

Generalized linear mixed models - computation of the likelihood function

Rasmus Waagepetersen
Department of Mathematics
Aalborg University
Denmark

November 6, 2025

Generalized linear mixed effects models

Consider stochastic vector $Y = (Y_1, \dots, Y_n)$ and vector of random effects $U = (U_1, \dots, U_m)$.

Two step formulation of GLMM:

- ▶ $U \sim N(0, \Sigma)$.
- ▶ Given realization u of U , Y_i independent and each follows density $f_i(y_i|u)$ with mean $\mu_i = g^{-1}(\eta_i)$ and linear predictor $\eta = X\beta + Zu$.

I.e. conditional on U , Y_i follows a generalized linear model.

NB: GLMM specified in terms of marginal density of U and conditional density of Y given U . But the likelihood is the marginal density $f(y)$ which can be hard to evaluate !

We already saw one example: logistic regression with random effects.

Another common example: Poisson-log normal. Here

$$U \sim N(0, \Sigma)$$

$$Y_i | U = u \sim \text{Pois}(\exp(\eta_i))$$

where $\eta_i = x_i\beta + z_i u$

Likelihood for generalized linear mixed model

Likelihood for a generalized linear mixed model given by integral:

$$f(y) = \int_{\mathbb{R}^m} f(y, u) du = \int_{\mathbb{R}^m} f(y|u)f(u) du$$

Difficult since $f(y|u)f(u)$ is a very complex function.

Huge statistical literature on how to compute good approximations of the likelihood: Laplace approximation, numerical quadrature, Monte Carlo, Markov chain Monte Carlo,...

EM-algorithm: method for maximizing $f(y; \theta)$ with respect to θ without computing $f(y; \theta)$

Note: equivalent to computing “observed data” likelihood in case of missing data.

Example: logistic regression with random intercepts

$$U_j \sim N(0, \tau^2), j = 1, \dots, m$$

$$Y_j | U_j = u_j \sim \text{binomial}(n_j, p_j)$$

$$\log(p_j / (1 - p_j)) = \eta_j = \beta + U_j$$

$$p_j = \exp(\eta_j) / (1 + \exp(\eta_j))$$

Conditional density:

$$f(y|u; \beta) = \prod_j p_j^{y_j} (1 - p_j)^{n_j - y_j} = \prod_j \frac{\exp(\beta + u_j)^{y_j}}{(1 + \exp(\beta + u_j))^{n_j}}$$

Likelihood function ($u = (u_1, \dots, u_m)$)

$$\int_{\mathbb{R}^m} f(y|u; \beta) f(u; \tau^2) du = \prod_j \int_{\mathbb{R}} \frac{\exp(\beta + u_j)^{y_j}}{(1 + \exp(\beta + u_j))^{n_j}} \frac{\exp(-u_j^2 / (2\tau^2))}{\sqrt{2\pi\tau^2}} du_j$$

Integrals can not be evaluated in closed form.

Hierarchical model with independent random effects

Suppose $U = (U_1, \dots, U_m)$ with the the U_i independent.

Moreover $Y = (Y_{ij})_{ij}$, $i = 1, \dots, m$ and $j = 1, \dots, n_i$ where the conditional distribution of the $Y_i = (Y_{ij})_j$ only depends on U_i .

Then we can factorize likelihood as

$$f(y) = \prod_{i=1}^m \int_{\mathbb{R}} f(y_i|u_i)f(u_i)du_i$$

That is, product of one-dimensional integrals.

Consider in the following computation of one-dimensional integral.

One-dimensional case

Compute

$$L(\theta) = \int_{\mathbb{R}} f(y|u; \beta) f(u; \tau^2) du$$

Some possibilities:

- ▶ Laplace approximation.
- ▶ Numerical integration/quadrature (e.g. Gaussian quadrature as in `glmer`, PROC NL MIXED (SAS) or GLLAM (STATA)) (one level of random effects, dimensions one or two).

Laplace approximation

Let $g(u) = \log(f(y|u)f(u))$ and choose \hat{u} so $g'(\hat{u}) = 0$
($\hat{u} = \arg \max g(u)$).

Taylor expansion around \hat{u} :

$$g(u) \approx \tilde{g}(u) =$$

$$g(\hat{u}) + (u - \hat{u})g'(\hat{u}) + \frac{1}{2}(u - \hat{u})^2 g''(\hat{u}) = g(\hat{u}) - \frac{1}{2}(u - \hat{u})^2 (-g''(\hat{u}))$$

i.e. $\exp(\tilde{g}(u))$ proportional to normal density $N(\mu_{LP}, \sigma_{LP}^2)$,
 $\mu_{LP} = \hat{u}$ $\sigma_{LP}^2 = -1/g''(\hat{u})$.

$$\begin{aligned} L(\theta) &= \int_{\mathbb{R}} \exp(g(u)) du \approx \int_{\mathbb{R}} \exp(\tilde{g}(u)) du \\ &= \exp(g(\hat{u})) \int_{\mathbb{R}} \exp\left(-\frac{1}{2\sigma_{LP}^2}(u - \mu_{LP})^2\right) du = \exp(g(\hat{u})) \sqrt{2\pi\sigma_{LP}^2} \end{aligned}$$

Laplace approximation also works for for higher dimensions (multivariate Taylor expansion).

NB:

$$f(u|y) = f(y|u)f(u)/f(y) \propto \exp(g(u)) \approx \text{const} \exp\left(-\frac{1}{2\sigma_{LP}^2}(u-\mu_{LP})^2\right)$$

where $\mu_{LP} = \hat{u}$ $\sigma_{LP}^2 = -1/g''(\hat{u})$.

Hence

$$U|Y = y \approx N(\mu_{LP}, \sigma_{LP}^2)$$

Note: μ_{LP} is mode of conditional distribution - used for prediction of random effects in `glmer(ranef())`.

Gaussian quadrature

Gauss-Hermite quadrature (numerical integration) is

$$\int_{\mathbb{R}} f(x)\phi(x)dx \approx \sum_{i=1}^n w_i f(x_i)$$

where ϕ is the standard normal density and $(x_i, w_i), i = 1, \dots, n$ are certain arguments and weights which can be looked up in a table.

We can replace \approx with $=$ whenever f is a polynomial of degree $2n - 1$ or less.

In other words $(x_i, w_i), i = 1, \dots, n$ is the solution of the system of $2n$ equations

$$\int_{\mathbb{R}} x^k \phi(x) dx = \sum_{i=1}^n w_i x_i^k, \quad k = 0, \dots, 2n - 1$$

where

$$\int_{\mathbb{R}} x^k \phi(x) dx = 1[k \text{ even}] (k-1)!! = 1[k \text{ even}] (k-1)(k-2)\dots(2)(1)$$

Adaptive Gauss-Hermite quadrature

Naive application of Gauss-Hermite ($U \sim N(0, \tau^2)$):

$$\int f(y|u)f(u)du = \int f(y|\tau x)\phi(x)dx$$

Now GH is applicable.

Adaptive GH:

$$\begin{aligned}\int f(y|u)f(u)du &= \int \frac{f(y|u)f(u)}{\phi(u; \mu_{LP}, \sigma_{LP}^2)}\phi(u; \mu_{LP}, \sigma_{LP}^2)du = \\ &= \int \frac{f(y|\sigma_{LP}x + \mu_{LP})f(\sigma_{LP}x + \mu_{LP})}{\phi(x)}\sigma_{LP}\phi(x)dx\end{aligned}$$

(change of variable: $x = (u - \mu_{LP})/\sigma_{LP}$)

In my experience, adaptive GH is way more accurate than naive GH.

Advantage

$$\frac{f(y|u)f(u)}{\phi(u; \mu_{LP}, \sigma_{LP}^2)} = \frac{f(y|\sigma_{LP}x + \mu_{LP})f(\sigma_{LP}x + \mu_{LP})}{\phi(x)} \quad x = (u - \mu_{LP})/\sigma_{LP}$$

close to constant ($f(y)$) – hence adaptive G-H quadrature very accurate.

GH scheme with $n = 5$:

x	2.020	0.959	0.0000000	-0.959	-2.020
w	0.011	0.222	0.533	0.222	0.011

(x 's are roots of Hermite polynomial computed using `ghq` in library `glmML`).

(GH schemes for $n = 5$ and $n = 10$ available on web page)

Prediction of random effects for GLMM

Conditional mean

$$\mathbb{E}[U|Y = y] = \int uf(u|y)du$$

is minimum mean square error predictor, i.e.

$$\mathbb{E}(U - \hat{U})^2$$

is minimal with $\hat{U} = H(Y)$ where $H(y) = \mathbb{E}[U|Y = y]$

Difficult to analytically evaluate

$$\mathbb{E}[U|Y = y] = \int uf(y|u)f(u)/f(y)du$$

Computation of conditional expectations (prediction)

$$\begin{aligned}\mathbb{E}[U|Y = y] &= \int u \frac{f(y|u)f(u)}{f(y)} du = \\ \frac{1}{f(y)} \int (\sigma_{LP}x + \mu_{LP}) \frac{f(y|\sigma_{LP}x + \mu_{LP})f(\sigma_{LP}x + \mu_{LP})}{\phi(x)} \sigma_{LP} \phi(x) dx\end{aligned}$$

Note:

$$(\sigma_{LP}x + \mu_{LP}) \frac{f(y|\sigma_{LP}x + \mu_{LP})f(\sigma_{LP}x + \mu_{LP})}{\phi(x)} \sigma_{LP}$$

behaves like a first order polynomial in x - hence GH still accurate.

Score function and Fisher information

Let

$$V_{\theta}(y, u) = \frac{d}{d\theta} \log f(y, u|\theta)$$

Then score and observed information are

$$u(\theta) = \frac{d}{d\theta} \log L(\theta) = \mathbb{E}_{\theta}[V_{\theta}(y, U)|Y = y] \quad (1)$$

and

$$\begin{aligned} j(\theta) &= -\frac{d^2}{d\theta^T d\theta} \log L(\theta) \\ &= -(\mathbb{E}_{\theta}[dV_{\theta}(y, U)/d\theta^T | Y = y] + \text{Var}_{\theta}[V_{\theta}(y, U)|Y = y]) \end{aligned} \quad (2)$$

Again the above expectations and variances can be evaluated using Laplace or adaptive GH.

Newton-Raphson iterations:

$$\theta_{l+1} = \theta_l + j(\theta_l)^{-1} u(\theta_l)$$

Comparison with EM-algorithm

Given current estimate θ_l :

1. (E) compute $Q(\theta_l, \theta) = \mathbb{E}_{\theta_l}[\log f(y, U|\theta) | Y = y]$
2. (M) $\theta_{l+1} = \operatorname{argmax}_{\theta} Q(\theta_l, \theta)$.

For GLMMs (E) step needs numerical integration or Monte Carlo.

Convergence of EM-algorithm can be quite slow. Maximization of likelihood using Newton-Raphson (or its variants) seems to me better alternative.

Difficult cases for numerical integration - dimension $m > 1$

- ▶ correlated random effects: multivariate density of U does not factorize
- ▶ crossed random effects: U_i and V_j independent $i = 1, \dots, m$
 $j = 1, \dots, n$ but Y_{ij} depends on both U_i and V_j .

Not possible to factorize likelihood into low-dimensional integrals

Number of quadrature points $\approx k^m$ where k is number of quadrature points for 1D and m number of random effects – hence numerical quadrature may not be feasible.

Alternatives: Laplace-approximation or *Monte Carlo computation*.

Wheeze results with different values of nAGQ

Default nAGQ=1 means Laplace approximation:

```
> fiter=glmer(resp~age+smoke+(1|id),family=binomial,data=ob  
> summary(fiter)
```

Random effects:

Groups Name	Variance	Std.Dev.
id (Intercept)	5.491	2.343

Number of obs: 2148, groups: id, 537

Fixed effects:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-3.37396	0.27496	-12.271	<2e-16 ***
age	-0.17677	0.06797	-2.601	0.0093 **
smoke	0.41478	0.28705	1.445	0.1485

5 quadrature points:

```
> fiter5=glmer(resp~age+smoke+(1|id),family=binomial,  
               data=ohio,nAGQ=5)
```

```
Groups Name          Variance Std.Dev.  
 id      (Intercept) 4.198    2.049
```

```
Number of obs: 2148, groups: id, 537
```

Fixed effects:

	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	-3.02398	0.20353	-14.857	< 2e-16	***
age	-0.17319	0.06718	-2.578	0.00994	**
smoke	0.39448	0.26305	1.500	0.13371	

10 quadrature points:

```
> fiter10=glmer(resp~age+smoke+(1|id),family=binomial
, data=ohio,nAGQ=10)
```

Random effects:

Groups Name	Variance	Std.Dev.
id (Intercept)	4.614	2.148

Number of obs: 2148, groups: id, 537

Fixed effects:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-3.08959	0.21557	-14.332	< 2e-16 ***
age	-0.17533	0.06762	-2.593	0.00952 **
smoke	0.39799	0.27167	1.465	0.14293

Some sensivity regarding variance estimate. Fixed effects results quite stable.

Results with 20 quadrature points very similar to those with 10 quadrature points.

Laplace - mathematical details in one-dimension

(one dimension to avoid technicalities of multivariate Taylor)

Let

$$I_n = \int_{\mathbb{R}} \exp(nh(x))g(x)dx$$

where $h(x)$ is three times differentiable and assume there exists \hat{x} so that

1. $H = -h''(\hat{x}) > 0$ and $h'(\hat{x}) = 0$
2. for any $\Delta > 0$ there exists an $\epsilon > 0$ so that $h(\hat{x}) - h(x) > \epsilon$ for $|x - \hat{x}| > \Delta$
3. there exists a $\delta > 0$ so that $|h^{(3)}(x)| < K$ and $|g(x)| < C$ for $|x - \hat{x}| \leq \delta$
4. a) $\int_{\mathbb{R}} |g(x)|dx < \infty$ or b) $\int_{\mathbb{R}} \exp(h(x))|g(x)|dx < \infty$

Then

$$\frac{I_n}{\exp(nh(\hat{x}))g(\hat{x})\sqrt{2\pi n^{-1}H^{-1}}} \rightarrow 1 \quad (3)$$

as $n \rightarrow \infty$.

Relative error of approximation

Absolute error of approximation is

$$E_n = I_n - \exp(nh(\hat{x}))g(\hat{x})\sqrt{2\pi n^{-1}H^{-1}}$$

Previous result says that relative error

$$\frac{E_n}{I_n} \rightarrow 0$$

Strong result in case I_n is a small quantity (may not be enough that absolute error is “small”)

We can say more:

$$\frac{I_n}{\exp(nh(\hat{x}))g(\hat{x})\sqrt{2\pi n^{-1}H^{-1}}} = 1 + O(n^{-1}).$$

That is, the relative error is of order n^{-1} .

Exercises

1. How does Laplace approximation look in the multivariate case ?
2. Show that adaptive GH with one quadrature point is equivalent to Laplace approximation.
3. Show the identities (1) and (2) (assuming differentiation and integration can be interchanged as needed).
4. Write down the likelihood in case of crossed random effects. What is the problem ?
5. Solve exercises on `exercises_lp_gh.pdf`
6. Carefully check the proof for Laplace approximation in Section 1 in note available on course webpage. If you like Taylor expansions you may also want to check Sections 2-3.
7. Consider the case of one Normal random effect $U \sim N(0, \tau^2)$ and observations Y_1, \dots, Y_n that are iid given U . Can you apply the formal result for the Laplace approximation to show convergence of the approximation of the likelihood ? Which problems do you face ?