

Fitting generalized linear mixed-effects models using `lme4`

Anna Ly
McMaster University

Rune Haubo Bojesen Christensen
Copenhagen Research Centre
for Biological and Precision Psychiatry

Douglas Bates
University of Wisconsin - Madison

Martin Mächler
ETH Zurich

Benjamin M. Bolker
McMaster University

Abstract

The `lme4` R package can be used to fit generalized linear mixed models (GLMMs), which extend the class of linear mixed models (LMMs). The two main extensions provided by GLMMs are (1) allowing for the conditional distribution of the response given the random effects to be non-Gaussian (e.g. binomial, Poisson) and (2) allowing the conditional mean to be a nonlinear function of a linear combination of the fixed and random effect coefficients, via an inverse link function. The conditional mode of the random effects given the observed data, the variance-covariance matrix of the random effects, and the fixed effect parameters are determined using penalized iteratively reweighted least squares (PIRLS). We compute an approximation of the integral over the distributions of the conditional modes to compute the MLE for a given set of parameters (by default we use the Laplace approximation or, alternatively, the more computationally expensive adaptive Gauss-Hermite quadrature). The package provides all the standard features available for GLMs in base R, including the the standard set of accessor functions as well as the possibility of user-specified distributions (within the exponential dispersion family) and link functions.

Keywords: sparse matrix methods, generalized linear mixed models, penalized least squares, Cholesky decomposition.

1. Introduction

1 `FIXME`: check digits of output?

2 The `lme4` package for R can be used to fit a broad range of mixed-effects models. One major
3 advantage of `lme4` over its predecessor, `nlme`, is that it can be used to fit generalized linear
4 mixed models (GLMMs), which combine the flexibility of linear mixed models (LMMs) and
5 generalized linear models (GLMs). In a companion paper, we have described the facilities
6 in `lme4` for fitting linear mixed models (LMMs). Here we describe the facilities for fitting

7 GLMMs.

2. Generalized Linear Mixed Models

8 Generalized linear mixed models extend the class of generalized linear models by allowing for
 9 both fixed and random effects. In a GL(M)M, the length- n vector-valued response variable,
 10 \mathcal{Y} , has a conditional distribution in the exponential dispersion family (e.g. Normal, binomial,
 11 Poisson).¹ The mean, $\mu_{\mathcal{Y}}$, of \mathcal{Y} depends on a linear predictor,

$$\eta = \mathbf{X}\beta. \quad (1)$$

12 where β is a p -dimensional coefficient vector and \mathbf{X} is an $n \times p$ model matrix. The mapping
 13 from $\mu_{\mathcal{Y}}$ to η , which is called the *link function* and written,

$$\mathbf{X}\beta = \eta = \mathbf{g}(\mu_{\mathcal{Y}}), \quad (2)$$

14 is a *diagonal mapping* in the sense that there is a scalar function, g , such that the i th com-
 15 ponent of η is g applied to the i th component of $\mu_{\mathcal{Y}}$. (The name “diagonal” reflects the fact
 16 that the Jacobian matrix, $\frac{d\eta}{d\mu}$, of such a mapping will be diagonal.) The scalar link function
 17 must be invertible and differentiable over its range. The vector-valued *inverse link* function,
 18 \mathbf{g}^{-1} , will be the scalar inverse link, g^{-1} , applied component-wise to η .

19 In the GLMM case, the vector of means of the exponential dispersion family distribution of \mathcal{Y}
 20 depends on an unobserved random vector, \mathcal{B} , of length q , called the random-effects coefficient
 21 vector. In particular, the conditional mean of \mathcal{Y} given that $\mathcal{B} = \mathbf{b}$, written $\mu_{\mathcal{Y}|\mathcal{B}=\mathbf{b}}$, depends
 22 on the linear predictor,

$$\eta = \mathbf{Z}\mathbf{b} + \mathbf{X}\beta. \quad (3)$$

23 where \mathbf{Z} is an $n \times q$ random-effects model matrix. Similar to the GLM case, the mapping
 24 from the conditional mean, $\mu_{\mathcal{Y}|\mathcal{B}=\mathbf{b}}$, to the linear predictor, η , is

$$\mathbf{Z}\mathbf{b} + \mathbf{X}\beta = \eta = \mathbf{g}(\mu_{\mathcal{Y}|\mathcal{B}=\mathbf{b}}), \quad (4)$$

25 The random vector \mathcal{B} is assumed to be distributed multivariate normally,

$$\mathcal{B} \sim \mathcal{N}(\mathbf{0}, \Sigma_{\theta}) \quad (5)$$

26 where Σ_{θ} is the covariance matrix of \mathcal{B} , which depends on a vector of covariance parameters,
 27 θ .

28 The optimization routines of **lme4** never actually compute Σ_{θ} directly, and instead use the
 29 elements of the covariance factor, Λ_{θ} , which is a matrix square root of Σ_{θ} (in practice a
 30 Cholesky factor),

$$\Sigma_{\theta} = \Lambda_{\theta}\Lambda_{\theta}'. \quad (6)$$

¹The **lme4** package supports fitting negative binomial GLMMs, which are an extension of GLMMs; the negative binomial family is within the exponential dispersion family if the dispersion parameter is fixed. The **glmer.nb** function uses a one-dimensional optimizer to estimate the dispersion parameter, fitting a negative binomial GLMM with a fixed dispersion parameter at each trial value. In practice this implementation is slower than those in other R packages (e.g., the **glmmTMB** package).

31 This characterization of the random-effects covariance structure allows us to write the linear
32 predictor as

$$\boldsymbol{\eta} = \mathbf{Z}\boldsymbol{\Lambda}_\theta\mathbf{u} + \mathbf{X}\boldsymbol{\beta}, \quad (7)$$

33 where the spherical random effects vector, \mathbf{u} (see Bates, Mächler, Bolker, and Walker (2015))
34 is a realization of the random vector, \mathcal{U} ,

$$\mathcal{U} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_q) \quad (8)$$

35 where \mathbf{I}_q is the identity matrix.

36 Common forms of the conditional distribution are Bernoulli, for binary responses, binomial
37 for binary responses that are recorded as the number of trials and the number of successes,
38 and Poisson, for count data. The combination of a distributional form and a link function
39 is called a *family*. For distributional forms in the exponential dispersion family there is a
40 *canonical link*. For Bernoulli or binomial forms the canonical link is the *logit* link function

$$\eta_i = \log\left(\frac{\mu_i}{1 - \mu_i}\right); \quad (9)$$

41 for the Poisson distribution the canonical link is the natural logarithm.

42 The form of the distribution determines the conditional variance, $\text{Var}(\mathcal{Y}|\mathcal{U} = \mathbf{u})$, as a function
43 of the conditional mean and, possibly, a separate scale factor. (In the most common cases
44 the conditional variance is completely determined by the conditional mean.)

45 Assuming that the conditional distribution belongs to the exponential dispersion family, we
46 can weight observations differently by scaling the dispersion parameter, ϕ . These **prior**
47 **weights**, \mathbf{w} , a vector of size n , are known positive constants (equal to the number of obser-
48 vations per trial in the particular case of a binomial GLMM).

49 By scaling the dispersion parameters, we also modify the conditional variance. Consider the
50 form, where $i = 1, 2, \dots, n$:

$$\text{Var}(\mathcal{Y}_i|\mathcal{U} = \mathbf{u}) = \frac{\phi}{w_i} \text{Var}(\mu_i) \quad (10)$$

51 where $\text{Var}(\mu_i)$ is the family-specific variance function. Higher weights indicate a smaller
52 variance.

53 If prior weights are not specified, by default \mathbf{w} is just a vector of 1's.

54 Another common modification when dealing with generalized linear mixed models is to allow
55 for an **offset**, which is an efficient way to include known scaling factors (like population or
56 area) without introducing additional parameters to the model.

57 In particular, the inclusion of an offset would modify the linear predictor as follows:

$$\boldsymbol{\eta} = \mathbf{Z}\mathbf{b} + \mathbf{X}\boldsymbol{\beta} + \text{offset}. \quad (11)$$

58 For GLMMs, the scale parameter is not typically estimated as part of the nonlinear mini-
59 mization of the negative likelihood over the $\boldsymbol{\beta}$ and θ parameters. As is common in generalized
60 linear model contexts, it is set equal to the (residual) deviance divided by the residual degrees
61 of freedom. (Except for the Binomial and Poisson families, where the scale parameter is fixed
62 at 1.)

63 The deviance itself is estimated using the method of moments estimator, which is the Pearson
 64 residual sum of squares. This approach is preferred by McCullagh and Nelder (1989) as it is
 65 less sensitive to small errors or model misspecification than the maximum likelihood estimator.

66 We therefore estimate the deviance using the penalized weighted Pearson residual sum of
 67 squares, as described in Section 3.2 of Bates *et al.* (2015).

68 The likelihood of the parameters, given the observed data, is now

$$L(\boldsymbol{\beta}, \boldsymbol{\theta} | \mathbf{y}_{\text{obs}}) = \int_{\mathbb{R}^q} f_{\mathcal{Y}, \mathcal{U}}(\mathbf{y}_{\text{obs}}, \mathbf{u}) d\mathbf{u} \quad (12)$$

69 where, as in the case of linear mixed models, $f_{\mathcal{Y}, \mathcal{U}}(\mathbf{y}_{\text{obs}}, \mathbf{u})$ is the unscaled conditional density
 70 of \mathcal{U} given $\mathcal{Y} = \mathbf{y}_{\text{obs}}$. The notation here is a bit blurred because, although the joint distribution
 71 of \mathcal{Y} and \mathcal{U} is always continuous with respect to \mathcal{U} , it can be (and often is) discrete with respect
 72 to \mathcal{Y} . However, when we condition on the observed value $\mathcal{Y} = \mathbf{y}_{\text{obs}}$, the resulting function is
 73 continuous with respect to \mathbf{u} so the unscaled conditional density is indeed well-defined as a
 74 density, up to a scale factor.

To evaluate the integrand in (12) we need the unscaled conditional density $f_{\mathcal{Y}, \mathcal{U}}$, which de-
 pends on the likelihood of \mathbf{y}_{obs} conditional on \mathbf{u} : $f(\mathbf{y}_{\text{obs}}, \mathbf{u}) = f(\mathbf{y}_{\text{obs}} | \mathbf{u})f(\mathbf{u})$. As usual in
 likelihood estimation we will eventually be working, for convenience, with negative log like-
 lihoods rather than likelihoods. However, base R's machinery for defining GLM models, or
 "families", does not provide a convenient expression for the negative log-likelihood: instead,
 objects of class `family` include a `dev.resids` function that returns a vector of observation-
 specific deviances, i.e. $-2(\log L(\mathbf{y}_{\text{obs}} | \mathbf{u}) - \log L_{\text{sat}})$ where L_{sat} is the likelihood corresponding
 to a saturated model with as many parameters as observations.² For example, the log-
 likelihood for a Poisson observation y with mean λ is (omitting the normalization constant)
 $\log L(y, \lambda) = y \log \lambda - \lambda$; the deviance is

$$-2(\log L(y, \lambda) - \log L(y, y)) = -2(y(\log \lambda - \log y) - (\lambda - y))$$

75 (with appropriate adjustments if $y = 0$). Because deviances differ from (negative) log likeli-
 76 hoods only by a parameter-independent constant (the log-likelihood of the saturated model)
 77 and a change of scale (multiplication by 2), this change does not affect our optimization
 78 procedure.

79 One advantage of using the pre-existing GLM family structure is that the software can thus
 80 fit models with user-specified families and link functions (although common families and link
 81 functions are hard-coded in C++ for computational speed), in contrast to early versions of
 82 `lme4` and some other R packages for GLMM fitting.

83 If `dev.resids` returns $\mathbf{d}(\mathbf{y}_{\text{obs}} | \mathbf{u})$, with elements $d_i(\mathbf{y}_{\text{obs}} | \mathbf{u}), i = 1, \dots, n$, the full deviance of a
 84 generalized linear model is

$$f(\mathbf{y}_{\text{obs}}, \mathbf{u}) = \sum_{i=1}^n d_i(\mathbf{y}_{\text{obs}} | \mathbf{u}) + \|\mathbf{u}\|^2.$$

²There is some confusion in R (and in its predecessor, S) about the exact definition of the deviance residuals of a family. As indicated above, we will use this name for the value of the `dev.resids` function for the family. The signed square root of this vector, using the signs of $\mathbf{y}_{\text{obs}} - \boldsymbol{\mu}$, is returned when the `residuals` method is applied to a fitted model of class `"glm"` when `type="deviance"`, the default, is specified. Both are called "deviance residuals" at different points in the documentation.

85 The likelihood can now be expressed as

$$L(\boldsymbol{\beta}, \boldsymbol{\theta} | \mathbf{y}_{\text{obs}}) = \int_{\mathbb{R}^q} \exp\left(-\frac{\sum_{i=1}^n d_i(\mathbf{y}_{\text{obs}} | \mathbf{u}) + \|\mathbf{u}\|^2}{2}\right) (2\pi)^{-q/2} d\mathbf{u} \quad (13)$$

86 As for linear mixed models, we simplify evaluation of the integral (12) by determining the
 87 value, $\tilde{\mathbf{u}}_{\beta, \theta}$, that maximizes the integrand. When the conditional density, $\mathcal{U} | \mathcal{Y} = \mathbf{y}_{\text{obs}}$, is
 88 multivariate Gaussian, this conditional mode will also be the conditional mean. However, for
 89 most families used in GLMMs, the mode and the mean need not coincide so we use the more
 90 general term and call $\tilde{\mathbf{u}}_{\beta, \theta}$ the *conditional mode*. We first describe the numerical methods for
 91 determining the conditional mode using the Penalized Iteratively Reweighted Least Squares
 92 (PIRLS) algorithm then return to the question of evaluating the integral (12).

93 2.1. Determining the conditional mode

94 The iteratively reweighted least squares (IRLS) algorithm is an efficient method of determining
 95 the maximum likelihood estimates of the coefficients in a generalized linear model. We extend
 96 it to a *penalized iteratively reweighted least squares* (PIRLS) algorithm for determining the
 97 conditional mode, $\tilde{\mathbf{u}}_{\beta, \theta}$. This algorithm has the form

98 1. Given parameter values, $\boldsymbol{\beta}$ and $\boldsymbol{\theta}$, and starting estimates, \mathbf{u}_0 , evaluate the linear pre-
 99 dictor, $\boldsymbol{\eta}$, the corresponding conditional mean, $\boldsymbol{\mu}_{\mathcal{Y} | \mathcal{U} = \mathbf{u}}$, and the conditional variance.
 100 Establish the weights as the inverse of the variance. We write these weights in the form
 101 of a diagonal weight matrix, \mathbf{W} , although they are stored and manipulated as a vector.

102 2. Let

$$Q(\mathbf{u}) = \left\| \mathbf{W}^{1/2} (\mathbf{y}_{\text{obs}} - \boldsymbol{\mu}_{\mathcal{Y} | \mathcal{U} = \mathbf{u}}) \right\|^2 + \|\mathbf{u}\|^2. \quad (14)$$

103 Solve the penalized, weighted, nonlinear least squares problem

$$\arg \min_{\mathbf{u}} (Q(\mathbf{u})). \quad (15)$$

104 3. Update the weights, \mathbf{W} , and check for convergence. If not converged, go to step 2.

105 We use a Gauss-Newton algorithm with an orthogonality convergence criterion (Bates and
 106 Watts 1988, §2.2.3) to solve the penalized, weighted, nonlinear least squares problem in step 2.
 107 At the i th iteration we determine an increment, $\boldsymbol{\delta}_i$, as the solution to the penalized, weighted,
 108 linear least squares problem

$$\boldsymbol{\delta}_i = \arg \min_{\boldsymbol{\delta}} \left\| \begin{bmatrix} \mathbf{W}^{1/2} (\mathbf{y}_{\text{obs}} - \boldsymbol{\mu}_i) \\ \mathbf{u}_i \end{bmatrix} - \begin{bmatrix} \mathbf{W}^{1/2} \mathbf{M}_i \mathbf{Z} \boldsymbol{\Lambda}_{\boldsymbol{\theta}} \\ \mathbf{I}_q \end{bmatrix} \mathbf{u} \right\|^2 \quad (16)$$

109 where \mathbf{u}_i is current value of \mathbf{u} , $\boldsymbol{\mu}_i$ is the corresponding conditional mean of $\mathcal{Y} | \mathcal{U} = \mathbf{u}_i$ and \mathbf{M}_i
 110 is the Jacobian matrix of the vector-valued inverse link, evaluated at $\boldsymbol{\mu}_i$. That is

$$\mathbf{M}_i = \left. \frac{d\boldsymbol{\mu}}{d\boldsymbol{\eta}'} \right|_{\boldsymbol{\eta}_i}, \quad (17)$$

111 which will be a diagonal matrix so, as for the weights, we store and manipulate the Jacobian
 112 as a vector.

113 When solving for the minimum of $Q(\mathbf{u})$, the proposed update for \mathbf{u} may increase rather than
 114 decrease the value of the objective function $Q(\mathbf{u})$, if the step size estimated by the Gauss-
 115 Newton step is too large. Let $Q(\mathbf{u}_i)$ represent the value of the objective function $Q(\mathbf{u})$ at the
 116 current iterate \mathbf{u}_i ; similarly, let $Q(\mathbf{u}_{i-1})$ represent the the value of the objective function of
 117 the previous iteration.

118 If $Q(\mathbf{u}_i) > Q(\mathbf{u}_{i-1})$, we proceed with a *step-halving* procedure. That is, we successively
 119 consider scaled updates of the form

$$\mathbf{u}_i^{(j)} = \mathbf{u}_{i-1} + \frac{\boldsymbol{\delta}_i}{2^j}, \quad j = 1, 2, \dots, \quad (18)$$

120 computing $Q(\mathbf{u}_i^{(j)})$ at each step. The step-halving procedure continues until a value of j is
 121 found such that $Q(\mathbf{u}_i^{(j)}) < Q(\mathbf{u}_{i-1})$, which the iterate is then updated by setting $\mathbf{u}_i = \mathbf{u}_i^{(j)}$.

122 If no such j is found after a predetermined number of halvings, then the step-halving procedure
 123 has failed and the algorithm terminates.

124 The minimizer, $\boldsymbol{\delta}_i$, of (16) satisfies

$$\mathbf{P} (\boldsymbol{\Lambda}'_\theta \mathbf{Z}' \mathbf{M}_i \mathbf{W} \mathbf{M}_i \mathbf{Z} \boldsymbol{\Lambda}_\theta + \mathbf{I}_q) \mathbf{P}' \boldsymbol{\delta}_i = \boldsymbol{\Lambda}'_\theta \mathbf{Z}' \mathbf{M}_i \mathbf{W} (\mathbf{y}_{\text{obs}} - \boldsymbol{\mu}_i) - \mathbf{u}_i \quad (19)$$

125 which we solve using the sparse Cholesky factor. At convergence, the factor, $\mathbf{L}_{\beta, \theta}$, satisfies

$$\mathbf{L}_{\beta, \theta} \mathbf{L}'_{\beta, \theta} = \mathbf{P} (\boldsymbol{\Lambda}'_\theta \mathbf{Z}' \mathbf{M} \mathbf{W} \mathbf{M} \mathbf{Z} \boldsymbol{\Lambda}_\theta + \mathbf{I}_q) \mathbf{P}' \quad (20)$$

126 As we show in the next section, the matrix $(\mathbf{L}_{\beta, \theta} \mathbf{L}'_{\beta, \theta})^{-1}$ is a Laplace approximation of the
 127 covariance matrix for the spherical random effects, conditional on the observed data. This fact
 128 is useful for constructing a nonlinear objective function for finding the approximate maximum
 129 likelihood estimates of θ and β .

130 2.2. Evaluating the likelihood for GLMMs using the Laplace approximation

131

132 Evaluating the likelihood for generalized linear mixed models requires approximating an in-
 133 tractable integral over the random effects distribution. The `glmer` function offers several
 134 approximations, controlled by the `nAGQ` argument. The default value of `nAGQ=1` specifies the
 135 *Laplace approximation* (Madsen and Thyregod 2011).

136 A second-order Taylor series approximation to $-2 \log[f_{\mathcal{Y}, \mathcal{U}}(\mathbf{y}_{\text{obs}}, \mathbf{u})]$ based at $\tilde{\mathbf{u}}$ provides an
 137 approximation of the unscaled conditional density as a multiple of the density for the multi-
 138 variate Gaussian $\mathcal{N}(\tilde{\mathbf{u}}, \mathbf{L}\mathbf{L}')$. The change of variable

$$\mathbf{u} = \tilde{\mathbf{u}} + \mathbf{L}\mathbf{z} \quad (21)$$

139 provides

$$\begin{aligned} L(\boldsymbol{\beta}, \boldsymbol{\theta} | \mathbf{y}_{\text{obs}}) &= \int_{\mathbb{R}^q} f_{\mathcal{Y}, \mathcal{U}}(\mathbf{y}_{\text{obs}}, \mathbf{u}) d\mathbf{u} \\ &\approx \tilde{f} |\mathbf{L}| \int_{\mathbb{R}^q} e^{-\|\mathbf{z}\|^2/2} (2\pi)^{-q/2} d\mathbf{z} \\ &= \tilde{f} |\mathbf{L}| \end{aligned} \quad (22)$$

140 or, on the deviance scale,

$$-2\ell(\boldsymbol{\beta}, \boldsymbol{\theta} | \mathbf{y}_{\text{obs}}) \approx \sum_{i=1}^n d_i(\mathbf{y}_{\text{obs}}, \tilde{\mathbf{u}}) + \|\tilde{\mathbf{u}}\|^2 + \log(|\mathbf{L}|^2) + \frac{q}{2} \log(2\pi) \quad (23)$$

141 The Laplace approximation normally conditions on both the fixed effects $\boldsymbol{\beta}$ and the variance-
 142 covariance parameters $\boldsymbol{\theta}$. A further approximation, which is denoted in `glmer` by `nAGQ=0`,
 143 profiles out the fixed effects by minimizing β and the conditional modes \mathbf{u} simultaneously in
 144 eq. 15. This approximation is exact when (1) $\partial(\log L)/\partial\beta$ is a linear function of the conditional
 145 modes \mathbf{u} and (2) when the conditional mode is equal to the conditional mean (typically,
 146 although not necessarily, implying a symmetric conditional distribution). Both assumptions
 147 hold for linear mixed models (although a Laplace approximation is not necessary there),
 148 consistent with Bates *et al.* (2015) showing that the fixed effects can be profiled out of the
 149 log-likelihood for LMMs. The Julia `MixedModels.jl` package offers the same approximation
 150 as the `fast` argument to the `pirls!` function (<https://juliastats.org/MixedModels.jl/stable/optimization/>); Template Model Builder (Kristensen, Nielsen, Berg, Skaug, and
 151 Bell 2016), and downstream packages such as `glmmTMB`, provide this function via a `profile`
 152 argument.
 153

154 By default, `glmer` uses a two-stage optimization procedure (described below) with `nAGQ=0` in
 155 the first stage; users can also specify `nAGQ=0` for faster, approximate model fits.

156 *Decomposing the deviance for simple models*

157 A common special case of mixed models is those where scalar (typically intercept) random
 158 effects are associated with levels of a single grouping factor, \mathbf{h} . In this case the dimension, q ,
 159 of the random effects is the number of levels of \mathbf{h} — i.e. there is exactly one random effect
 160 associated with each level of \mathbf{h} . We will write the vector of variance-covariance parameters,
 161 which is one-dimensional, as a scalar, θ . The matrix $\boldsymbol{\Lambda}_\theta$ is a multiple of the identity, $\theta\mathbf{I}_q$,
 162 and \mathbf{Z} is the $n \times q$ matrix of indicators of the levels of \mathbf{f} . The permutation matrix, \mathbf{P} , can
 163 be set to the identity and \mathbf{L} is diagonal, although not necessarily homogeneous (i.e., a scalar
 164 multiple of the identity matrix).

165 Because each element of $\boldsymbol{\mu}$ depends on only one element of \mathbf{u} and the elements of \mathcal{Y} are
 166 conditionally independent, given $\mathcal{U} = \mathbf{u}$, the conditional densities of the $u_j, j = 1, \dots, q$ given
 167 $\mathcal{Y} = \mathbf{y}_{\text{obs}}$ are independent. We partition the indices $1, \dots, n$ as $\mathbb{I}_j, j = 1, \dots, q$ according to
 168 the levels of \mathbf{h} . That is, the index i is in \mathbb{I}_j if $h_i = j$. This partitioning also applies to the
 169 deviance residuals in that the i th deviance residual depends only on u_j when $i \in \mathbb{I}_j$.

170 Writing the univariate conditional densities as

$$f_j(\mathbf{y}_{\text{obs}}, u_j) = \exp\left(-\frac{\sum_{i \in \mathbb{I}_j} d_i(\mathbf{y}_{\text{obs}}, u_j) + u_j^2}{2}\right) (2\pi)^{-1/2} \quad (24)$$

171 we have

$$f_{\mathcal{Y}, \mathcal{U}}(\mathbf{y}_{\text{obs}}, \mathbf{u}) = \prod_{j=1}^q f_j(\mathbf{y}_{\text{obs}}, u_j) \quad (25)$$

172 and

$$L(\boldsymbol{\beta}, \boldsymbol{\theta} | \mathbf{y}_{\text{obs}}) = \prod_{j=1}^q \int_{\mathbb{R}} f_j(\mathbf{y}_{\text{obs}}, u) du \quad (26)$$

173 We consider this special case both because it occurs frequently and because, for some software,
 174 it is the only type of GLMM that can be fit. Also, in this particular case we can graphically
 175 assess the quality of the Laplace approximation by comparing the actual integrand to its
 176 approximation.

177 Consider the `cbpp` data on contagious bovine pleuropneumonia (CBPP) incidence according
 178 to season and herd, available in the `lme4` package (see 5.1 for more details), and the model

```
> print(m1 <- glmer(cbind(incidence, size-incidence) ~ period + (1|herd),
+   cbpp, binomial), corr=FALSE)
```

```
Generalized linear mixed model fit by maximum likelihood (Laplace
  Approximation) [glmerMod]
Family: binomial ( logit )
Formula: cbind(incidence, size - incidence) ~ period + (1 | herd)
Data: cbpp
      AIC      BIC    logLik -2*log(L)  df.resid
194.0531 204.1799 -92.0266  184.0531      51
Random effects:
Groups Name      Std.Dev.
herd (Intercept) 0.6421
Number of obs: 56, groups: herd, 15
Fixed Effects:
(Intercept)      period2      period3      period4
      -1.3983      -0.9919      -1.1282      -1.5797
```

179 This model has been fit by minimizing the Laplace approximation to the deviance. We can
 180 assess the quality of this approximation by evaluating the unscaled conditional density at
 181 $u_j(z) = \tilde{u}_j + z/\mathbf{L}_{j,j}$ and comparing the ratio, $f_j(\mathbf{y}_{\text{obs}}, u)/(\tilde{f}_j\sqrt{2\pi})$, to the standard normal
 182 density, $\phi(z) = e^{-z^2/2}/\sqrt{2\pi}$, as shown in Figure 1.

183 As Figure 1 shows, the univariate integrands are very close to the standard normal density,
 184 indicating that the Laplace approximation to the deviance is a good approximation in this
 185 case.

3. Adaptive Gauss-Hermite quadrature for GLMMs

186 When the integral (12) can be expressed as a product of low-dimensional integrals, we can
 187 use Gauss-Hermite quadrature to provide a closer approximation to the integral. Univariate
 188 Gauss-Hermite quadrature evaluates the integral of a function that is multiplied by a “kernel”
 189 where the kernel is a multiple of e^{-z^2} or $e^{-z^2/2}$. For statisticians the natural candidate is the
 190 standard normal density, $\phi(z) = e^{-z^2/2}/\sqrt{(2\pi)}$. A k th-order Gauss-Hermite formula provides
 191 knots, $z_i, i = 1, \dots, k$, and weights, $w_i, i = 1, \dots, k$, such that

$$\int_{\mathbb{R}} t(z)\phi(z) dz \approx \sum_{i=1}^k w_i t(z_i)$$

192 The function `GHrule` in `lme4` (based on code in the `SparseGrid` package) provides knots and
 193 weights relative to the standard normal kernel for orders k from 1 to 100. For example,

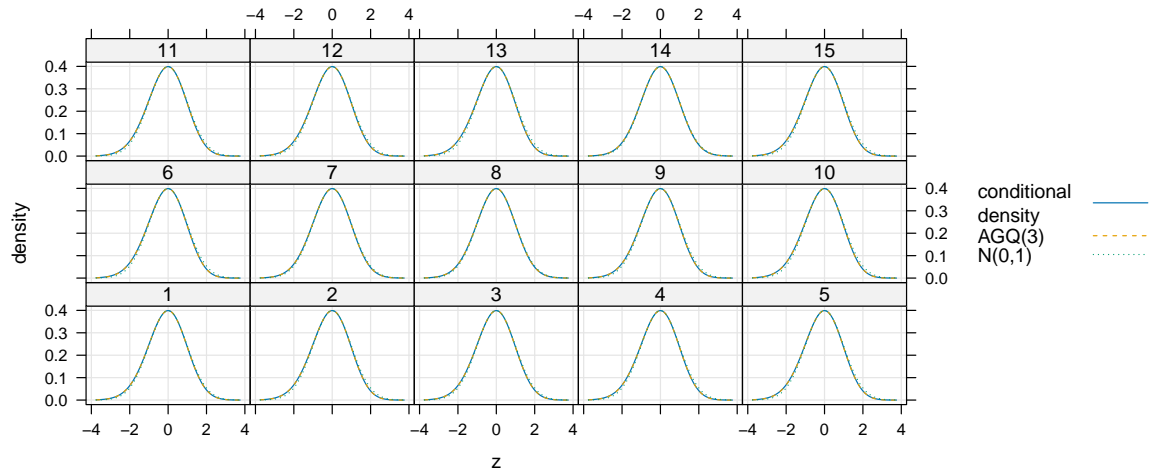


Figure 1: Comparison of univariate integrands (solid, blue line); 3-point Gauss-Hermite quadrature (dashed, yellow); and standard normal density function (green, dotted) for the CBPP model. For this model, the three approximations are nearly indistinguishable.

```
> GHrule(5)
```

```

      z      w      ldnorm
[1,] -2.856970 0.01125741 -5.0000774
[2,] -1.355626 0.22207592 -1.8377997
[3,]  0.000000 0.53333333 -0.9189385
[4,]  1.355626 0.22207592 -1.8377997
[5,]  2.856970 0.01125741 -5.0000774

```

194 where \mathbf{z} is the vector of knots, \mathbf{w} is the vector of weights, and \mathbf{dlnorm} is the log-density of the
 195 standard normal distribution at z .

196 The choice of the value of k depends on the behavior of the function $t(z)$. If $t(z)$ is a polynomial
 197 of degree $2k - 1$ then the Gauss-Hermite formula for orders k or greater provides an exact
 198 answer. The fact that we want $t(z)$ to behave like a low-order polynomial is often neglected
 199 in the formulation of a Gauss-Hermite approximation to a quadrature. The quadrature knots
 200 on the u scale are chosen as

$$u_{i,j}(z) = \tilde{u}_j + z_i / \mathbf{L}_{j,j}, \quad i = 1, \dots, k; \quad j = 1, \dots, q \quad (27)$$

201 exactly so that the function $t(z)$ should behave like a low-order polynomial over the region of
 202 interest, which is to say the region where quadrature knots with large weights are located. The
 203 term “adaptive Gauss-Hermite quadrature” reflects the fact that the approximating Gaussian
 204 density is scaled and shifted to provide a second order approximation to the logarithm of the
 205 unscaled conditional density.

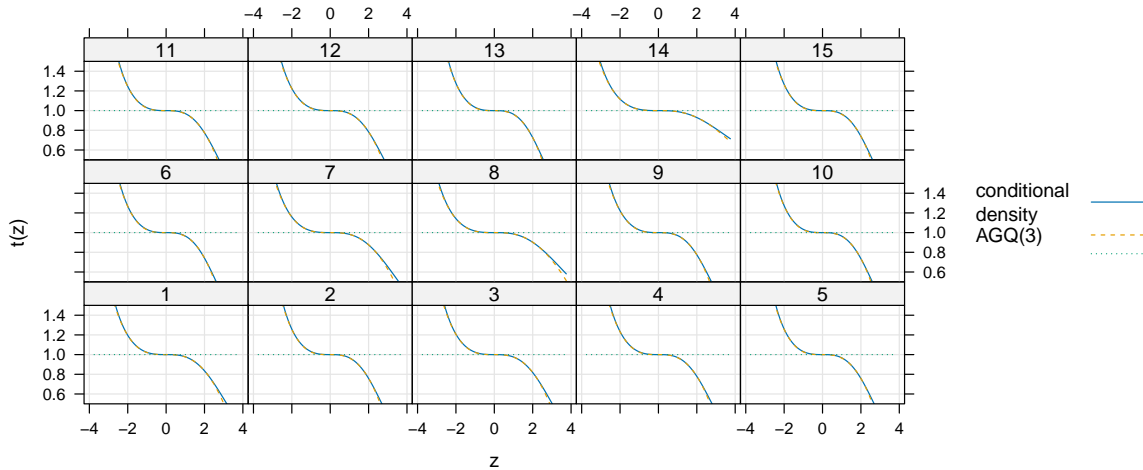


Figure 2: The function $t(z)$, which is the ratio of the normalized unscaled conditional density to the standard normal density, for each of the univariate integrals in the evaluation of the deviance for model `m1` (blue, solid line). As in Figure 1, the dashed yellow line shows the approximation for 3-point adaptive Gauss-Hermite quadrature. The horizontal green reference line ($z = 1$) gives the reference value that would apply for the Laplace approximation, which assumes the conditional density is exactly equal to the standard normal. The y -axis limits are truncated to $(0.5, 1.5)$.

206 Figure 2 shows $t(z)$ for each of the unidimensional integrals in the likelihood for the model `m1`
 207 at the parameter estimates. While this view shows the deficiency of the Laplace approximation
 208 (deviation from a horizontal reference line at $z = 1$ clearly, The AGQ(3) approximation fits
 209 extremely well over the range shown. The tails of the polynomials implied by the AGQ
 210 approximation can fluctuate widely, but these fluctuations are suppressed by multiplying by
 211 the thin tails of the Gaussian distribution.

212 The CBPP data set is a relatively well-behaved data set, where Laplace approximation works
 213 well. In contrast, a widely used data set on toenail onychomycosis (De Backer, De Vroey,
 214 Lesaffre, Scheys, and De Keyser 1998), which has a very low effective sample size per cluster
 215 — an average of about 6.5 binary observations (“moderate or severe” vs “none or mild”
 216 disease) per patient — represents an example where Gauss-Hermite quadrature is necessary
 217 for reliable results. In this case the conditional densities depart much more clearly from the
 218 standard normal (Figure 3; note the scale of the density ratios goes from 0 to 10, in contrast
 219 the maximum of 4 in Figure 2). Stringer, Bilodeau, and Tang (2022) and Stringer (2024)
 220 further explore the limitations of Laplace approximation.

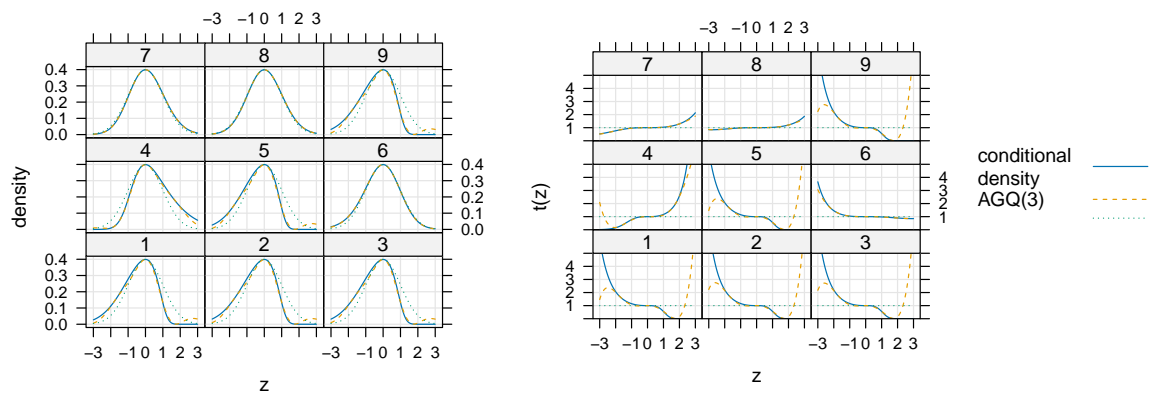


Figure 3: Normalized unscaled conditional density (left) and ratio of density to the standard normal density (right) for a random sample of 9 patients from the toenail onychomycosis data set. The y -axis limits in the right panel ($t(z)$) are truncated to $(0, 5)$.

221 To use adaptive Gauss-Hermite quadrature for model fitting in `glmer` models, set the argu-
 222 ment `nAGQ`, the number of quadrature points, to a value greater than 1. Increasing the number
 223 of nodes generally improves the accuracy of the likelihood approximation at the expense of
 224 computation time — although rounding errors may accumulate when using large numbers of
 225 quadrature points. At present, AGQ is only available for models with a single scalar random
 226 effect. (The `GLMMadaptive` package implements AGQ for vector-valued random effect models,
 227 although it is still restricted to models with a single random effect.)

4. Model fitting

228 Once we can calculate the deviance by PIRLS for specified values of θ and β (or only θ
 229 if profiling out the fixed effects via `nAGQ=0`), we then estimate the parameters by nonlinear
 230 optimization. This procedure largely follows the description in [Bates et al. \(2015\)](#), using
 231 derivative-free optimizers with box constraints to prevent non-positive-(semi)definite covari-
 232 ance matrices. Specifically, the elements of θ corresponding to the diagonal of Λ_θ are currently
 233 constrained to be non-negative.

234 The only difference is that by default `glmer` uses a two-step fitting procedure, using `nAGQ=0`
 235 at the first stage to get preliminary estimates which are then used as starting points for a
 236 second optimization with Laplace approximation or Gauss-Hermite quadrature as specified
 237 by the user. Different nonlinear optimizers can be used at each stage: the current default, as
 238 specified in `glmerControl`, is to use Powell’s BOBYQA followed by a box-constrained variant
 239 of the Nelder-Mead simplex algorithm. For faster, approximate fitting, the second stage can
 240 be omitted; in the rare cases where the initial `nAGQ=0` fit gives poor results, the first stage
 241 can be skipped via `glmerControl(nAGQ0initStep = FALSE)`.

242 However, because the profiled log-likelihood is an even function of the diagonal elements of
 243 θ and $\Sigma_\theta = \Lambda_\theta \Lambda_\theta'$ depends on them only through their squares, positive and negative values
 244 yield identical likelihoods. This symmetry means that constrained optimization is not strictly
 245 necessary — an unconstrained optimizer will converge to a correct solution, approaching zero
 246 from either side in the boundary case of a singular random-effects covariance matrix. (Once a
 247 solution is found, we can map it to a unique solution where the diagonal elements are all non-
 248 negative.) A similar approach to removing constraints could work for structured covariance
 249 matrices where correlation parameters are constrained to $(-1, 1)$, e.g. by parameterizing the
 250 model in terms of a phase parameter p where $\rho = \sin(p)$ (and then mapping p to $(0, 2\pi)$). Re-
 251 moving these constraint would permit the use of a broader class of unconstrained optimizers,
 252 though this feature has not yet been incorporated into `lme4`.

5. Examples

253 5.1. CBPP

254 The `?cbpp` help page describes the CBPP data set ([Lesnoff et al. 2004](#)) as follows:

255 Contagious bovine pleuropneumonia (CBPP) is a major disease of cattle in Africa,
 256 caused by a mycoplasma. This dataset describes the serological incidence of CBPP
 257 in zebu cattle during a follow-up survey implemented in 15 commercial herds

258 located in the Boji district of Ethiopia. The goal of the survey was to study
 259 the within-herd spread of CBPP in newly infected herds. Blood samples were
 260 quarterly collected from all animals of these herds to determine their CBPP status.
 261 These data were used to compute the serological incidence of CBPP (new cases
 262 occurring during a given time period). Some data are missing (lost to follow-up).

263 Lesnoff *et al.* (2004) estimated the effects of different treatments using (1) ordinary logistic
 264 regression incorporating a variance-inflation factor, also known as a quasi-binomial model
 265 (“logistic regression” is sometimes used specifically to describe analyses of Bernoulli responses,
 266 but in this case there are multiple trials per observation [cows that could become seropositive],
 267 and so a dispersion or scale parameter can be estimated); (2) a GLMM implemented in `lme4`;
 268 and a (3) Markov chain Monte Carlo algorithm Zeger and Karim (1991), which as they state
 269 allows for a non-parametric rather than a Normal model for the random effects. The authors
 270 did not find any significant effects of treatment, ascribing the null results to “a lack of power in
 271 the statistical analyses or to a quality problem for the medications used (and more generally,
 272 for health-care delivery in the Boji district).”

273 (Note that Table 1 of Lesnoff *et al.* (2004) contains a known typographical error for herd 6.
 274 Consequently, results obtained using the `cbpp` data set may not exactly reproduce some of
 275 the findings reported in that paper.)

276 The `lme4` package includes two variants of the `cbpp` data set. The second variant, `cbpp2`,
 277 contains corrected values corresponding to Table 1 of Lesnoff *et al.* (2004). Although the
 278 precise provenance of these data sets is unclear, the `cbpp` data set matches the version held
 279 by the corresponding author of Lesnoff *et al.* (2004).

280 We model proportional incidence as a binomial response depending on the additive fixed
 281 effects of period, treatment and average herd size; to account for repeated measures we fit a
 282 model with a random effect of herd. In practice we might also be interested in a period by
 283 treatment interaction, but we neglect that term here.

284 As in `glm`, we can specify a binomial response as a proportion and use the `weights` argument to
 285 specify the sample size, instead of the more typical two-column `cbind(successes, failures)`
 286 format:

```
> gm1 <- glmer(incidence/size ~ period + treatment + avg_size + (1 | herd),
+             family = binomial,
+             data = cbpp2, weights = size)
```

287 It is also worth considering adding an observation-level random effect to the model, which
 288 we can do by creating a new factor based on observation number and using `update()` on the
 289 previous model:

```
> cbpp2 <- transform(cbpp2, obs=factor(seq(nrow(cbpp2))))
> gm2 <- update(gm1, .~.(1|obs)) ## herd and observation-level REs
> gm3 <- update(gm1, .~.(1|herd)+(1|obs)) ## observation-level REs only
```

290 Model summary

291 The first part of the summary reiterates the family and link function used, the model formula,

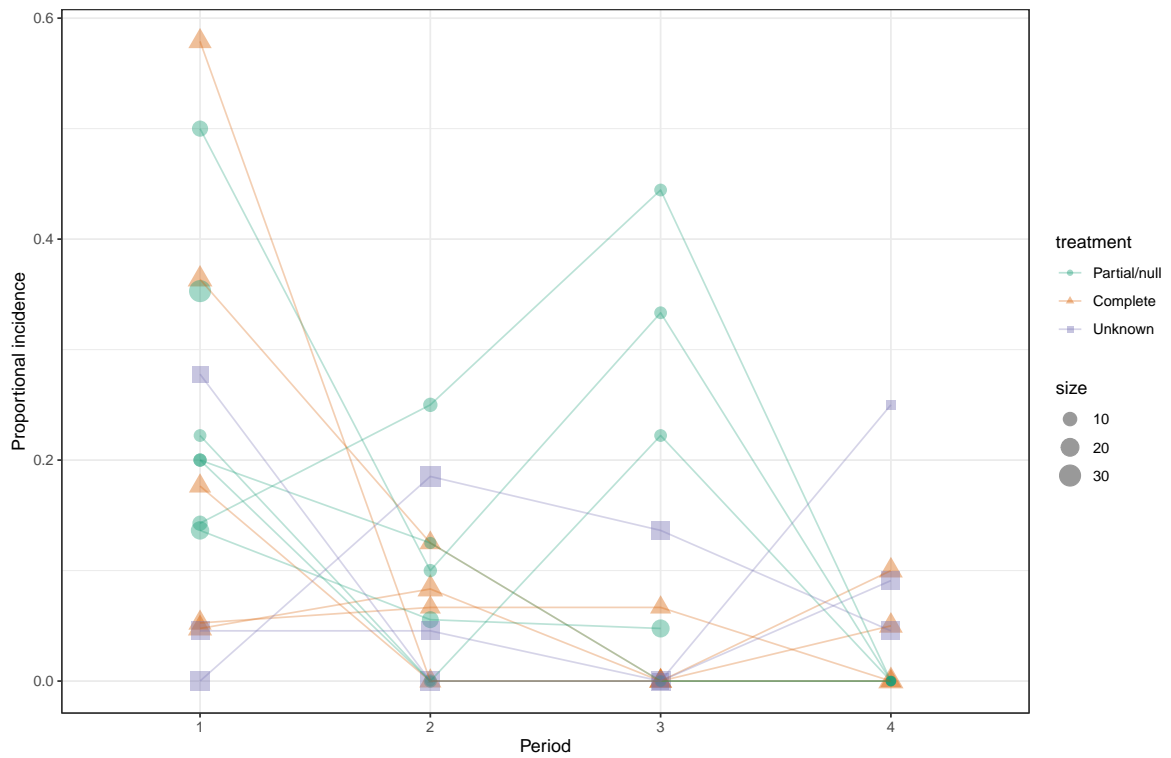


Figure 4: Incidence (proportion of cows becoming seropositive per observation period) vs. period. Colours show treatment category for each herd; point sizes reflect the number of seronegative cows at the start of each period. Lines connect the sets of observations from each herd.

292 and gives various summary statistics (log-likelihood etc.), as well as quantiles of the scaled
 293 (Pearson) residuals:

```
Generalized linear mixed model fit by maximum likelihood (Laplace
Approximation) [glmerMod]
Family: binomial ( logit )
Formula: incidence/size ~ period + treatment + avg_size + (1 | herd)
Data: cbpp2
Weights: size

      AIC      BIC    logLik -2*log(L)  df.resid
197.8    214.0    -90.9     181.8      48

Scaled residuals:
   Min      1Q  Median      3Q      Max
-2.2311 -0.7967 -0.3732  0.4684  2.7557
```

294 These quantities are also accessible via standard accessors (`AIC()`, `BIC()`, `logLik()`).

295 The next chunk of `summary()` describes the random effects and the number of levels associated
 296 with each grouping factor (the latter is useful for checking that random-effects formulae have
 297 been specified correctly):

```
Random effects:
 Groups Name          Variance Std.Dev.
 herd   (Intercept) 0.3116   0.5582
Number of obs: 56, groups: herd, 15
```

298 This information is also accessible via `VarCorr()`, which returns a list of variance-covariance
 299 matrices (the `print` method for `VarCorr` objects allows control of whether the variance, or
 300 standard deviation, or both, are printed).

301 Next come the estimates of the fixed effects, along with Wald estimates of the standard error,
 302 Z statistic, and p -value:

```
Fixed effects:
              Estimate Std. Error z value Pr(>|z|)
(Intercept)  -1.005623   0.708418  -1.420 0.155743
period2      -0.986283   0.303381  -3.251 0.001150
period3      -1.125147   0.323142  -3.482 0.000498
period4      -1.561098   0.422631  -3.694 0.000221
treatmentComplete -0.376225   0.503000  -0.748 0.454483
treatmentUnknown -0.683246   0.645179  -1.059 0.289599
avg_size     -0.006135   0.045608  -0.135 0.893002
```

303 One can use `coef(summary())` to retrieve this information, and optionally format it with
 304 `printCoefmat()`.

305 The last component of `summary()` gives the estimated correlations among the fixed-effect
 306 parameters, which can be useful for assessing multicollinearity (it can also be overwhelming:

307 it is suppressed by default for models with more than 20 fixed-effect parameters, and can also
 308 be suppressed by using `print(summary(.),correlation=FALSE)`).

```
Correlation of Fixed Effects:
              (Intr) perid2 perid3 perid4 trtmnC trtmnU
period2      -0.135
period3      -0.130  0.278
period4      -0.086  0.210  0.195
trtmntCmplt  0.289 -0.017 -0.021 -0.059
trtmntUnknw  0.431 -0.053 -0.045 -0.043  0.588
avg_size     -0.910  0.026  0.028  0.020 -0.547 -0.649
```

309 *Diagnostics*

310 A range of graphical diagnostic tools is available for `merMod` objects. The plot methods in
 311 the `lme4` package are inspired by those in the `nlme` package, using `lattice` plots to provide
 312 a reasonable blend of convenience and flexibility.

313 `merMod` objects are also compatible with the `performance` package and the `DHARMA` package,
 314 both commonly used for model checking.

315 The following code produces a standard range of diagnostic plots (Figure 5), similar to the
 316 ones in base R's `plot.lm` method. These diagnostics will generally be useful for models
 317 where the conditional density is approximately normal (but heteroscedastic) — e.g., Poisson
 318 responses with large mean or binomial responses with large numbers of successes — and less
 319 so otherwise, e.g. for binary responses.

```
> ## basic residual plot
> plot(gm1)
> ## scale-location plot
> plot(gm1,sqrt(abs(resid(.)))~fitted(.),type=c("p","smooth"))
> ## boxplot of residuals grouped by a categorical predictor
> plot(gm1,period~resid(.))
> ## Q-Q plot
> qqmath(gm1)
```

320 The `ranef()` accessor extracts the conditional modes; the argument `condVar=TRUE` addition-
 321 ally extracts the variances of the conditional modes, which are stored as an attribute labelled
 322 "`postVar`" — a three-dimensional array that gives the variance-covariance matrix of the con-
 323 ditional modes for each level of the grouping variable. The plotting methods `dotplot()` and
 324 `qqmath()` return lists of graphical objects showing *caterpillar plots* (ordered values of the
 325 random effects with confidence bars); in the case of the Q-Q plot (`qqmath`) the *y*-axis shows
 326 corresponding values of the standard normal quantiles (Figure 6).

327 `performance::check_model()` checks a variety of classical assumptions of GLMs. For mixed-
 328 effect models, it also assesses the normality of the distribution of conditional modes (Figure 7).

329 The `DHARMA` package is also compatible with `merMod` objects. `DHARMA` generates simulation-
 330 based residuals for generalized linear (mixed) models and uses them for graphical and statis-
 331 tical tests of model assumptions (Figure 8).

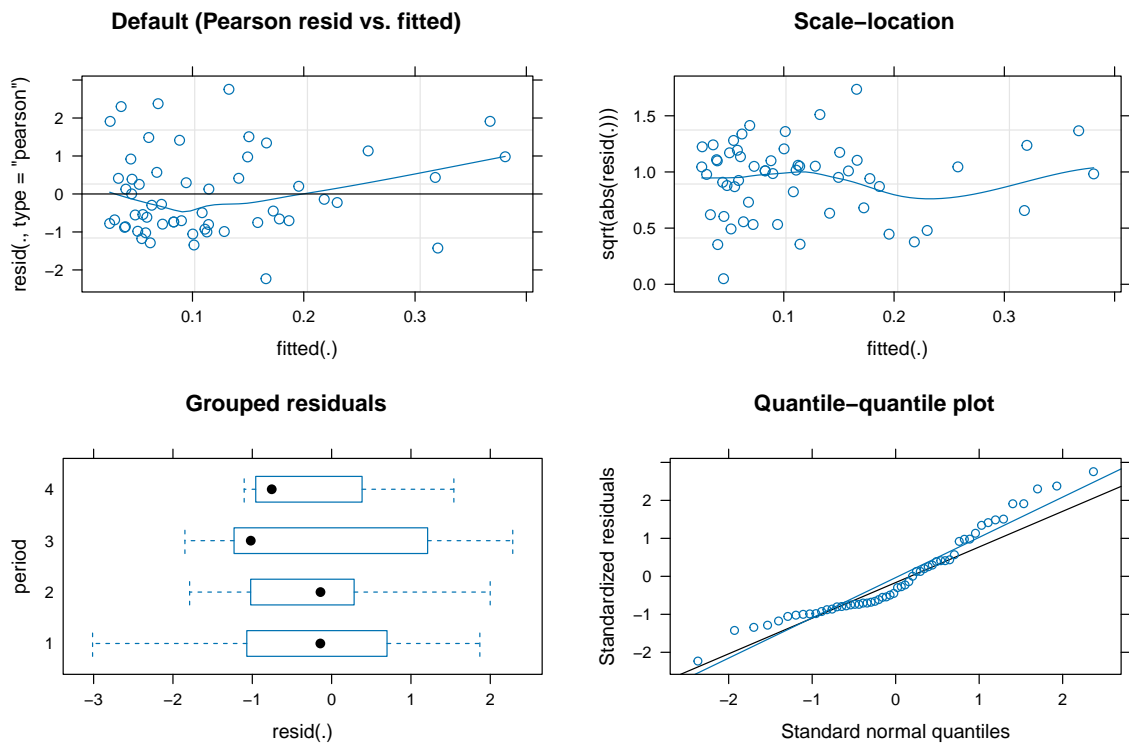


Figure 5: Graphical diagnostics (using package built-in functions)

332 Having checked the diagnostics, we would now like to compare the three models we have fitted.
 333 Inspecting the `VarCorr` components, we see that when we fit both herd- and observation-level
 334 random effects, the among-herd variance is estimated as zero. The appropriate procedure at
 335 this point (e.g. whether one drops non-significant terms, or those with small scaled magni-
 336 tudes, or those that worsen the AIC or BIC of the model) depends on the goals of the analysis
 337 and one's philosophy of model-building (Barr, Levy, Scheepers, and Tily 2013; Matuschek,
 338 Kliegl, Vasishth, Baayen, and Bates 2017; Scandola and Tidoni 2024).

339 One might either stick with the full model, or continue with the reduced model with observation-
 340 level random effects only (as it has exactly the same likelihood as the full model but uses an
 341 additional parameter, it would be chosen according to either an information-theoretic or a
 342 hypothesis-testing model selection framework).

343 Here we will start by computing likelihood profiles and confidence intervals for the model
 344 incorporating both random effects; although it has the same point estimates and maximum
 345 likelihood as the reduced model, confidence intervals that incorporate non-local information
 346 (i.e. profile- or parametric bootstrap-based) will give different, more conservative results for
 347 the full model.

348 The `profile` method computes profile likelihoods. The computation can be slow, since
 349 complete profiling for a model with p random- and fixed-effect parameters requires fitting p
 350 profiles, each of which requires many $p - 1$ -dimensional optimizations. The machinery for
 351 generating likelihood profiles for GLMMs is similar to that for LMMs (see Bates *et al.* (2015),
 352 § 5.1).

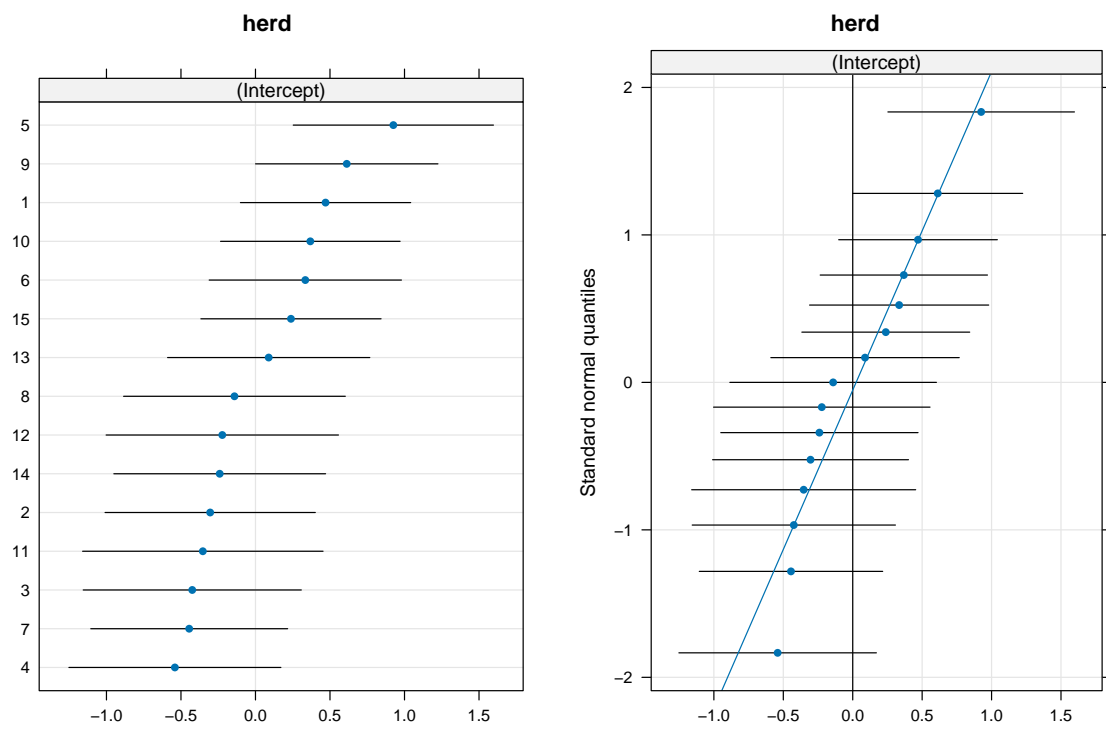


Figure 6: Graphical display of random effects. *Left:* conditional modes $\pm 1.96 \times$ conditional standard deviation, ordered by magnitude. *Right:* quantile-quantile plot, with linear regression line overlaid.

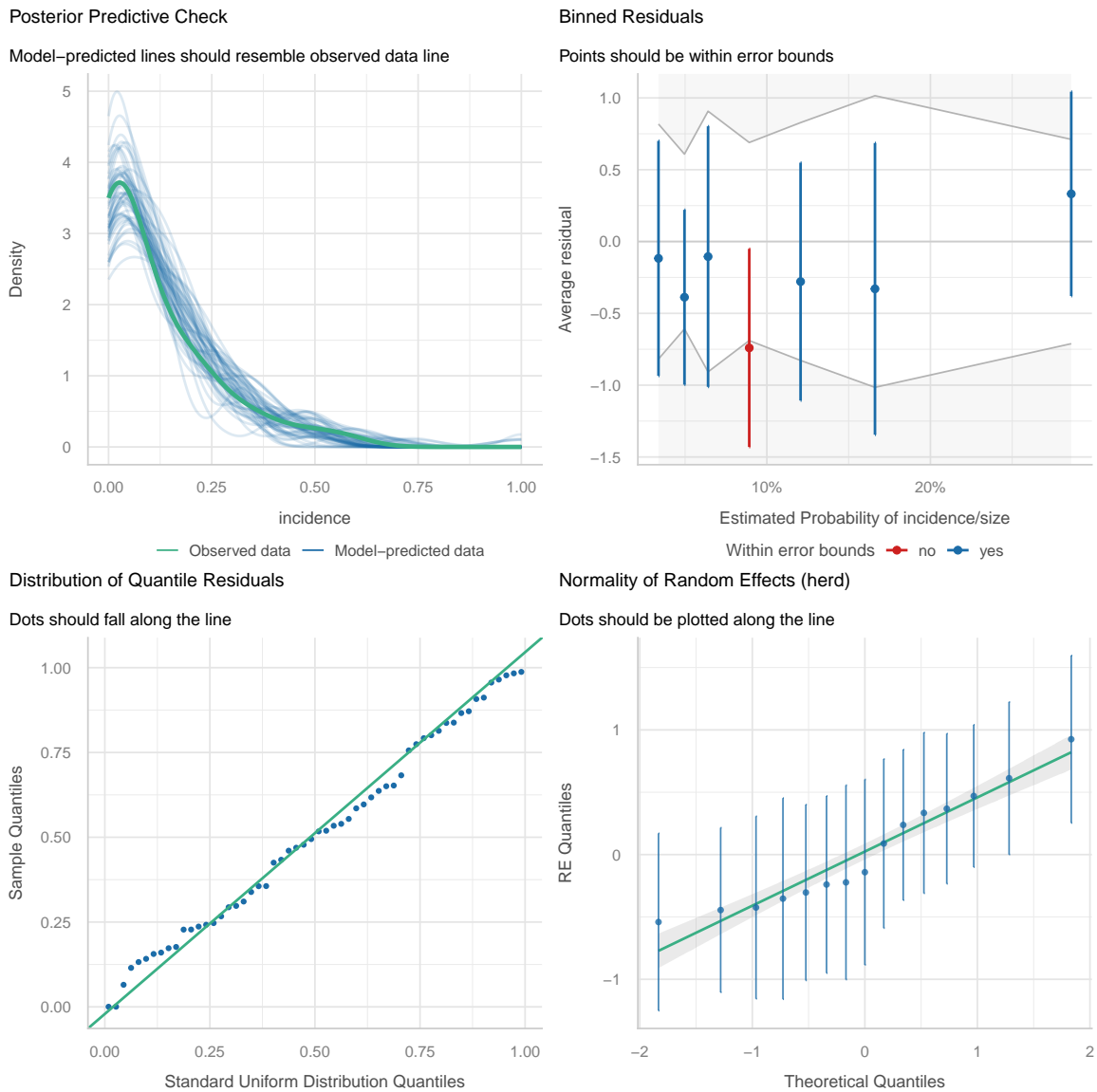


Figure 7: Model diagnostics using `performance::check_model`. The plot showcasing the normality of random effects (lower right panel) is the transpose of the right panel in Figure 6.

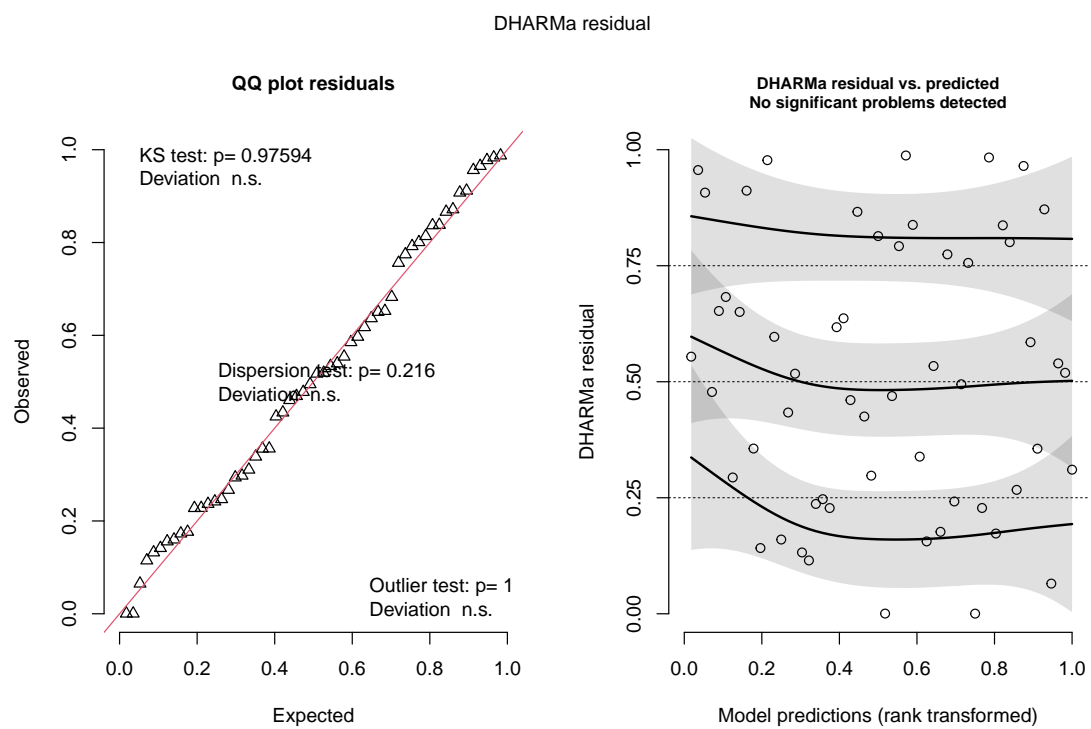


Figure 8: Model assumption checks using the DHARMA package. *n.s.* denotes “not significant”.

353 The `profile` method returns an object of class `thpr` — a data frame containing the profiles,
 354 augmented with attributes containing interpolation splines for each parameter profile and
 355 their inverses (using `splines::interpSpline` and `splines::backSpline`); the latter are
 356 used for plotting profiles and computing confidence intervals. An `as.data.frame` method
 357 adds `.focal` and `.par` variables to the data frame, useful for customized plots.

358 Profiles can be used for univariate (`xyplot`) and bivariate (`splom`) profile plots, and to com-
 359 pute profile confidence intervals (`confint`). (`confint` applied to a `glmer` fit will first fit the
 360 profile, then use it to compute profile confidence intervals. Given the computational cost of
 361 profiling, it makes sense to compute and save the profile as an intermediate step if one plans
 362 to do anything other than computing confidence intervals.)

363 Two other common methods for computing confidence intervals are parametric bootstrapping
 364 (`method = "boot"`) and the classical Wald approximation (`method = "Wald"`). Parametric
 365 bootstrapping is much slower, but more accurate (and, via the `FUN` argument, can generate
 366 confidence intervals for any quantity that can be derived from a fitted model). The Wald
 367 approximation is faster and less accurate than profile confidence intervals. By default `glmer`
 368 only returns estimates for the fixed-effects parameters, as the assumptions of the Wald ap-
 369 proximation are often violated badly for random-effects (co)variances and correlations. In the
 370 examples below we use the finite-difference Hessian (second derivative matrix of the estimated
 371 parameters) and the delta method to compute Wald confidence intervals for random-effects
 372 standard deviations and correlations, when possible.

373 Figure 9 compares all three of these confidence intervals across all three of the models fitted.

374 If the default optimizers (Nelder-Mead followed by BOBYQA) do not perform well, one could
 375 attempt to re-fit a [g]lmer model with a variety of different optimizers using `allFit()`. To see
 376 which optimizers `lme4` currently supports, one can use the command `allFit(show.meth.tab=TRUE)`
 377 to show all of the different methods.

378 To use `allFit()`, supply the initial fitted model as input. Useful results are obtained from
 379 `summary(allFit())`:

- 380 • `$which.OK` returns which optimizers worked. Below only applies to optimizers that were
 381 successful.
- 382 • `$l1lik` returns log-likelihoods
- 383 • `$fixef` fixed-effect estimates
- 384 • `$sdcor` random-effect standard deviations and correlations
- 385 • `$theta` random-effect parameters on the Cholesky scale

386 We will show the summary results for `$sdcor`; the other components are similar.

```
> gm_all <- allFit(gm1)
```

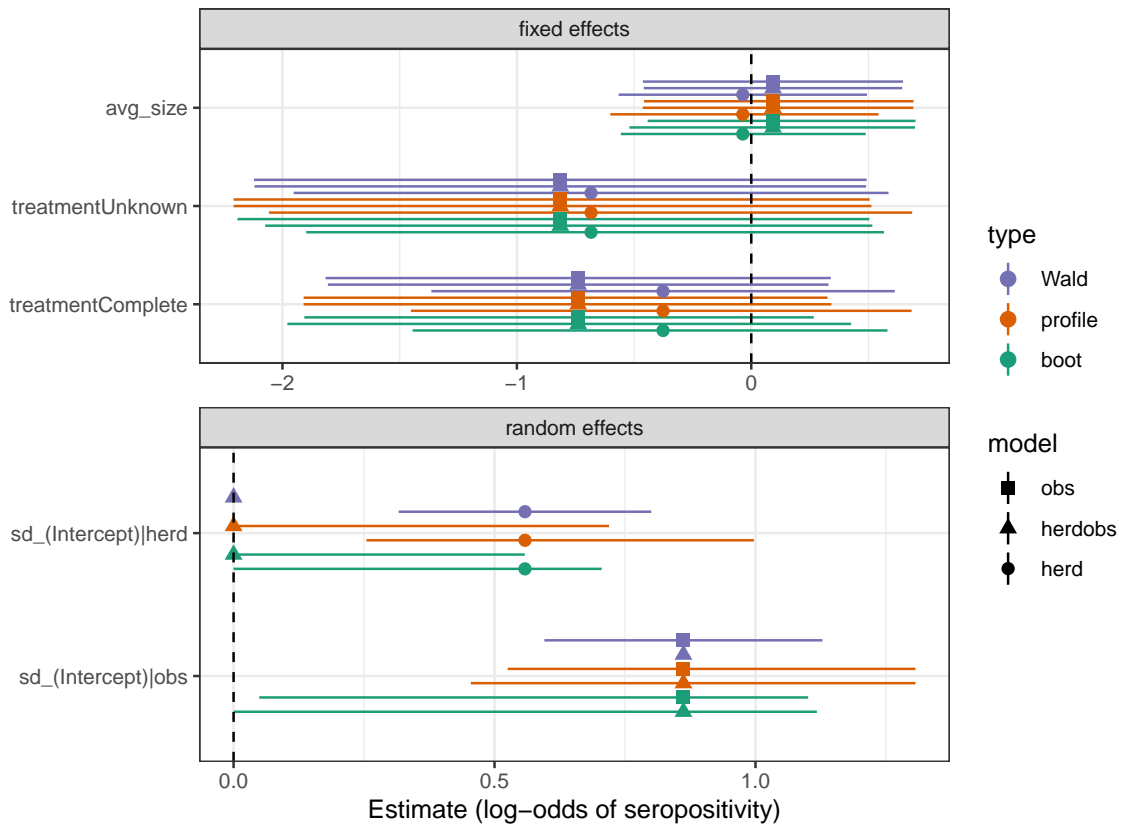


Figure 9: CBPP example: comparison of point and confidence interval estimation for different methods. Wald CIs are missing for random-effects parameters when the fitted model is singular.

```

bobyqa : [OK]
Nelder_Mead : [OK]
nlminbwrap : [OK]
nmkbw :
[OK]
optimx.L-BFGS-B : [OK]
nloptwrap.NLOPT_LN_NELDERMEAD : [OK]
nloptwrap.NLOPT_LN_BOBYQA : [OK]

```

```

> ss <- summary(gm_all)
> ss$sdcor

                herd.(Intercept)
bobyqa                0.5581743
Nelder_Mead           0.5581777
nlminbwrap            0.5581814
nmkbw                 0.5580124
optimx.L-BFGS-B      0.5581753
nloptwrap.NLOPT_LN_NELDERMEAD 0.5581391
nloptwrap.NLOPT_LN_BOBYQA    0.5581637

```

387 As of version 2.0, `lme4` can also specify structured variance-covariance matrices for (gen-
388 eralized) linear mixed models. `lme4` now supports unstructured (general positive definite),
389 diagonal, compound symmetry, and first-order autoregressive (AR1) structures. By default,
390 AR1 models assume a homogeneous-variance model (the variance is the same for all time
391 steps), while the other models assume heterogeneous-variance models (variances differ for
392 every level of the varying term); users can adjust this with the `hom` argument (e.g. `ar1(..., hom = FALSE)`).

394 The unstructured covariance structure is the default for mixed models. Here we illustrate
395 fitting an AR1 model; the next example will show diagonal and compound symmetric models.

```

> gm.ar1 <- glmer(incidence/size ~ ar1(1 + herd | period),
+               family = binomial,
+               data = cbpp, weights = size)
> print(VarCorr(gm.ar1))

```

```

Groups Name      Std.Dev. Corr
period (Intercept) 0.78959  -0.041 (ar1)

```

396 For this particular model, a heterogeneous AR1 model (`ar1(..., hom = FALSE)`) results in
397 a singular fit.

398 See the [lme4 covariance structures vignette](#) for more detail.

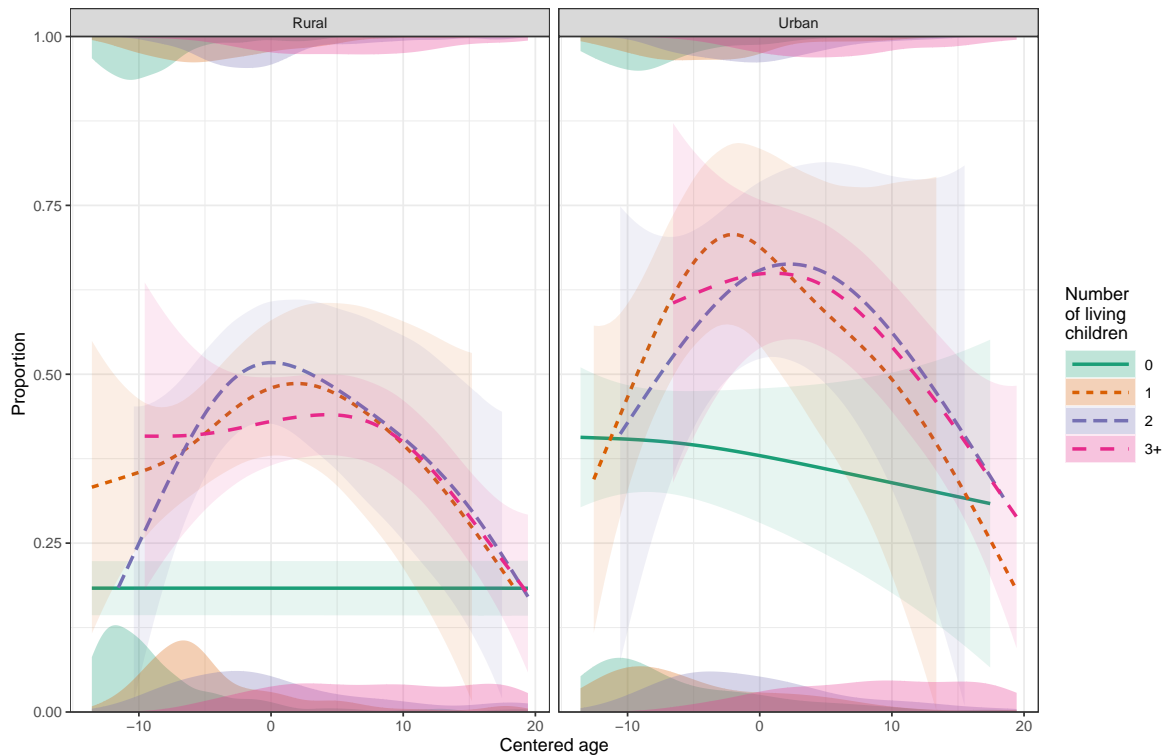


Figure 10: Contraception example: proportion of contraceptive use by centered age, number of living children (line type), and urban/rural residence (facet). Curves fitted using generalized additive models with cubic spline smoothing. Density plots along bottom and top margins show the age distributions of women not using contraception (bottom) and using contraception (top).

399 5.2. Contraception

400 [Huq and Cleland \(1990\)](#) use multilevel models to analyze data from a fertility survey of
 401 women in Bangladesh. These data are available as the `Contraception` object in the `mlmRev`
 402 package. The response variable is binary and indicates whether or not each woman was using
 403 contraception at the time of the survey. Covariates included the woman's age, the number of
 404 live children she had, whether she lived in an urban or rural setting, and the district in which
 405 she lived.

406 Figure 10 shows exploratory smooth curves of contraceptive use as a function of centered
 407 age, stratified by number of living children and urban/rural residence. In rural areas, women
 408 with two living children show the highest rates of contraceptive use, peaking near 50% at the
 409 average age, while women in urban areas with one living child show the highest rates near the
 410 average age but with a more pronounced decline at older ages. In both settings, women with
 411 no living children consistently show the lowest rates of contraceptive use. The nonmonotonic
 412 effect of age on contraceptive use and the apparent dependence of this age trend on child
 413 status motivate the model specifications explored below. (These nonmonotonic trends were
 414 not noticed in the original analysis of the data, which concluded that there was no significant
 415 effect of age.)

	df	Δ negloglik	Δ AIC
binary_child \times age + (1 district:urban)	7	0.472	0.00
binary_child \times age + (1 district/urban)	8	0.467	1.99
binary_child \times age + (1 + urban district)	9	0.000	3.06
binary_child \times age + (1 district)	7	5.825	10.71
binary_child + age + (1 district)	6	9.828	16.71
int_child + age + (1 district)	8	9.599	20.25

Table 1: Model comparison for Contraception fits. Note that for likelihood ratio tests, or for AIC comparisons restricted to nested models (Ripley 2004), the nesting sequence for the random-effects models is $(1+urban|district) > (1|district/urban) > \{(1|district), (1|district:urban)\}$.

416 We construct six `glmer` models with varying choices of explanatory variables and random
417 effects structure, comparing them via `anova`. The models considered different variables.

418 For instance, there are two different ways to encode the number of living children: `livch`,
419 a four-level factor distinguishing 0, 1, 2, or 3+ children, while other models use `ch` instead
420 (later we label as `binary_child`), which is a binary indicator for whether the woman has any
421 living children.

422 Second, we vary whether child status interacts with the woman’s age: two models include `ch`
423 (or `livch`) and `age` as additive terms, while the other four include a `ch` and `age` interaction.
424 A quadratic effect of age ($I(\text{age}^2)$) was added to account for the nonlinear effect of age.

425 Third, we explore different random effects structures at the district level. Three of them use
426 a single random intercept per district (`1 | district`). One of them extends this with a
427 random slope for urban status (`urban | district`), allowing the urban/rural difference to
428 vary by district. The next uses nested random intercepts for district and site within district
429 (`1 | district/urban`) separating district-level from urban-within-district variation. (Scan-
430 dola and Tidoni (2024) refer to this formulation as a “complex random intercepts” model).
431 The remaining model uses only the `urban:district` grouping. To reduce convergence warn-
432 ings and facilitate interpretability, age (already approximately centered in the data set) was
433 standardized by scaling by twice the standard deviation (Gelman 2008).

```

> cm1 <- glmer(use ~ age_s + I(age_s^2) + urban + livch + (1|district),
+               Contraception, binomial)
> ## switch from livch (ordinal) to ch (binary)
> cm2 <- update(cm1, . ~ . - livch + ch)
> ## add age by children interaction
> cm3 <- update(cm2, . ~ . + age_s:ch)
> ## allow urban effect to vary across districts (correlated)
> cm4 <- update(cm3, . ~ . - (1|district) + (1+urban|district))
> ## compound symmetric/nested formulation
> cm5 <- update(cm3, . ~ . - (1|district) + (1 | district/urban))
> ## as above but drop district effect
> cm6 <- update(cm3, . ~ . - (1|district) + (1 | district:urban))

```

434 Table 1 shows that the top three models, all of which include a child-by-age interaction

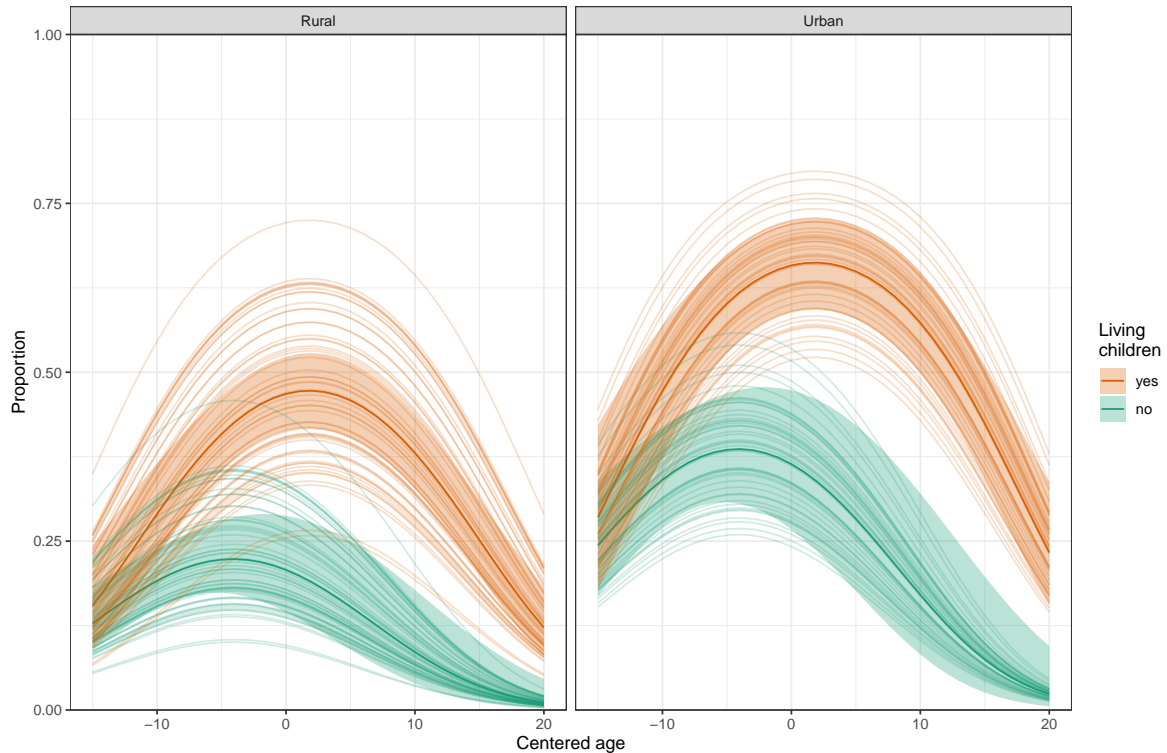


Figure 11: Predictions for contraception fit ($(1|\text{district}/\text{urban})$ model). Heavy lines and ribbons show population-level predictions; light lines show district-level predictions.

435 and some effect of urbanization, fit approximately equally well (Δ negative log-likelihood
 436 < 0.5). The $(1|\text{district}/\text{urban})$ model barely improves on the fit of $(1|\text{district}:\text{urban})$
 437 (0.005 log-likelihood units), at the cost of an extra variance parameter, so it is almost 2
 438 AIC units worse. $(1 + \text{urban} | \text{district})$ is a bit better (≈ 0.5 log-likelihood units), but
 439 includes a covariance parameter. Models that drop the interaction between child status and
 440 age or replace the binary child indicator with an integer count perform substantially worse
 441 ($\Delta\text{AIC} > 10$).

442 Overall, the results suggest that accounting for urban/rural variation at the district level and
 443 including the age and child interaction are both important, but the precise random effects
 444 structure for urban/rural variation is relatively unimportant.

445 As with the CBPP data set, we can also compute profile, Wald, and parametric bootstrap
 446 confidence intervals for all of the models to understand the effects of each variable and visualize
 447 the among-model variation.

448 The point estimates and confidence intervals of the explanatory variables are similar across
 449 the six models, and across different methods for confidence interval construction, with a few
 450 exceptions.

451 Urban residence, presence of living children, and age were all strong predictors of contraceptive
 452 use among women in Bangladesh (because we are fitting a quadratic model for the effect of
 453 age, the non-significant effect of `age_s` simply means that the marginal effect of age *at the*
 454 *mean age* is not clearly negative or positive).

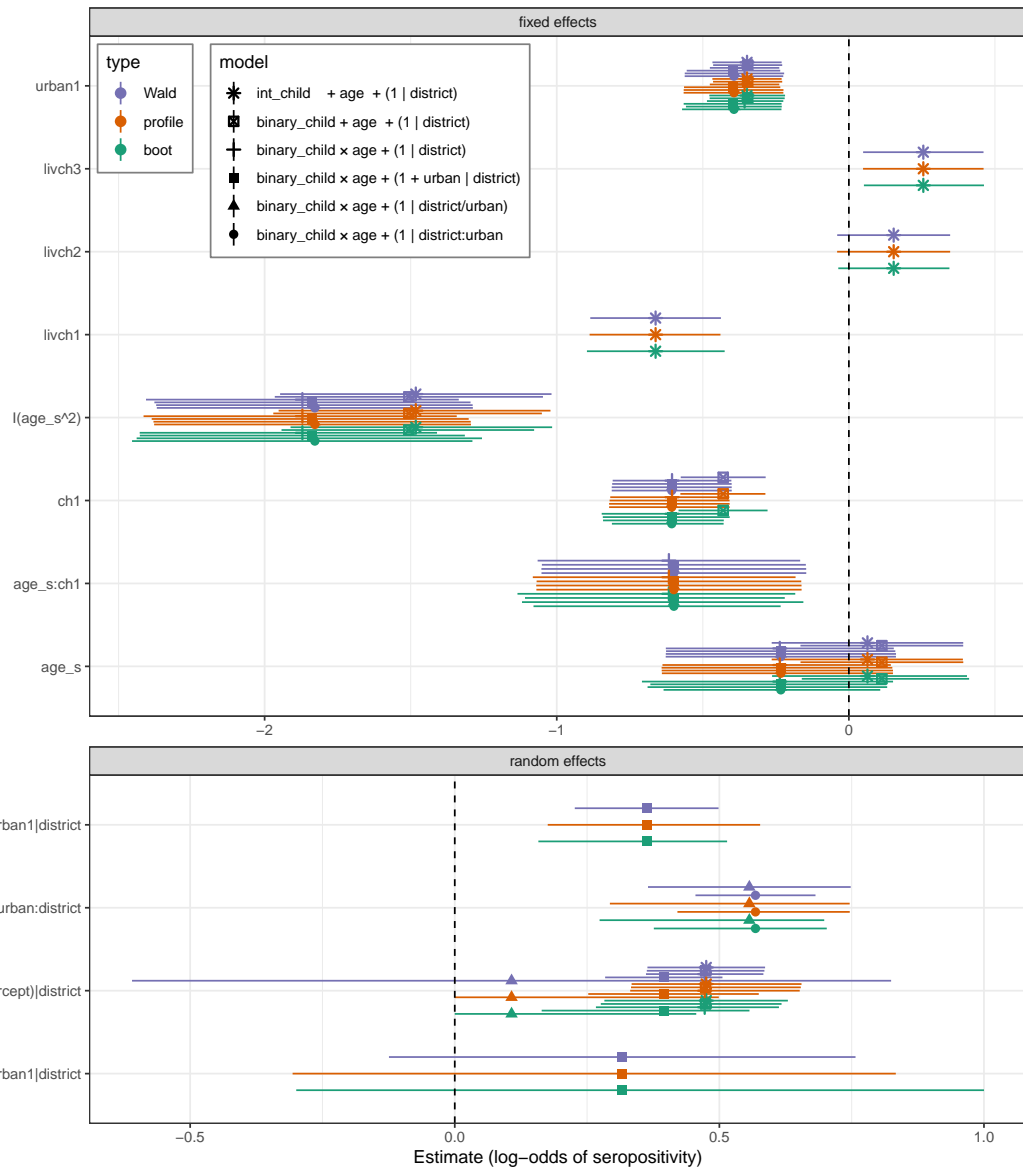


Figure 12: Contraception example: comparison of point and confidence interval estimation for different methods. (Note that the Wald CIs for variation in the intercept across districts, in the (1|district/urban) model, include negative values.)

455 With `lme4` versions greater than 2.0, we can also set up models with structured random
 456 effects covariance matrices. In the models below `diag` specifies a diagonal variance-covariance
 457 structure, while `cs` specifies a compound symmetric model. The `hom` argument specifies
 458 whether the model should allow different variances for each varying effect (e.g. for intercepts
 459 and slopes in the models below, which allow the effects of the continuous covariate of age to
 460 vary across districts).

461 Thus instead of (e.g.) `(1+urban|district)`, we can specify the random effect with diago-
 462 nal, heterogeneous variances (`diag(1+urban|district)`); diagonal, homogeneous variances
 463 (`diag(1+urban|district, hom = TRUE)`); compound symmetric, heterogeneous variances
 464 (`cs(1+urban|district)`); or compound symmetric, homogeneous variances (`cs(1+urban|district,`
 465 `hom = TRUE)`).

466 We can use `VarCorr` and other accessor methods as we would for models with default, un-
 467 structured covariance matrices, e.g.:

```
> VarCorr(cm.cs)
```

Groups	Name	Std.Dev.	Corr
district	(Intercept)	0.615	-0.79 (cs)
	urbanY	0.725	

References

- 468 Barr DJ, Levy R, Scheepers C, Tily HJ (2013). “Random Effects Structure for Confirmatory
 469 Hypothesis Testing: Keep It Maximal.” *Journal of Memory and Language*, **68**(3), 255–278.
 470 [doi:10.1016/j.jml.2012.11.001](https://doi.org/10.1016/j.jml.2012.11.001).
- 471 Bates D, Mächler M, Bolker B, Walker S (2015). “Fitting linear mixed-effects models using
 472 `lme4`.” *Journal of Statistical Software*, **67**, 1–48.
- 473 Bates DM, Watts DG (1988). *Nonlinear Regression Analysis and Its Applications*. Wiley,
 474 Hoboken, NJ. ISBN 0-471-81643-4.
- 475 De Backer M, De Vroey C, Lesaffre E, Scheys I, De Keyser P (1998). “Twelve weeks of
 476 continuous oral therapy for toenail onychomycosis caused by dermatophytes: a double-
 477 blind comparative trial of terbinafine 250 mg/day versus itraconazole 200 mg/day.” *Journal*
 478 *of the American Academy of Dermatology*, **38**(5, Supplement 2), S57–S63. [doi:10.1016/](https://doi.org/10.1016/S0190-9622(98)70486-4)
 479 [S0190-9622\(98\)70486-4](https://doi.org/10.1016/S0190-9622(98)70486-4).
- 480 Gelman A (2008). “Scaling regression inputs by dividing by two standard deviations.” *Statis-*
 481 *tics in medicine*, **27**(15), 2865–2873.
- 482 Huq N, Cleland J (1990). *Bangladesh fertility survey 1989*. National Institute of Population
 483 Research and Training, Dhaka.
- 484 Kristensen K, Nielsen A, Berg CW, Skaug H, Bell BM (2016). “TMB : Automatic Dif-
 485 ferentiation and Laplace Approximation.” *Journal of Statistical Software*, **70**(5). [doi:](https://doi.org/10.18637/jss.v070.i05)
 486 [10.18637/jss.v070.i05](https://doi.org/10.18637/jss.v070.i05).

- 487 Lesnoff M, Laval G, Bonnet P, Abdicho S, Workalemahu A, Kifle D, Peyraud A, Lancelot
 488 R, Thiaucourt F (2004). “Within-herd spread of contagious bovine pleuropneumonia in
 489 Ethiopian highlands.” *Preventive Veterinary Medicine*, **64**(1), 27–40. doi:10.1016/j.
 490 [prevetmed.2004.03.005](https://doi.org/10.1016/j.prevetmed.2004.03.005).
- 491 Madsen H, Thyregod P (2011). *Introduction to General and Generalized Linear Models*. CRC
 492 Press. ISBN 978-1-4200-9155-7.
- 493 Matuschek H, Kliegl R, Vasishth S, Baayen H, Bates D (2017). “Balancing Type I Error
 494 and Power in Linear Mixed Models.” *Journal of Memory and Language*, **94**, 305–315.
 495 doi:10.1016/j.jml.2017.01.001.
- 496 McCullagh P, Nelder JA (1989). *Generalized Linear Models*. Chapman and Hall, London.
- 497 Ripley BD (2004). “Selecting amongst Large Classes of Models.” In NM Adams, M Crowder,
 498 D Hand, D Stephens (eds.), *Methods and models in statistics: In honor of Professor John*
 499 *Nelder, FRS*, pp. 155–170. Imperial College Press.
- 500 Scandola M, Tidoni E (2024). “Reliability and Feasibility of Linear Mixed Models in Fully
 501 Crossed Experimental Designs.” *Advances in Methods and Practices in Psychological Sci-*
 502 *ence*, **7**(1), 25152459231214454. doi:10.1177/25152459231214454.
- 503 Stringer A (2024). “Exact Gradient Evaluation for Adaptive Quadrature Approximate
 504 Marginal Likelihood in Mixed Models for Grouped Data.” *Statistics and Computing*, **35**(1),
 505 4. ISSN 1573-1375. doi:10.1007/s11222-024-10536-z.
- 506 Stringer A, Bilodeau B, Tang Y (2022). “Asymptotics of Numerical Integration for Two-Level
 507 Mixed Models.” doi:10.48550/arXiv.2202.07864.
- 508 Zeger SL, Karim MR (1991). “Generalized linear models with random effects: a Gibbs
 509 sampling approach.” *Journal of the American Statistical Association*, **86**(413), 79–86.

6. Package versions used

510 Compiled with R version 4.6.0 (2026-04-24) and package versions lme4: 2.0.2, performance:
 511 0.17.0, DHARMA: 0.5.0, see: 0.14.0.

7. Appendix: derivation of PIRLS

512 We seek to maximize the unscaled conditional log density for a GLMM over the conditional
 513 modes, \mathbf{u} . This problem is very similar to maximizing the log-likelihood for a GLM, which is a
 514 very thoroughly studied problem (e.g. McCullagh and Nelder 1989). The standard algorithm
 515 for dealing with this kind of problem is iteratively reweighted least squares (IRLS). Here
 516 we modify IRLS by incorporating a penalty term that accounts for variation in the random
 517 effects; we call the resulting algorithm penalized iteratively reweighted least squares (PIRLS).
 518 The unscaled conditional log-density takes the form,

$$f(\mathbf{u}) = \log p(\mathbf{y}, \mathbf{u} | \boldsymbol{\beta}, \boldsymbol{\theta}) = \boldsymbol{\psi}^\top \mathbf{A} \mathbf{y} - \mathbf{a}^\top \boldsymbol{\phi} + \mathbf{c} - \frac{1}{2} \mathbf{u}^\top \mathbf{u} - \frac{q}{2} \log 2\pi \quad (28)$$

519 where $\boldsymbol{\psi}$ is the n -by-1 canonical parameter of an exponential family, $\boldsymbol{\phi}$ is the n -by-1 vector of
 520 cumulant functions, \boldsymbol{c} an n -by-1 vector of normalizing constants, and \boldsymbol{A} is an n -by- n diagonal
 521 matrix of prior weights, \boldsymbol{a} . Both \boldsymbol{a} and \boldsymbol{c} could depend on a dispersion parameter, although
 522 we ignore this possibility for now.

523 The canonical parameter, $\boldsymbol{\psi}$, and vector of cumulant functions, $\boldsymbol{\phi}$, depend on a linear predic-
 524 tor,

$$\boldsymbol{\eta} = \boldsymbol{o} + \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{Z}\boldsymbol{\Lambda}_\theta\boldsymbol{u} \quad (29)$$

525 where \boldsymbol{o} is an n -by-1 vector of *a priori* offsets. The specific form of this dependency is specified
 526 by the choice of the exponential family. The mean of this distribution, $\boldsymbol{\mu}$, is the *inverse link*
 527 *function* g^{-1} applied to $\boldsymbol{\eta}$.

528 Our goal is to find the values of \boldsymbol{u} that maximize the unscaled conditional density, for given $\boldsymbol{\theta}$
 529 and $\boldsymbol{\beta}$ vectors. These maximizers are the conditional modes, which we require for the Laplace
 530 approximation and adaptive Gauss-Hermite quadrature. To do this maximization we use a
 531 variant of the Fisher scoring method, which is the basis of the iteratively reweighted least
 532 squares algorithm for generalized linear models. Fisher scoring is itself based on Newton's
 533 method, which we apply first.

534 7.1. Newton's method

535 To apply Newton's method, we need the gradient and the Hessian of the unscaled conditional
 536 log-likelihood. Following standard GLM theory (McCullagh and Nelder 1989), we use the
 537 chain rule,

$$\frac{dL(\boldsymbol{\beta}, \boldsymbol{\theta}|\boldsymbol{y}, \boldsymbol{u})}{d\boldsymbol{u}} = \frac{dL(\boldsymbol{\beta}, \boldsymbol{\theta}|\boldsymbol{y}, \boldsymbol{u})}{d\boldsymbol{\psi}} \frac{d\boldsymbol{\psi}}{d\boldsymbol{\mu}} \frac{d\boldsymbol{\mu}}{d\boldsymbol{\eta}} \frac{d\boldsymbol{\eta}}{d\boldsymbol{u}}$$

538 The first derivative in this chain follow from basic results in GLM theory,

$$\frac{dL(\boldsymbol{\beta}, \boldsymbol{\theta}|\boldsymbol{y}, \boldsymbol{u})}{d\boldsymbol{\psi}} = (\boldsymbol{y} - \boldsymbol{\mu})^\top \boldsymbol{A}$$

539 Again from standard GLM theory, the next two derivatives define the inverse diagonal variance
 540 matrix,

$$\frac{d\boldsymbol{\psi}}{d\boldsymbol{\mu}} = \boldsymbol{V}^{-1}$$

541 and the diagonal Jacobian matrix,

$$\frac{d\boldsymbol{\mu}}{d\boldsymbol{\eta}} = \boldsymbol{M} \quad .$$

542 Finally, because $\boldsymbol{\beta}$ affects $\boldsymbol{\eta}$ only linearly,

$$\frac{d\boldsymbol{\eta}}{d\boldsymbol{u}} = \boldsymbol{Z}\boldsymbol{\Lambda}_\theta$$

543 Therefore we have,

$$\frac{dL(\boldsymbol{\beta}, \boldsymbol{\theta}|\boldsymbol{y}, \boldsymbol{u})}{d\boldsymbol{u}} = (\boldsymbol{y} - \boldsymbol{\mu})^\top \boldsymbol{A}\boldsymbol{V}^{-1}\boldsymbol{M}\boldsymbol{Z}\boldsymbol{\Lambda}_\theta + \boldsymbol{u}^\top \quad . \quad (30)$$

544 This is very similar to the gradient for GLMs with respect to fixed effects coefficients, $\boldsymbol{\beta}$.
 545 The only difference induced by differentiating with respect to the random effects, \boldsymbol{u} , is the
 546 addition of the \boldsymbol{u}^\top term.

547 Again we apply the chain rule to take the Hessian,

$$\frac{d^2 L(\boldsymbol{\beta}, \boldsymbol{\theta} | \mathbf{y}, \mathbf{u})}{d\mathbf{u} d\mathbf{u}} = \frac{d^2 L(\boldsymbol{\beta}, \boldsymbol{\theta} | \mathbf{y}, \mathbf{u})}{d\mathbf{u} d\boldsymbol{\mu}} \frac{d\boldsymbol{\mu} d\boldsymbol{\eta}}{d\boldsymbol{\eta} d\mathbf{u}} + \mathbf{I}_q \quad (31)$$

548 which leads to,

$$\frac{d^2 L(\boldsymbol{\beta}, \boldsymbol{\theta} | \mathbf{y}, \mathbf{u})}{d\mathbf{u} d\mathbf{u}} = \frac{d^2 L(\boldsymbol{\beta}, \boldsymbol{\theta} | \mathbf{y}, \mathbf{u})}{d\mathbf{u} d\boldsymbol{\mu}} \mathbf{M} \mathbf{Z} \boldsymbol{\Lambda}_\theta + \mathbf{I}_q \quad (32)$$

549 The first derivative in this chain can be expressed as,

$$\frac{d^2 L(\boldsymbol{\beta}, \boldsymbol{\theta} | \mathbf{y}, \mathbf{u})}{d\mathbf{u} d\boldsymbol{\mu}} = -\boldsymbol{\Lambda}_\theta^\top \mathbf{Z}^\top \mathbf{M} \mathbf{V}^{-1} \mathbf{A} + \boldsymbol{\Lambda}_\theta^\top \mathbf{Z}^\top \left[\frac{d\mathbf{M} \mathbf{V}^{-1}}{d\boldsymbol{\mu}} \right] \mathbf{A} \mathbf{R} \quad (33)$$

550 where \mathbf{R} is a diagonal residuals matrix with $\mathbf{y} - \boldsymbol{\mu}$ on the diagonal. The two terms arise from
551 a type of product rule, where we first differentiate the residuals, $\mathbf{y} - \boldsymbol{\mu}$, and then the diagonal
552 matrix, $\mathbf{M} \mathbf{V}^{-1}$, with respect to $\boldsymbol{\mu}$.

553 The Hessian can therefore be expressed as,

$$\frac{d^2 L(\boldsymbol{\beta}, \boldsymbol{\theta} | \mathbf{y}, \mathbf{u})}{d\mathbf{u} d\mathbf{u}} = -\boldsymbol{\Lambda}_\theta^\top \mathbf{Z}^\top \mathbf{M} \mathbf{A}^{1/2} \mathbf{V}^{-1/2} \left(\mathbf{I}_n - \mathbf{V} \mathbf{M}^{-1} \left[\frac{d\mathbf{M} \mathbf{V}^{-1}}{d\boldsymbol{\mu}} \right] \mathbf{R} \right) \mathbf{V}^{-1/2} \mathbf{A}^{1/2} \mathbf{M} \mathbf{Z} \boldsymbol{\Lambda}_\theta + \mathbf{I}_q \quad (34)$$

554 This result can be simplified by expressing it in terms of a weighted random-effects design
555 matrix, $\mathbf{U} = \mathbf{A}^{1/2} \mathbf{V}^{-1/2} \mathbf{M} \mathbf{Z} \boldsymbol{\Lambda}_\theta$,

$$\frac{d^2 L(\boldsymbol{\beta}, \boldsymbol{\theta} | \mathbf{y}, \mathbf{u})}{d\mathbf{u} d\mathbf{u}} = -\mathbf{U}^\top \left(\mathbf{I}_n - \mathbf{V} \mathbf{M}^{-1} \left[\frac{d\mathbf{V}^{-1} \mathbf{M}}{d\boldsymbol{\mu}} \right] \mathbf{R} \right) \mathbf{U} + \mathbf{I}_q \quad (35)$$

556 7.2. Fisher-like scoring

557 There are two ways to further simplify this expression for $\mathbf{U}^\top \mathbf{U}$. The first is to use the
558 canonical link function for the family being used. Canonical links have the property that
559 $\mathbf{V} = \mathbf{M}$, which means that for canonical links,

$$\frac{d^2 L(\boldsymbol{\beta}, \boldsymbol{\theta} | \mathbf{y}, \mathbf{u})}{d\mathbf{u} d\mathbf{u}} = -\mathbf{U}^\top \left(\mathbf{I}_n - \mathbf{I}_n \left[\frac{d\mathbf{I}_n}{d\boldsymbol{\mu}} \right] \mathbf{R} \right) \mathbf{U} + \mathbf{I}_q = \mathbf{U}^\top \mathbf{U} + \mathbf{I}_q \quad (36)$$

560 The second way to simplify the Hessian is to take its expectation with respect to the dis-
561 tribution of the response, conditional on the current values of the spherical random effects
562 coefficients, \mathbf{u} . The diagonal residual matrix, \mathbf{R} , has expectation 0. Therefore, because the
563 response only enters into the expression for the Hessian via \mathbf{R} , we have that,

$$E \left(\frac{d^2 L(\boldsymbol{\beta}, \boldsymbol{\theta} | \mathbf{y}, \mathbf{u})}{d\mathbf{u} d\mathbf{u}} \middle| \mathbf{u} \right) = -\mathbf{U}^\top \left(\mathbf{I}_n - \mathbf{U} \mathbf{M}^{-1} \left[\frac{d\mathbf{V}^{-1} \mathbf{M}}{d\boldsymbol{\mu}} \right] E(\mathbf{R}) \right) \mathbf{U} + \mathbf{I}_q = \mathbf{U}^\top \mathbf{U} + \mathbf{I}_q \quad (37)$$

564 **Affiliation:**

565 Anna Ly
566 Department of Mathematics & Statistics
567 McMaster University
568 1280 Main Street W
569 Hamilton, ON L8S 4K1, Canada
570

571 Rune Haubo Bojesen Christensen
572 Copenhagen Research Centre for Biological and Precision Psychiatry
573 Mental Health Centre Copenhagen, Copenhagen University Hospital - Bispebjerg and Fred-
574 eriksberg
575 Copenhagen, Denmark
576 Email: Rune.Haubo@pm.me

577 Douglas Bates
578 Department of Statistics, University of Wisconsin - Madison
579 1205 University Ave.
580 Madison, WI 53706, U.S.A.
581 E-mail: bates@stat.wisc.edu

582 Martin Mächler
583 Seminar für Statistik, HG G 16
584 ETH Zurich
585 8092 Zurich, Switzerland
586 E-mail: maechler@stat.math.ethz.ch
587

588 Benjamin M. Bolker
589 Departments of Mathematics & Statistics and Biology
590 McMaster University
591 1280 Main Street W
592 Hamilton, ON L8S 4K1, Canada
593 E-mail: bolker@mcmaster.ca