses the `summary` and `print.summary` functions for that class of objects.

## 8 PQL

Before we get started on describing PQL, we need to change notation to avoid constrained optimization. Recall that  $D = \text{Var}(u)$  and is (for now) assumed to be diagonal. Let  $A = D^{1/2}$  so that A has diagonal components that are positive or negative. Using A rather than D enables unconstrained optimization. If  $\sigma$  is a vector of the distinct standard deviations with components  $\sigma_t$ , we can write A as a function of  $\sigma$  by

$$A = \sum_t E_t \sigma_t.$$

Recall that  $E_t$  has a diagonal of indicators to show which random effects have the same variance components, and  $\sum_{t=1}^T E_t$  is the identity matrix. PQL will estimate the components contained



on the diagonal of  $A$ . Taking the absolute value of those components will provide the standard deviations (the square root of the variance components).

In addition to using  $A$  rather than  $D$ , the other change is that we'll be using  $s$  where  $u = As$ . The purpose of this is to avoid  $D^{-1/2}$  in the function we're trying to optimize.

There are two ways to do PQL, both of which are described in the vignette `re.pdf` in the R package `aster` (Geyer, 2014). In either case, there is an inner optimization and an outer optimization. The inner optimization is well behaved while the other optimization is a little tougher. I will be using the method that is not quite PQL but is pretty close and is better behaved. In this version, the inner optimization finds  $\tilde{\beta}$  and  $\tilde{s}$  given  $X$ ,  $Z$  and  $A$ . Then, given  $\tilde{\beta}$  and  $\tilde{s}$ , the outer optimization finds  $A$ .

The **inner** optimization will be done with the trust function in R. We elect to use trust because it requires two derivatives, which will make the optimization more precise. We would like more accuracy in the inner maximization because any sloppiness will carry into the outer optimization.

The inner optimization maximizes the penalized log likelihood. After defining

$$\eta = X\beta + ZAs \tag{73}$$

we calculate the log likelihood as

$$l(\eta) = Y'\eta - c(\eta) \tag{74}$$

and the penalized likelihood as:

$$l(\eta) - \frac{1}{2}s's. \tag{75}$$

This function will be maximized using the `trust` package. We need to give `trust` derivatives with respect to  $s$  and  $\beta$ . We're going to express these via the multivariate chain rule, taking advantage of  $\eta_i$ .

Create vector  $\mu$  from the components  $\mu_i = c'(\eta_i)$ . Since

$$l(\eta) = \sum_i Y_i \eta_i - c(\eta_i) \tag{76}$$

then

$$\frac{\partial l(\eta)}{\partial \eta_i} = Y_i - c'(\eta_i) = Y_i - \mu_i. \tag{77}$$

Now we can write the following expression:

$$\frac{\partial}{\partial \beta_k} \left[ l(\eta) - \frac{1}{2} s' s \right] = \frac{\partial l(\eta)}{\partial \beta_k} \quad (78)$$

$$= \frac{\partial l(\eta)}{\partial \eta_i} \frac{\partial \eta_i}{\partial \beta_k} \quad (79)$$

$$= \sum_i (Y_i - \mu_i) \frac{\partial \eta_i}{\partial \beta_k} \quad (80)$$

$$= \sum_i (Y_i - \mu_i) X_{ik} \quad (81)$$

We find the function's derivative with respect to  $s$  as follows:

$$\frac{\partial}{\partial s} \left[ l(\eta) - \frac{1}{2} s' s \right] = \frac{\partial l(\eta)}{\partial s} - \frac{1}{2} \frac{\partial s' s}{\partial s} \quad (82)$$

$$= \frac{\partial l(\eta)}{\partial \eta} \frac{\partial \eta}{\partial s} - s \quad (83)$$

$$= (Y - \mu)' \left[ \frac{\partial}{\partial s} Z A s \right] - s \quad (84)$$

$$= (Y - \mu)' Z A - s \quad (85)$$

This gives us the following derivatives of the penalized log likelihood:

$$\frac{\partial}{\partial \beta} [l(\eta) - (1/2) s' s] = X' (Y - \mu) \quad (86)$$

$$\frac{\partial}{\partial s} [l(\eta) - (1/2) s' s] = A Z' (Y - \mu) - s \quad (87)$$

Lastly, we need the Hessian of the penalized likelihood. This matrix can be broken down into four pieces:

$$1. \frac{\partial^2}{\partial s^2}$$

$$2. \frac{\partial^2}{\partial \beta^2}$$

$$3. \frac{\partial^2}{\partial s \partial \beta}$$

$$4. \left( \frac{\partial^2}{\partial s \partial \beta} \right)' = \frac{\partial^2}{\partial \beta \partial s}$$

We'll start at the top with  $\frac{\partial^2}{\partial s^2}$ , which is a  $q \times q$  matrix.

$$\frac{\partial^2}{\partial s^2} [l(\eta) - (1/2)s's] = \frac{\partial}{\partial s} [(Y - c'(\eta))'ZA - s] \quad (88)$$

$$= \left[ -\frac{\partial}{\partial s} c'(\eta) \right]' ZA - I_q \quad (89)$$

$$= \left[ -\frac{\partial c'(\eta)}{\partial \eta} \frac{\partial \eta}{\partial s} \right]' ZA - I_q \quad (90)$$

$$= [-c''(\eta)ZA]' ZA - I_q \quad (91)$$

$$= -AZ' [c''(\eta)] ZA - I_q \quad (92)$$

Note that  $c''(\eta)$  is a  $q \times q$  diagonal matrix with diagonal elements  $c''(\eta_i)$ .  $I_q$  is the identity matrix of dimension  $q$ . This makes  $\frac{\partial^2}{\partial s^2}$  a  $q \times q$  matrix.

Next up is  $\frac{\partial^2}{\partial \beta^2}$ , which is a  $p \times p$  matrix.

$$\frac{\partial^2}{\partial \beta^2} [l(\eta) - (1/2)s's] = \frac{\partial}{\partial \beta} [X'(Y - c'(\eta))] \quad (93)$$

$$= \frac{\partial}{\partial \beta} [X'Y - X'(c'(\eta))] \quad (94)$$

$$= \frac{\partial}{\partial \beta} [-X'(c'(\eta))] \quad (95)$$

$$= -X' \left[ \frac{\partial}{\partial \beta} c'(\eta) \right] \quad (96)$$

$$= -X' \left[ \frac{\partial c'(\eta)}{\partial \eta} \frac{\partial \eta}{\partial \beta} \right] \quad (97)$$

$$= -X' [c''(\eta)] X \quad (98)$$

Next is the  $p \times q$  mixed partial  $\frac{\partial^2}{\partial \beta \partial s}$ .

$$\frac{\partial^2}{\partial \beta \partial s} [l(\eta) - (1/2)s's] = \frac{\partial}{\partial \beta} \left\{ [Y - c'(\eta)]' Z A \right\} \quad (99)$$

$$= \frac{\partial}{\partial \beta} \left\{ Y' Z A - [c'(\eta)]' Z A \right\} \quad (100)$$

$$= -\frac{\partial}{\partial \beta} \left\{ [c'(\eta)]' Z A \right\} \quad (101)$$

$$= -\left[ \frac{\partial c'(\eta)}{\partial \beta} \right]' Z A \quad (102)$$

$$= -\left[ \frac{\partial c'(\eta)}{\partial \eta} \frac{\partial \eta}{\partial \beta} \right]' Z A \quad (103)$$

$$= -[c''(\eta)X]' Z A \quad (104)$$

$$= -X' [c''(\eta)] Z A \quad (105)$$

And last we have the  $q \times p$  mixed partial  $\frac{\partial^2}{\partial \beta \partial s} = -AZ'[c''(\eta)]X$ . These four pieces specify the hessian matrix for the penalized likelihood. Trust can now perform the inner optimization to find  $\tilde{\beta}$  and  $\tilde{s}$ . Trust will maximize the penalized likelihood.

The **outer** optimization is done using **optim** with the default method of “Nelder-Mead.” This requires just the function value and no derivatives. This optimization method was chosen because the optimization function already contains second derivatives of the cumulant function; requiring derivatives of the optimization function would in turn require higher-order derivatives of the cumulant.

The default of **optim** is to minimize, but we’d like to do maximization. Reversing the sign of the optimization function will turn the maximization into minimization.

If  $\tilde{\beta}$  and  $\tilde{s}$  are available from previous calls to the inner optimization function, then they are used here. Otherwise, the values are taken to be 0 and 1. The outer optimization’s function is evaluated by first defining

$$\tilde{\eta} = X\tilde{\beta} + ZA\tilde{s} \quad (106)$$

$$l(\tilde{\eta}) = Y'\tilde{\eta} - c(\tilde{\eta}) \quad (107)$$

$$A = \sum_k E_k \sigma_k. \quad (108)$$

Let  $W$  be a diagonal matrix with elements  $c''(\eta_i)$  on the diagonal. Then the quantity we’d like to

maximize is the penalized quasi-likelihood:

$$l(\tilde{\eta}) - \frac{1}{2}\tilde{s}'\tilde{s} - \frac{1}{2}\log|AZ'WZA + I| \quad (109)$$

$$(110)$$

Again, `optim` minimizes, so we have to minimize the negative value of the penalized quasi-likelihood in order to maximize the value of it.

We do need to be careful about the determinant. Let  $AZ'WZA + I = LL'$  where  $L$  is the lower triangular matrix resulting from a Cholesky decomposition. Then

$$\frac{1}{2}\log|AZ'WZA + I| = \frac{1}{2}\log|LL'| \quad (111)$$

$$= \frac{1}{2}\log(|L|)^2 \quad (112)$$

$$= \log|L| \quad (113)$$

Since  $L$  is triangular, the determinant is just the product of the diagonal elements. Let  $l_i$  be the diagonal elements of  $L$ . Then

$$\frac{1}{2}\log|AZ'WZA + I| = \log\prod l_i \quad (114)$$

$$= \sum \log l_i \quad (115)$$

If I am worried about all variance components being zero, I could implement an eigendecomposition using the R function `eigen`. This would be more numerically stable, but is slower. Let  $O$  be the matrix containing the eigenvectors and let  $\Lambda$  be a diagonal matrix with the eigenvalues  $\lambda_i$  on the diagonal. Then

$$AZ'WZA = O\Lambda O' \quad (116)$$

Then we can rewrite the argument of the determinant as follows:

$$AZ'WZA + I = O\Lambda O' + I = O\Lambda O + OO' = O(I + \Lambda)O' \quad (117)$$

This leads to the careful calculation of our determinant as follows:

$$|AZ'WZA + I| = 1 * \prod_{i=1}^n (1 + \lambda_i) * 1 \quad (118)$$

$$\Rightarrow \log|AZ'WZA + I| = \sum \log(1 + \lambda_i) \quad (119)$$

The last quantity can be accurately calculated using the `log1p` function in R.

The outer optimization uses  $\tilde{\beta}$  and  $\tilde{s}$  provided by the inner optimization, but it does not return them. To keep track of the most recent  $\tilde{s}$ , so store them in an environment that I call “cache.” The

purpose of  $\tilde{\beta}$  and  $\tilde{s}$  is two-fold. First, if they are available from a previous iteration of the inner optimization, then they are used in the outer optimization of PQL. Second, after PQL is finished,  $\tilde{s}$  is used to help center the generated random effects.

The point of doing PQL is to construct a decent importance sampling distribution. Thus, the estimates don't have to be perfect. It is possible that one of the  $\sigma_k$  will be 0 according to PQL. If this happens, then I'll just use  $\sigma_k = .01$  or something like that for the importance sampling distribution.

## 9 Checks

### 9.1 Checking the MCLA finite differences

To check the function that calculates the MCLA  $l_m(\theta)$ , I use finite differences on the Booth and Hobert (1999) example. To do this, I chose a value of  $\theta = (\beta, \sigma)$  and a relatively small value of  $\delta$ , where  $\delta$  is a vector of length 2. We can check that the value and first derivative of the MCLA function are calculated correctly by making sure the following approximation holds

$$\nabla l_m(\theta) \cdot \delta \approx l_m(\theta + \delta) - l_m(\theta). \quad (120)$$

Then, we can check the first and second derivatives of the MCLA are calculated correctly by checking for the following approximation:

$$\nabla^2 l_m(\theta) \delta \approx \nabla l_m(\theta + \delta) - \nabla l_m(\theta) \quad (121)$$

### 9.2 Checking functions using the Booth and Hobert example

I also test the objective function by taking the Booth and Hobert (1999) example and rewriting the functions for this specific example. I then compare the values produced by the check functions and the original functions. Functions checked this way are:

1. log of the data density (`e1`)
2. log of the density for the random effects (both `distRand` and `distRandGeneral`)
3. the Monte Carlo likelihood approximation (`objfun`)

### 9.3 Checking a putative MLE using MCMC

Suppose  $\hat{\theta}$  is claimed to be the MLE. For example,  $\hat{\theta}$  could be the MCMLE. Ideally, we would like to check whether the true likelihood  $l(\theta)$  achieves a local max at  $\hat{\theta}$ . Due to the intractable

integral in the likelihood expression, the best we can do is make sure that  $\nabla l_m(\theta)$  evaluated at  $\hat{\theta}$  is very close to 0.

Suppose  $u_k, k = 1, \dots, m$  are sampled from importance sampling distribution  $\tilde{f}(u_k)$ . Recall that the gradient of the MCLA is:

$$\nabla l_m(\theta) = \frac{\sum_{k=1}^m \frac{\nabla f_\theta(y|u_k) f_\theta(u_k)}{\tilde{f}(u_k)}}{\sum_{k=1}^m \frac{f_\theta(u_k)}{\tilde{f}(u_k)}} \quad (122)$$

We can choose  $\tilde{f}(u_k) = f_{\hat{\theta}}(u_k, y)$ . If the putative MLE is truly the MLE, then  $\theta = \hat{\theta}$ , which in turn implies that  $\tilde{f}(u_k) = f_\theta(u_k, y)$ . To be clear, in this check we no longer use a mixture of normals for  $\tilde{f}(u_k)$ . With the selected importance sampling distribution, each of the importance sampling weights is equal to 1 and the sum of the weights is  $m$ . Therefore, the gradient of the MCLA simplifies as follows:

$$\nabla l_m(\theta) = \frac{1}{m} \sum_{k=1}^m \frac{\nabla f_\theta(y|u_k) f_\theta(u_k)}{\tilde{f}(u_k, y)} \quad (123)$$

$$= \frac{1}{m} \sum_{k=1}^m \frac{\nabla f_\theta(u_k, y)}{\tilde{f}(u_k, y)} \quad (124)$$

$$= \frac{1}{m} \sum_{k=1}^m \nabla \log f_\theta(u_k, y) \quad (125)$$

This shows that the gradient of the MCLA is the average of the gradient of the complete data log likelihood, as long as  $\hat{\theta}$  is truly an MLE. We can produce  $u_k, k = 1, \dots, m$ , using Markov chain Monte Carlo. Using  $u$  as the variable, we run the Markov chain (perhaps **metrop** from the R package **mcmc**) on the complete data log likelihood. We then use these samples to calculate the gradient of the complete data log likelihood, which in turn calculates the gradient of the MCLA.

If we split the MCMC runs into batches, we can calculate the batch means of the MCLA gradient, the grand mean of the MCLA gradient, and the corresponding Monte Carlo standard error. If the putative MLE truly maximizes the likelihood, then the MCLA gradient's components should be close to 0. We can check that they are close enough to 0 by comparing them to the Monte Carlo standard error.

## 10 Likelihood ratio test

This is another S3 generic function that the user can choose to implement if they'd like to do likelihood ratio tests for nested models. Eventually I'll want this to handle an arbitrary number of models, but for now we'll stick to comparing two models: *mod2* nested in *mod1*. I will assume these models have already been fit.

There are a few ways that the two models could differ. In other words, we could be testing:

1. whether one or more fixed effects are zero.
2. whether one variance component is zero.
3. whether multiple variance components are zero.
4. whether one or more fixed effects is zero and one or more variance components are zero.

The hypothesis testing procedure and method for calculating pvalues are different for each of these cases, so I'll go through each one in the following subsubsections.

### 10.1 Testing whether one or more fixed effects are zero

Consider the case of testing whether one or more fixed effects are zero and the random effects are the same between the two models. I think we just do a likelihood ratio test as we're familiar with for linear models. Let the number of fixed effect parameters in the larger model be  $p_1$  and the number in the nested model be  $p_2$ . I will need the log likelihood values of the models outputted (from the main function at the same time as when I output the MCMLEs and the Fisher information). The log likelihood is simply the "value" from the objective function that trust uses. Let's call the log likelihood from the larger model  $l_1$  and the log likelihood from the nested model  $l_2$ .

Then the likelihood ratio test statistic is

$$t_{LRT} = -2l_1 + 2l_2 \quad (126)$$

and this follows a  $\chi^2$  distribution with  $p_1 - p_2$  degrees of freedom. The calculation

$$P(t_{\chi^2_{p_1-p_2}} > t_{LRT}). \quad (127)$$

gives us a two-sided pvalue to fit with the two-sided alternative hypothesis.

### 10.2 Testing whether one variance component is zero

Hypothesis testing for one variance component is easy but not what most people expect. In this setup, the two models have the exact same fixed effects. The only difference is the larger model has one more random effect  $\nu_{t_1}$ . This means we want to test whether  $\nu_{t_1} = 0$ . A variance component must be nonnegative, meaning the alternative hypothesis is that  $\nu_{t_1} > 0$ . In other words, the hypotheses are:

$$H_0 : \nu_{t_1} = 0 \quad (128)$$

$$H_1 : \nu_{t_1} > 0. \quad (129)$$



Note the alternative hypothesis is one-sided. This means we need to calculate a one-sided pvalue. Let  $t_{LRT} = 2l_2 - 2l_1$ . If we naively follow the pvalue calculation of

$$P(\chi_1^2 > t_{LRT}) \quad (130)$$

we will end up calculating a two-sided pvalue. To fix this, we cut this pvalue in half. In other words, use the standard normal distribution  $Z$  and think of calculating the one-sided pvalue this way:

$$P(Z > \sqrt{|t_{LRT}|}) \quad (131)$$

However, this is not obvious; I never would have thought about having a test statistic of

$$\sqrt{2|l_2 - l_1|}. \quad (132)$$

### 10.3 Testing whether multiple variance components are zero

In this case, the fixed effects are identical between the larger and nested model. The only difference is that the larger model has more than one additional variance components. This gets a lot more complicated (for example, there is no clear way to count parameters so calculating the degrees of freedom is out the window). Luckily, someone's worked out the details already. Charlie suggested a few papers to look at. I'll read about and add this a bit later. Charlie is hoping that we'll have a mixture of  $\chi^2$  distributions, such as

$$\frac{1}{2}P(\chi_1^2 > t_{LRT}) + \frac{1}{4}P(\chi_2^2 > t_{LRT}). \quad (133)$$

The reason this is complicated is because the variance components are restricted to be nonnegative. This means that the likelihood is only defined for nonnegative variance components. Then differentiating the likelihood at 0 becomes tricky because that's the boundary.

### 10.4 Testing whether one or more fixed effects is zero and one or more variance components are zero

This seems very complicated. I don't know if it's any more complicated than the previous case of testing multiple variance components. I'll have to think about it later when I have time.

### 10.5 Determining the hypothesis test

To follow convention, this command will be called "anova" and the two arguments will be the two models to be compared. The command would look like

```
anova(mod1,mod2)
```

R will first need to figure out if the fixed effects are different and, if they are, which model is the larger model.

```
if(length(coef(mod1))>length(coef(mod2)))
  {bigmod<-mod1; smallmod <-mod2}
if(length(coef(mod1))<length(coef(mod2)))
  {bigmod<-mod2; smallmod <-mod1}
```

Next, if there is a difference in the fixed effects, I'm going to make sure the fixed effects of the small model are nested in those of the big model. Otherwise, produce an error. Then, continue with the testing. Note: this check for nesting isn't perfect, but no other anova checks which model is bigger and whether they're nested.

```
if(big mod is defined) {
  pnames1<-names(coef(bigmod))
  pnames2<-names(coef(smallmod))
  if(sum(pnames2 %in% pnames1) != length(pnames2)) {
    stop("The models you provided are not nested.")}

  if(the variance components differ between the two models){
    check big model has more variance components (ow: error)
    check variance components of the small model are nested (ow: error)
    calculate pvalue according to section 3.3.4. return it.
  }
  if(the variance components do not differ between the two models){
    calculate pvalue according to 3.3.1. return it.
  }
}
```

If and only if we've gotten to the end of this chunk of code without returning anything, bigmod has not been defined, meaning the fixed effects are same. Therefore, we now need to figure out whether the variance components differ by one or by more than one. In order to compare the number of random effects, we need to figure out how to get the number of variance components from each model. I don't yet have a solid grasp of how formula works its magic, but I have the impression that I can count the number of model matrices returned in order to figure out the number of variance components because there is a model matrix for each one.

```

T1<- number of variance components for mod1
T2<- number of variance components for mod2
if(T1>T2){call mod1 the bigmod and mod2 the smallmod}
if(T1<T2){call mod2 the bigmod and mod1 the smallmod}
if(bigmod is defined){
  check small model is nested in the big model (ow: error)
  (maybe do this by looking at the call?)
  if(T1==T2+1 || T2==T1+1)
    {calculate and return pvalue according to sec 3.3.2}
  calculate and return pvalue according to section 3.3.3
}
if still haven't returned anything, produce error.

```

Why the last error? If we have gotten to this point in the code with nothing returned, it ends up that the number of variance components in each model are equal, meaning the user has made some kind of mistake. Either the user provided two identical models or they provided non-nested models (e.g. they have the same number of variance components, but different components in each model). Either way, we can't help them.

This command would return something that first reminds the user what predictors are in each model, then has a table. Each model would have one row of the table. The columns would contain the model name, the log likelihood for each model, the number of parameters in the model. Then the next columns would have one entry each: the calculated difference between the log likelihoods, the calculated difference between the number of parameters, and the pvalue of the test.

## 11 Confidence intervals

The user can implement this command to create confidence intervals after fitting a model *mod* using the main function. This is an S3 generic. The command would look like

```
confint(mod,parm,level=.95)
```

The only required argument would be the fitted model (the first argument). The third argument (the confidence level) has a default of .95.

If the second argument is omitted, confidence intervals would be created for all of the parameters. There are two options to calculate confidence intervals for a subset of the parameters. The user can provide either the names of the coefficients or a vector with length equal to the number of

parameters (entries being 1 if they would like that intervals for that parameter and 0 otherwise). The code to do this is taken from “confint.lm” and is:

```
function (object, parm, level = 0.95, ...)
{
  cf <- coef(object)
  pnames <- names(cf)
  if (missing(parm))
    parm <- pnames
  else if (is.numeric(parm))
    parm <- pnames[parm]
    stopifnot(parm %in% pnames)

  rest of the code to actually do the confidence intervals goes here
}
```

What this says: write down the coefficient names of the object we’re given. If `parm` is missing, we’re going to assume they want to create intervals for all the parameters. If `parm` is numeric (a vector of 0s and 1s), use that to select the parameter names we want and call that selection `parm`. Finally, stop the whole process if `parm` (the parameter names we want intervals for) don’t appear in the model.

The form for a confidence interval is the point estimate plus or minus the standard error of the point estimate times some cutoff. The point estimate and the standard error are from the summary of the model. The reference distribution used for the cutoff is the standard normal. This will produce an asymptotic confidence interval.

The output would either be a vector of length 2 (if creating confidence interval for only one parameter) or a matrix with 2 columns (one for the lower bound of the confidence interval, one for the upper bound). The column names will be “ $(100 - level)/2$ ” and “ $50 + level/2$ .” (in the default case, this is “2.5%” and “97.5%”).

The only potential hangup is when creating confidence intervals for the variance components  $\nu_t, t = 1, \dots, T$ . We know that  $\nu_t > 0$  for all  $t = 1, \dots, T$ . We’ll find this boundary again makes things complicated. I don’t yet know how to deal with this.

To illustrate the problem, consider the scenario that the margin of error for  $\nu_t$  is greater than the estimate of  $\nu_t$  itself. It wouldn’t make sense to produce a confidence interval with a negative lower bound. It could make sense to truncate the interval so that the lower bound starts at 0 and the upper bound remains untouched.

I thought for three or four minutes about a one sided confidence interval, but I don't think that captures what we want. We want the range of all likely values for the variance component, not just an upper bound.

## References

- Booth, J. G. and Hobert, J. P. (1999). Maximizing generalized linear mixed model likelihoods with an automated Monte Carlo EM algorithm. *Journal of the Royal Statistical Society, Series B*, 61:265–285.
- Coull, B. and Agresti, A. (2000). Random effects modeling of multiple binomial responses using the multivariate binomial logit-normal distribution. *biometrics*, 56:73–80.
- Geyer, C. J. (2014). R package `aster` (aster models), version .8-30. <http://cran.r-project.org/package=aster>.
- Geyer, C. J. and Thompson, E. (1992). Constrained Monte Carlo maximum likelihood for dependent data. *Journal of the Royal Statistical Society, Series B*, 54:657–699.

## A MCLA calculations

Let  $u_k, k = 1, \dots, m$  be a sample from  $\tilde{f}(u_k)$ . The Monte Carlo log likelihood approximation is

$$l_m(\theta) = \log \frac{1}{m} \sum_{k=1}^m f_\theta(y|u_k) \frac{f_\theta(u_k)}{\tilde{f}(u_k)} \quad (134)$$

$$= \log \frac{1}{m} \sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}. \quad (135)$$

The gradient vector and Hessian matrix calculations depend on whether  $\tilde{f}(u_k)$  contains  $\theta$ . The calculations in A.2 are for  $\tilde{f}(u_k)$  containing  $\theta$  and the calculations in A.1 are for  $\tilde{f}(u_k)$  not containing  $\theta$ .

### A.1 MCLA derivatives when $\tilde{f}$ independent of $\theta$

When  $\tilde{f}$  is independent of  $\theta$ , the gradient vector of the MCLA with respect to  $\theta$  is

$$\nabla l_m(\theta) = \frac{\sum_{k=1}^m \left( \nabla \log f_\theta(u_k, y) \frac{f_\theta(u_k, y)}{\tilde{f}(u_k)} \right)}{\sum_{k=1}^m \left( \frac{f_\theta(u_k, y)}{\tilde{f}(u_k)} \right)} \quad (136)$$

and the Hessian matrix of the MCLA is

$$\nabla^2 l_m(\theta) = \frac{\sum_{k=1}^m \left( \nabla^2 \log f_\theta(u_k, y) \frac{f_\theta(u_k, y)}{\tilde{f}(u_k)} \right)}{\sum_{k=1}^m \left( \frac{f_\theta(u_k, y)}{\tilde{f}(u_k)} \right)} - \nabla l_m(\theta) (\nabla l_m(\theta))' \quad (137)$$

$$+ \frac{\sum_{k=1}^m \left( \nabla \log f_\theta(u_k, y) (\nabla \log f_\theta(u_k, y))' \frac{f_\theta(u_k, y)}{\tilde{f}(u_k)} \right)}{\sum_{k=1}^m \left( \frac{f_\theta(u_k, y)}{\tilde{f}(u_k)} \right)}. \quad (138)$$

To reduce the risk of catastrophic cancellation, we can combine the last two terms of the Hessian:

$$\nabla^2 l_m(\theta) = \frac{\sum_{k=1}^m \nabla^2 \log f_\theta(y, u_k) \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}}{\sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}} \quad (139)$$

$$+ \frac{\sum_{k=1}^m [\nabla \log f_\theta(y, u_k) - \nabla l_m(\theta)] [\nabla \log f_\theta(y, u_k) - \nabla l_m(\theta)]' \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}}{\sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}} \quad (140)$$

We are able to combine the last two terms because  $\nabla l_m(\theta)$  is a weighted mean of  $\nabla \log f_\theta(y, u_k) - \nabla \log \tilde{f}(u_k)$ . Letting

$$Z = \nabla \log f_\theta(y, u_k) \quad (141)$$

we can use the following equality:

$$E(ZZ') - E(Z)E(Z)' = [E(Z - EZ)] [E(Z - EZ)]'. \quad (142)$$

## A.2 MCLA derivatives when $\tilde{f}$ depends on $\theta$

When  $\tilde{f}$  contains  $\theta$ , the gradient of the MCLA is

$$\nabla l_m(\theta) = \frac{\nabla \left[ \sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)} \right]}{\sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}} \quad (143)$$

$$\begin{aligned} & \sum_{k=1}^m \frac{\nabla f_\theta(y, u_k)}{\tilde{f}(u_k)} - \frac{f_\theta(y, u_k) \nabla \tilde{f}(u_k)}{(\tilde{f}(u_k))^2} \\ &= \frac{\sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}}{\sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}} \end{aligned} \quad (144)$$

$$\begin{aligned} & \sum_{k=1}^m \frac{\nabla f_\theta(y, u_k)}{f_\theta(y, u_k)} \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)} - \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)} \frac{\nabla \tilde{f}(u_k)}{\tilde{f}(u_k)} \\ &= \frac{\sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}}{\sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}} \end{aligned} \quad (145)$$

$$\begin{aligned} & \sum_{k=1}^m \nabla \log f_\theta(y, u_k) \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)} - \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)} \nabla \log \tilde{f}(u_k) \\ &= \frac{\sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}}{\sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}} \end{aligned} \quad (146)$$

$$\begin{aligned} & \sum_{k=1}^m \left[ \nabla \log f_\theta(y, u_k) - \nabla \log \tilde{f}(u_k) \right] \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)} \\ &= \frac{\sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}}{\sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}} \end{aligned} \quad (147)$$

Then the Hessian of the MCLA is

$$\nabla^2 l_m(\theta) = \nabla \left[ \frac{\sum_{k=1}^m \left[ \nabla \log f_\theta(y, u_k) - \nabla \log \tilde{f}(u_k) \right] \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}}{\sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}} \right] \quad (148)$$

$$= \frac{\sum_{k=1}^m \left[ \nabla^2 \log f_\theta(y, u_k) - \nabla^2 \log \tilde{f}(u_k) \right] \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}}{\sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}} \quad (149)$$

$$+ \frac{\sum_{k=1}^m \left[ \nabla \log f_\theta(y, u_k) - \nabla \log \tilde{f}(u_k) \right] \nabla \left[ \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)} \right]}{\sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}} \quad (150)$$

$$- \left[ \frac{\sum_{k=1}^m \left[ \nabla \log f_\theta(y, u_k) - \nabla \log \tilde{f}(u_k) \right] \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}}{\left( \sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)} \right)^2} \right] \nabla \left[ \sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)} \right] \quad (151)$$

$$= \frac{\sum_{k=1}^m \left[ \nabla^2 \log f_\theta(y, u_k) - \nabla^2 \log \tilde{f}(u_k) \right] \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}}{\sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}} \quad (152)$$

$$+ \frac{\sum_{k=1}^m \left[ \nabla \log f_\theta(y, u_k) - \nabla \log \tilde{f}(u_k) \right] \left[ \nabla \log f_\theta(y, u_k) - \nabla \log \tilde{f}(u_k) \right]' \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}}{\sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}} \quad (153)$$

$$- \left[ \frac{\sum_{k=1}^m \left[ \nabla \log f_\theta(y, u_k) - \nabla \log \tilde{f}(u_k) \right] \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}}{\left( \sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)} \right)^2} \right] \left[ \nabla \log f_\theta(y, u_k) - \nabla \log \tilde{f}(u_k) \right]' \quad (154)$$

$$= \frac{\sum_{k=1}^m \left[ \nabla^2 \log f_\theta(y, u_k) - \nabla^2 \log \tilde{f}(u_k) \right] \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}}{\sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}} \quad (155)$$

$$+ \frac{\sum_{k=1}^m \left[ \nabla \log f_\theta(y, u_k) - \nabla \log \tilde{f}(u_k) \right] \left[ \nabla \log f_\theta(y, u_k) - \nabla \log \tilde{f}(u_k) \right]' \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}}{\sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}} \quad (156)$$

$$- [\nabla l_m(\theta)] [\nabla l_m(\theta)]' \quad (157)$$



To reduce the risk of catastrophic cancellation, we can combine the last two terms of the Hessian:

$$\nabla^2 l_m(\theta) = \frac{\sum_{k=1}^m \left[ \nabla^2 \log f_\theta(y, u_k) - \nabla^2 \log \tilde{f}(u_k) \right] \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}}{\sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}} \quad (158)$$

$$+ \frac{\sum_{k=1}^m \left[ \nabla \log f_\theta(y, u_k) - \nabla \log \tilde{f}(u_k) - \nabla l_m(\theta) \right] \left[ \nabla \log f_\theta(y, u_k) - \nabla \log \tilde{f}(u_k) - \nabla l_m(\theta) \right]' \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}}{\sum_{k=1}^m \frac{f_\theta(y, u_k)}{\tilde{f}(u_k)}} \quad (159)$$

We are able to combine the last two terms because  $\nabla l_m(\theta)$  is a weighted mean of  $\nabla \log f_\theta(y, u_k) - \nabla \log \tilde{f}(u_k)$ . Letting

$$Z = \nabla \log f_\theta(y, u_k) - \nabla \log \tilde{f}(u_k) \quad (160)$$

we can use the following equality:

$$E(ZZ') - E(Z)E(Z)' = [E(Z - EZ)] [E(Z - EZ)]'. \quad (161)$$

## B Central Limit Theorem for MCLA

In this section, we focus on the gradient of the MCLA because the trust criterion for finding the maximum is based on the gradient. Define

$$\gamma_1 = \int f_\theta(u, y) du \quad (162)$$

$$\gamma_2 = \int \tilde{f}(u) du. \quad (163)$$

Recall the calculation for the MCLA gradient originally stated in (4):

$$\nabla l_m(\theta) = \frac{\sum_{k=1}^m \nabla \log f_\theta(u_k, y) \frac{f_\theta(u_k, y)}{\tilde{f}(u_k)}}{\sum_{k=1}^m \left( \frac{f_\theta(u_k, y)}{\tilde{f}(u_k)} \right)} \quad (164)$$

Note that both the numerator and denominator are sample means. By the law of large numbers,

the numerators and denominator each converge to their true means. That is,

$$\frac{1}{m} \sum_{k=1}^m \nabla \log f_{\theta}(u_k, y) \frac{f_{\theta}(u_k, y)}{\tilde{f}(u_k)} \rightarrow E_{\tilde{f}} \left[ \nabla \log f_{\theta}(u, y) \frac{f_{\theta}(u, y)}{\tilde{f}(u)} \right] \quad (165)$$

$$= \int \nabla \log f_{\theta}(u, y) \frac{f_{\theta}(u, y)}{\tilde{f}(u)} \frac{\tilde{f}(u)}{\gamma_2} du \quad (166)$$

$$= \frac{\gamma_1}{\gamma_2} \int \nabla \log f_{\theta}(u, y) \frac{f_{\theta}(u, y)}{\gamma_1} du \quad (167)$$

$$= \frac{\gamma_1}{\gamma_2} E_f [\nabla \log f_{\theta}(u, y)] \quad (168)$$

and

$$\frac{1}{m} \sum_{k=1}^m \frac{f_{\theta}(u_k, y)}{\tilde{f}(u_k)} \rightarrow E_{\tilde{f}} \left[ \frac{f_{\theta}(u, y)}{\tilde{f}(u)} \right] \quad (169)$$

$$= \int \frac{f_{\theta}(u, y)}{\tilde{f}(u)} \frac{\tilde{f}(u)}{\gamma_2} du \quad (170)$$

$$= \frac{\gamma_1}{\gamma_2} \int \frac{f_{\theta}(u, y)}{\gamma_2} du \quad (171)$$

$$= \frac{\gamma_1}{\gamma_2} \quad (172)$$

Then, by Slutsky's theorem,

$$\nabla l_m(\theta) \rightarrow \frac{\frac{\gamma_1}{\gamma_2} E_f [\nabla \log f_{\theta}(u, y)]}{\frac{\gamma_1}{\gamma_2}} \quad (173)$$

$$= E_f [\nabla \log f_{\theta}(u, y)] \quad (174)$$

$$= \nabla l(\theta) \quad (175)$$

In addition to a law of large numbers, we would like a Central Limit Theorem for  $\nabla l_m(\theta)$ . In other

words, the quantity

$$\sqrt{m} \left[ \frac{\frac{1}{m} \sum_{k=1}^m \nabla \log f_{\theta}(u_k, y) \frac{f_{\theta}(u_k, y)}{\tilde{f}(u_k)}}{\frac{1}{m} \sum_{k=1}^m \frac{f_{\theta}(u_k, y)}{\tilde{f}(u_k)}} - l(\theta) \right] \quad (176)$$

$$= \frac{\frac{1}{\sqrt{m}} \sum_{k=1}^m \nabla \log f_{\theta}(u_k, y) \frac{f_{\theta}(u_k, y)}{\tilde{f}(u_k)}}{\frac{1}{m} \sum_{k=1}^m \frac{f_{\theta}(u_k, y)}{\tilde{f}(u_k)}} - \sqrt{m} l(\theta) \quad (177)$$

$$= \frac{\frac{1}{\sqrt{m}} \sum_{k=1}^m \nabla \log f_{\theta}(u_k, y) \frac{f_{\theta}(u_k, y)}{\tilde{f}(u_k)} - l(\theta) \frac{1}{\sqrt{m}} \sum_{k=1}^m \frac{f_{\theta}(u_k, y)}{\tilde{f}(u_k)}}{\frac{1}{m} \sum_{k=1}^m \frac{f_{\theta}(u_k, y)}{\tilde{f}(u_k)}} \quad (178)$$

will have a normal distribution if and only if the variances of

$$\sum_{k=1}^m \frac{f_{\theta}(u_k, y)}{\tilde{f}(u_k)} \quad (179)$$

and

$$\sum_{k=1}^m \nabla \log f_{\theta}(u_k, y) \frac{f_{\theta}(u_k, y)}{\tilde{f}(u_k)} \quad (180)$$

are finite. First, we want to show that

$$\sum_{k=1}^m \frac{f_{\theta}(u_k, y)}{\tilde{f}(u_k)} \quad (181)$$

has finite variance. This is true if and only if

$$\int \frac{f_{\theta}(u, y)^2}{\tilde{f}(u)} du < \infty. \quad (182)$$

We see:

$$\int \frac{f_{\theta}(u, y)^2}{\tilde{f}(u)} du = \int \frac{f_{\theta}(y|u)^2 f_{\theta}(u)^2}{\tilde{f}(u)} du \quad (183)$$

$$\leq \int \frac{f_{\theta}(y|u)^2 f_{\theta}(u)^2}{p_1 \tilde{f}(u)} du \quad (184)$$

$$\propto \frac{1}{p_1} \int \frac{\exp(-u' D^{-1} u) f_{\theta}(y|u)^2}{[1 + u' u / \zeta]^{-(\zeta+q)/2}} du \quad (185)$$

Since

$$f_{\theta}(y|u)^2 = \left( e^{y' \eta - c(\eta)} \right)^2 = e^{2y' \eta - 2c(\eta)}, \quad (186)$$

$$\int \frac{f_\theta(u, y)^2}{\tilde{f}(u)} du \propto \frac{1}{p_1} \int \frac{\exp(-u' D^{-1} u) \exp(2y' \eta)}{[1 + u' u / \zeta]^{-(\zeta+q)/2} \exp(2c(\eta))} du \quad (187)$$

$$= \frac{1}{p_1} \int \frac{[1 + u' u / \zeta]^{(\zeta+q)/2} \exp(2y' \eta)}{\exp(u' D^{-1} u) \exp(2c(\eta))} du \quad (188)$$

$$= \frac{1}{p_1} \int \frac{[1 + u' u / \zeta]^{(\zeta+q)/2} \exp(2y' Z u)}{\exp(u' D^{-1} u) \exp(2c(\eta))} du \quad (189)$$

When  $y|u \sim \text{Bernoulli}$ , then

$$c(\eta) = \log(1 + e^\eta) \Rightarrow \exp(2c(\eta)) = \exp(2 \log(1 + e^\eta)) \quad (190)$$

$$= (\exp(\log(1 + e^\eta)))^2 \quad (191)$$

$$= (1 + e^\eta)^2 \quad (192)$$

and

$$\int \frac{f_\theta(u, y)^2}{\tilde{f}(u)} du \propto \frac{1}{p_1} \int \frac{[1 + u' u / \zeta]^{(\zeta+q)/2} \exp(2y' Z u)}{\exp(u' D^{-1} u) (1 + e^\eta)^2} du, \quad (193)$$

which converges since  $\exp(u' D^{-1} u)$  grows faster than any term in the numerator of the integrand.

On the other hand, when  $y|u \sim \text{Poisson}$ , then

$$c(\eta) = \exp(\eta) \Rightarrow \exp(2c(\eta)) = (\exp(\exp(\eta)))^2 \quad (194)$$

and

$$\int \frac{f_\theta(u, y)^2}{\tilde{f}(u)} du \propto \frac{1}{p_1} \int \frac{[1 + u' u / \zeta]^{(\zeta+q)/2} \exp(2y' Z u)}{\exp(u' D^{-1} u) (\exp(\exp(\eta)))^2} du, \quad (195)$$

which converges since  $(\exp(\exp(\eta)))^2$  grows faster than any term in the numerator of the integrand.

Therefore, whether the data are Bernoulli or Poisson distributed, the importance sampling weights

$$\sum_{k=1}^m \frac{f_\theta(u_k, y)}{\tilde{f}(u_k)} \quad (196)$$

have finite variance.

Now, we want to show that

$$\sum_{k=1}^m \nabla \log f_\theta(u_k, y) \frac{f_\theta(u_k, y)}{\tilde{f}(u_k)} = \sum_{k=1}^m \nabla \log f_\theta(y|u_k) \frac{f_\theta(u_k, y)}{\tilde{f}(u_k)} + \sum_{k=1}^m \nabla \log f_\theta(u_k) \frac{f_\theta(u_k, y)}{\tilde{f}(u_k)} \quad (197)$$

has finite variance. We need to show the following are finite:

$$\text{Var} \left( \sum_{k=1}^m \nabla \log f_{\theta}(y|u_k) \frac{f_{\theta}(u_k, y)}{\tilde{f}(u_k)} \right) \quad (198)$$

$$\text{Var} \left( \sum_{k=1}^m \nabla \log f_{\theta}(u_k) \frac{f_{\theta}(u_k, y)}{\tilde{f}(u_k)} \right) \quad (199)$$

$$\text{Cov} \left( \sum_{k=1}^m \nabla \log f_{\theta}(y|u_k) \frac{f_{\theta}(u_k, y)}{\tilde{f}(u_k)}, \sum_{k=1}^m \nabla \log f_{\theta}(u_k) \frac{f_{\theta}(u_k, y)}{\tilde{f}(u_k)} \right). \quad (200)$$

That is, we want to show the existence of the following:

$$\int \frac{[\nabla \log f_{\theta}(y|u)]^2 [f_{\theta}(u, y)]^2}{\tilde{f}(u)} du \quad (201)$$

$$\int \frac{[\nabla \log f_{\theta}(u)]^2 [f_{\theta}(u, y)]^2}{\tilde{f}(u)} du \quad (202)$$

$$\int \frac{\nabla \log f_{\theta}(u) \nabla \log f_{\theta}(y|u) [f_{\theta}(u, y)]^2}{\tilde{f}(u)} du \quad (203)$$

First, we start with (201):

$$\int \frac{[\nabla \log f_{\theta}(y|u)]^2 [f_{\theta}(u, y)]^2}{\tilde{f}(u)} du = \int \frac{[\nabla \log f_{\theta}(y|u)]^2 [f_{\theta}(y|u)]^2 [f_{\theta}(u)]^2}{\tilde{f}(u)} du \quad (204)$$

$$\leq \int \frac{[\nabla \log f_{\theta}(y|u)]^2 [f_{\theta}(y|u)]^2 [f_{\theta}(u)]^2}{p_1 \tilde{f}(u)} du \quad (205)$$

$$\propto \int \frac{[\nabla \log f_{\theta}(y|u)]^2 [f_{\theta}(y|u)]^2 \exp(-u' D^{-1} u)}{p_1 [1 + u' u / \zeta]^{-(\zeta+q)/2}} du \quad (206)$$

$$= \int \frac{[\nabla \log f_{\theta}(y|u)]^2 [f_{\theta}(y|u)]^2 [1 + u' u / \zeta]^{(\zeta+q)/2}}{p_1 \exp(u' D^{-1} u)} du \quad (207)$$

$$= \int \frac{[\nabla \log f_{\theta}(y|u)]^2 \exp(2y'(X\beta + Zu)) [1 + u' u / \zeta]^{(\zeta+q)/2}}{p_1 \exp(2c(\eta)) \exp(u' D^{-1} u)} du \quad (208)$$

By the Cauchy-Schwartz inequality, the integral in (208) exists if

$$\begin{aligned} & \int \frac{\left[ \frac{\partial}{\partial \beta} \log f_{\theta}(y|u) \right]^2 \exp(2y'(X\beta + Zu)) [1 + u' u / \zeta]^{(\zeta+q)/2}}{p_1 \exp(2c(\eta)) \exp(u' D^{-1} u)} du \\ & \propto \int \frac{[X' c'(\eta)]^2 \exp(2y'(X\beta + Zu)) [1 + u' u / \zeta]^{(\zeta+q)/2}}{p_1 \exp(2c(\eta)) \exp(u' D^{-1} u)} du \end{aligned} \quad (209)$$

exists, where  $c'(\eta)$  is a vector with components  $c'(\eta_i)$ . This will converge as long as no term in the numerator of the integrand grows more quickly than  $\exp(u'D^{-1}u)$ . The second and third terms grow less quickly than  $\exp(u'D^{-1}u)$ , so  $[c'(\eta)]^2$  is the only term to check. If  $y|u \sim \text{Bernoulli}$ , then

$$c'(\eta_i) = \frac{\exp(\eta_i)}{1 + \exp(\eta_i)}, \quad (210)$$

which does not grow more quickly than  $\exp(u'D^{-1}u)$ . If  $y|u \sim \text{Poisson}$ , then

$$c'(\eta_i) = e^{\eta_i}, \quad (211)$$

which does not grow more quickly than  $\exp(u'D^{-1}u)$ . Therefore, (201) exists and is finite, and

$$\text{Var} \left( \sum_{k=1}^m \nabla \log f_\theta(y|u_k) \frac{f_\theta(u_k, y)}{\tilde{f}(u_k)} \right) < \infty. \quad (212)$$

is finite.

Next, we move on to showing the existence of (202).

$$\int \frac{[\nabla \log f_\theta(u)]^2 [f_\theta(u, y)]^2}{\tilde{f}(u)} du = \int \frac{[\nabla \log f_\theta(u)]^2 [f_\theta(y|u)]^2 [f_\theta(u)]^2}{\tilde{f}(u)} du \quad (213)$$

$$\leq \int \frac{[\nabla \log f_\theta(u)]^2 [f_\theta(y|u)]^2 [f_\theta(u)]^2}{p_1 \dot{f}(u)} du \quad (214)$$

By the Cauchy-Schwartz inequality, the above integral exists as long as

$$\int \frac{\left[ \frac{\partial}{\partial \nu_t} \log f_\theta(u) \right]^2 [f_\theta(y|u)]^2 [f_\theta(u)]^2}{p_1 \dot{f}(u)} du < \infty \quad (215)$$

for every  $t = 1, \dots, T$ . We see

$$\int \frac{\left[ \frac{\partial}{\partial \nu_t} \log f_\theta(u) \right]^2 [f_\theta(y|u)]^2 [f_\theta(u)]^2}{p_1 \dot{f}(u)} du \quad (216)$$

$$= \int \frac{\left[ \frac{\partial}{\partial \nu_t} \log f_\theta(u) \right]^2 \exp(2y'(X\beta + Zu)) [f_\theta(u)]^2}{p_1 \exp(2c(\eta)) \dot{f}(u)} du \quad (217)$$

$$\propto \int \frac{\left[ \frac{\partial}{\partial \nu_t} \log f_\theta(u) \right]^2 \exp(2y'(X\beta + Zu)) [1 + u'u/\zeta]^{(\zeta+q)/2}}{p_1 \exp(2c(\eta)) \exp(u'D^{-1}u)} du \quad (218)$$

$$= \int \frac{\left[ -\frac{q_t}{2\nu_t} + \frac{u'_t u_t}{2\nu_t^2} \right]^2 \exp(2y'(X\beta + Zu)) [1 + u'u/\zeta]^{(\zeta+q)/2}}{p_1 \exp(2c(\eta)) \exp(u'D^{-1}u)} du \quad (219)$$

Because every term in the numerator of the integrand grows less quickly than  $\exp(u'D^{-1}u)$ , the integral exists. Therefore, (202) exists and

$$\text{Var} \left( \sum_{k=1}^m \nabla \log f_{\theta}(u_k) \frac{f_{\theta}(u_k, y)}{\tilde{f}(u_k)} \right) < \infty. \quad (220)$$

Lastly, we want to show the existence of (203). We see

$$\int \frac{\nabla \log f_{\theta}(u) \nabla \log f_{\theta}(y|u) [f_{\theta}(u, y)]^2}{\tilde{f}(u)} du = \int \frac{\nabla \log f_{\theta}(u) \nabla \log f_{\theta}(y|u) [f_{\theta}(y|u)]^2 [f_{\theta}(u)]^2}{\tilde{f}(u)} du \quad (221)$$

$$\leq \int \frac{\nabla \log f_{\theta}(u) \nabla \log f_{\theta}(y|u) [f_{\theta}(y|u)]^2 [f_{\theta}(u)]^2}{p_1 \hat{f}(u)} du \quad (222)$$

Again, we use the Cauchy-Schwartz inequality to check the existence of (223). If, for every  $t = 1, \dots, T$ ,

$$\int \frac{\left[ \frac{\partial}{\partial \nu_t} \log f_{\theta}(u) \right] \left[ \frac{\partial}{\partial \beta} \log f_{\theta}(y|u) \right] [f_{\theta}(y|u)]^2 [f_{\theta}(u)]^2}{p_1 \hat{f}(u)} du \quad (223)$$

exists, then (223) also exists. Continuing, we see

$$\int \frac{\left[ \frac{\partial}{\partial \nu_t} \log f_{\theta}(u) \right] \left[ \frac{\partial}{\partial \beta} \log f_{\theta}(y|u) \right] [f_{\theta}(y|u)]^2 [f_{\theta}(u)]^2}{p_1 \hat{f}(u)} du \quad (224)$$

$$\propto \int \frac{\left[ \frac{\partial}{\partial \nu_t} \log f_{\theta}(u) \right] \left[ \frac{\partial}{\partial \beta} \log f_{\theta}(y|u) \right] [f_{\theta}(y|u)]^2}{p_1 \exp(u'D^{-1}u) \hat{f}(u)} du \quad (225)$$

$$\propto \int \frac{\left[ \frac{\partial}{\partial \nu_t} \log f_{\theta}(u) \right] \left[ \frac{\partial}{\partial \beta} \log f_{\theta}(y|u) \right] [f_{\theta}(y|u)]^2 [1 + u'u/\zeta]^{(\zeta+q)/2}}{p_1 \exp(u'D^{-1}u)} du \quad (226)$$

$$\propto \int \frac{\left[ -\frac{q_t}{2\nu_t} + \frac{u'_t u_t}{2\nu_t^2} \right] [X'(y - c'(\eta))] [f_{\theta}(y|u)]^2 [1 + u'u/\zeta]^{(\zeta+q)/2}}{p_1 \exp(u'D^{-1}u)} du \quad (227)$$

$$= \int \frac{\left[ -\frac{q_t}{2\nu_t} + \frac{u'_t u_t}{2\nu_t^2} \right] [X'(y - c'(\eta))] \exp(2y'(X\beta + Zu)) [1 + u'u/\zeta]^{(\zeta+q)/2}}{p_1 \exp(u'D^{-1}u) \exp(2c(\eta))} du. \quad (228)$$

We have already shown that  $-c'(\eta)$  does not grow more quickly than  $\exp(u'D^{-1}u)$ . The other terms in the numerator of the integrand shown in (228) also grow more slowly than  $\exp(u'D^{-1}u)$ .

We conclude that (228) exists, which implies that (203) exists and

$$\text{Cov} \left( \sum_{k=1}^m \nabla \log f_{\theta}(y|u_k) \frac{f_{\theta}(u_k, y)}{\tilde{f}(u_k)}, \sum_{k=1}^m \nabla \log f_{\theta}(u_k) \frac{f_{\theta}(u_k, y)}{\tilde{f}(u_k)} \right) < \infty. \quad (229)$$

This shows that

$$\text{Var} \left( \sum_{k=1}^m \nabla \log f_{\theta}(u_k, y) \frac{f_{\theta}(u_k, y)}{\tilde{f}(u_k)} \right) < \infty. \quad (230)$$

Therefore, the gradient of the MCLA has finite variance and a central limit theorem.